

NEGATIVE-ENERGY MATTER AND THE DIRECTION OF TIME

J.C. Lindner*

Department of Physics, Université de Montréal
Montréal, QC, Canada

17 June 2016 (updated 11 February 2024)[†]

Abstract

This report offers a modern perspective on the problem of negative energy, based on a reexamination of the concept of time direction as it arises in a classical and quantum-mechanical context. From this analysis emerges an improved understanding of the general-relativistic stress-energy of matter as being a manifestation of local variations in the energy density of zero-point vacuum fluctuations. Based on those developments, a set of axioms is proposed from which are derived generalized gravitational field equations which actually constitute a simplification of relativity theory in the presence of negative-energy matter and a non-zero cosmological constant. Important clarifications are also achieved regarding the nature of the binary degrees of freedom of matter in the final stages of a gravitational collapse. Those results are then applied to provide original solutions to several long-standing problems in cosmology, including the problem of the nature of dark matter and dark energy, that of the origin of thermodynamic time asymmetry, and several other issues traditionally approached using inflation theory. Finally, we draw on those developments to provide significant new insights into the foundations of quantum theory, regarding, in particular, the problem of quantum non-locality, that of the emergence of time in quantum cosmology, as well as the question of the persistence of quasiclassicality following decoherence.

*Email address: research@jclindner.ca

[†]Preprint coordinates: arxiv.org/abs/1608.01531

Contents

1	Introduction	8
1.1	Motivations	10
1.2	Approach	12
1.3	Historical context	17
1.4	Organizing principles	20
2	Negative Energy and Gravitation	24
2.1	The negative-energy problem	24
2.2	The time-direction degree of freedom	31
2.3	Our current understanding	44
2.4	The negative-mass concept	52
2.5	The equivalence principle with negative mass	74
2.6	An effect of voids in the matter distribution	91
2.7	Six problems for negative-energy matter	111
2.8	The origin of repulsive gravitational forces	116
2.9	No energy out of nothing	123
2.10	The problem of vacuum decay	128
2.11	Energy and momentum conservation	134
2.12	Absolute inertial mass	139
2.13	A few other misconceptions	142
2.14	An axiomatic formulation	151
2.15	Generalized gravitational field equations	158
3	Time Reversal and Information	182
3.1	The problem of discrete symmetries	182
3.2	The constraint of relational definition	184
3.3	The concept of bidirectional time	189
3.4	Alternative definition of C , P , and T	197

3.5	The time-reversal operation	203
3.6	The charge-conjugation operation	210
3.7	Invariance under combined reversals	214
3.8	The significance of classical equations	218
3.9	Reversal of action	220
3.10	Black-hole entropy	229
3.11	Negative temperatures	269
4	Cosmology and Irreversibility	273
4.1	The outstanding problems of cosmology	273
4.2	The cosmological-constant problem	276
4.3	Missing mass and dark matter	300
4.4	Large-scale structure	341
4.5	The flatness problem and matter creation	350
4.6	The problem of time asymmetry	374
4.7	Gravitational entropy	385
4.8	The initial singularity	408
4.9	The horizon problem and irreversibility	422
4.10	A criticism of inflation theory	446
5	Quantum Theory and Causality	454
5.1	The problem of interpretation	454
5.2	A simple analogy	457
5.3	Time-symmetric causality	460
5.4	Closed causal chains and time travel	468
5.5	Advanced waves and time asymmetry	478
5.6	Early interpretations	483
5.7	The constraint of scientific realism	491
5.8	Time-symmetric quantum theory	518
5.9	Quantum entanglement and non-locality	538
5.10	The quantum measurement problem	549
5.11	The emergence of time in quantum cosmology	570
5.12	Universal causal chain and quasiclassicality	583
5.13	A possible role for gravitation	615
6	Conclusion	633
6.1	Historical perspective	642

Acknowledgements	646
Bibliography	647
Index	653

List of Figures

3.1	Variation of physical parameters under the proposed alternative definition of P , T , and C , as described from the bidirectional-time viewpoint	193
3.2	Variation of physical parameters under the proposed alternative definition of P , T , and C , as apparent from the unidirectional-time viewpoint	196
3.3	Four different outcomes of applying each of the relationally distinct action-reversal symmetry operations, as described from the bidirectional-time viewpoint	223
4.1	Alternative Feynman diagrams for flavor- and color-changing interactions between quarks	312
4.2	Alternative Feynman diagrams for flavor-conserving electroweak interactions between fermions	313
5.1	The four possible, combined, retarded and advanced histories of a double slit or simple interferometer experiment	524

List of Tables

3.1	Variation of the physical parameters associated with a process transformed by the discrete P , T , and C symmetry operations, as they are traditionally defined	199
3.2	Variation of physical parameters under the redefined discrete P , T , and C symmetry operations, as described from the bidirectional-time viewpoint	200
3.3	Variation of physical parameters under the redefined discrete P , T , and C symmetry operations, as described from the unidirectional-time viewpoint	200
3.4	Variations of physical parameters under the four relationally distinct action-sign-reversing symmetry operations, as described from the bidirectional-time viewpoint	222

Note to the reader

For reading convenience, an adapted version of this report will soon be published in print and shall be available for order at selected online retailers.

Chapter 1

Introduction

The reflection which gave rise to the developments that will be introduced in this report started with a very simple question: could gravitation be a repulsive force under certain circumstances and what would it mean for gravitational mass to be negative? Even though there appears to be important difficulties associated with the possibility that a gravitationally-repulsive body may exist, particularly in the context of a general-relativistic theory, the idea of a symmetry which would have to do with the sign of mass or energy is certainly quite appealing aesthetically. Indeed, if the electric charge and all the other charges turning up in particle physics are allowed to be both positive and negative, why should mass or energy be restricted to positive values? What I came to realize, through a careful analysis of the assumptions behind the common idea that gravitationally-repulsive matter does not exist, is that there is actually a general misunderstanding surrounding the whole idea of negative energy in modern physical theory and that this is the single most important stumbling block that is preventing necessary progress from being achieved in several fields of fundamental theoretical physics. The objective of this essay is to clear up the misunderstanding and to provide a detailed account of the most crucial advances which are made possible by adopting a more consistent approach regarding some essential concepts related to time directionality and their relationships with classical gravitation theory and quantum theory.

I will, therefore, begin by revisiting the old problem of negative energy states and by explaining the difficulties which arise in the context of the current conception of negative mass. This will allow me to achieve a more consistent integration of the concept of negative-energy matter to the classical

theory of gravitation, by drawing on the analogy provided by the gravitational dynamics of voids in a matter distribution. I will show that traditional expectations, regarding the interaction of negative-energy matter with itself and with positive-energy matter, are inappropriate, because they violate the requirement that all physical properties be defined in a relational way. From this analysis will emerge an improved understanding of the notion of gravitational repulsion involving negative-energy matter as a form of dark matter whose existence must, under certain circumstances, be considered unavoidable from both a theoretical and an empirical viewpoint. An alternative set of axioms, which allows an appropriate and at last, consistent integration of negative energy states to physical theory will then be proposed. I will conclude this portion of my analysis with a reformulation of the relativistic gravitational field equations that provides the foundation for the first-ever bi-metric theory of gravitation that is truly symmetric under exchange of positive and negative energy states and which actually simplifies the original theory in the presence of a non-zero cosmological constant.

What allowed me to achieve a better understanding of the concept of negative-energy matter is the acknowledgment that there must exist a fundamental time-direction degree of freedom, independent from the thermodynamic concept of time direction. In such a context, it emerges that only the sign of energy defined in relation to a given direction of propagation in time is significant from a gravitational viewpoint. Once the significance of this insight was properly assimilated, it became possible to develop an alternative concept of time reversal that allows a reformulation of the discrete symmetry operations and a more consistent description of the changes occurring under a reversal of space- and time-related parameters. In order to achieve full consistency, it was necessary to introduce an additional set of discrete symmetry operations of a kind which had never been considered and which transforms a positive energy state into various negative energy states. Those developments then allowed the derivation of an exact binary measure for the entropy of the matter contained within the event horizon of a black hole that reproduces the results of the semi-classical theory in the case of elementary (Planck-mass) black holes.

As a consequence of the relatively long period of gestation during which the mere intuitive insights from which this work originates evolved into a revised, classical theory of gravitation, I was able to explore the consequences of some of the most decisive results which were reached in the course of that process for a rather large number of questions of fundamental interest. Thus,

I can now provide a complete account of the implications of this improved understanding of gravitational physics for classical cosmology theory and in the process achieve a better understanding of several issues related to time directionality. I will, in particular, provide significant new insight regarding the whole question of dark energy and dark matter and the related problem of the formation of large-scale structures. Still by making use of the results derived in the first portion of this report, I will then propose alternative solutions to some outstanding problems in theoretical cosmology which were originally addressed using inflation theory. I will conclude this part of my analysis by providing a definitive solution to the problem of the origin of time irreversibility which relies on a more accurate assessment of the measures of entropy associated with the gravitational fields of homogeneous and inhomogeneous matter distributions.

In the fifth chapter of the report, I will then offer a fresh perspective on several aspects of the problem of the interpretation of quantum theory which centers around a reconsideration of the significance of the requirement of time-reversal symmetry. Following a critical review of some early time-symmetric formulations of quantum mechanics, I will argue that a more consistent approach must overcome the contradictions of the orthodox interpretation of quantum theory that follow from its rejection of scientific realism. I will also show that the condition of time-reversal invariance provides strong enough a constraint to allow a realist interpretation of quantum theory to satisfy the principle of local causality in the face of quantum entanglement. Finally, in the second portion of my discussion concerning the foundations of quantum theory, I will explain that the existence of a maximum quasiclassical domain can only be predicted to arise and to persist, following measurement, once we consider the problem of the emergence of time in quantum cosmology from the perspective of the solution provided in the first portion of the report to the problem of the origin of thermodynamic time asymmetry.

1.1 Motivations

It must be mentioned that, even though I became interested in the basic idea underlying the developments discussed in this report based on mostly aesthetic motives, the actual reasons that later fueled my interest in developing a viable model around it were of a more pragmatic nature. In particular, I

saw the need that existed, but that few authors recognized, to reformulate the current classical theory of gravitation in a way that would be consistent with the possibility for elementary particles to be found in the negative energy states allowed by special relativistic quantum theories. Indeed, I had come to understand that the current interpretation of negative energy states as merely being those of particles propagating backward in time (the antiparticles), whose behavior is identical to that of ordinary matter from a gravitational viewpoint, was dependent on the *a priori* assumption that only some of those energy states were allowed. In other words, we had solved the puzzling problem of the prediction of negative energy states by postulating that those states were not allowed, without justifying this very assumption. But if we recognize that the whole spectrum of energy states predicted to exist by quantum theory can, in effect, be occupied, even if transitions between positive and negative energy states may not be allowed, then we need a classical theory of gravitation that is consistent with this requirement. However, further considerations indicated that the general theory of relativity is not entirely compatible with an appropriate notion of negative energy obeying certain theoretical requirements which must be imposed in order to achieve consistency.

Despite those difficulties, I believe that the imperative to provide an appropriate description of negative-energy matter should prevail over our willingness to leave untouched the current theory of gravitation, because I have recognized the inadequacy of the arguments against the physical nature of negative energy states, while I also understand that quantum theory constitutes a more appropriate basis to decide what states are allowed for elementary particles. Thus, I persisted in seeking to achieve this integration and as it turned out, this insistence was vindicated, given that I was able to develop an alternative framework that merely generalizes relativity theory in a very elegant manner, without affecting its basic mathematical structure, while allowing an appropriate description of negative-energy matter.

But I was also motivated by the desire to obtain a better agreement between theoretical predictions and astronomical observations concerning certain aspects of the gravitational dynamics of the universe. In particular, there was the exceptionally severe disagreement between most theoretical derivations of the expected average value of vacuum energy density and observational constraints on the upper (positive or negative) value of the cosmological constant. Very early on, I saw that the hypothesis that there should exist a usually ignored portion of zero-point vacuum fluctuations that would

interact (other than gravitationally) only with matter in a negative energy state could potentially provide a whole new class of contributions to vacuum energy density, which would be the exact opposite of those already considered in conventional calculations, and which could naturally allow an overall cancellation of all contributions, if some level of symmetry exists between the viewpoint of positive-energy observers and that of negative-energy observers. Here, again, I chose not to ignore, as most people did, what seemed to be the necessary conclusion that matter must be allowed to occupy the currently forbidden, negative energy states if we are to obtain a compensation for the known contributions to vacuum energy. Despite the apparent difficulties, perceived or real, associated with negative energy as a possible state of matter, it had become very clear to me that this was a hypothesis which had become unavoidable.

Finally, I also wanted to bring some much-needed clarity to the theoretical context in which we are to address the problem of the elaboration of a theory of the gravitational interaction compatible with the basic principles of quantum theory. Here I will show the essential role played by the discrete spacetime and momentum-energy symmetry operations (appropriately redefined and extended to comply with an improved concept of time reversal) in characterizing states of matter at the spatial scale and energy level at which we can expect the gravitational interaction among elementary particles to be as strong as the other known interactions. This will be achieved by demonstrating the relevance of those symmetry operations for a definition of the microscopic states of matter that must be taken into consideration in order to provide an appropriate measure of black-hole information and entropy. But I will also explain that one of the main consequences of the solution I have developed concerning the problem of the origin of thermodynamic time asymmetry is that it allows one to understand how a uniformly flowing time variable can emerge from the timeless equations of quantum gravity, thereby providing the metric of spacetime with a unique signature from which originates the causal structure of relativity theory.

1.2 Approach

Basically, the approach I will follow in this report consists in explaining how some specific aspect of the quantum world, namely the ignored possibility for both positive and negative energy states to propagate forward and backward

in time, changes our understanding of the *classical* theory of gravitation and allows to actually improve and simplify its formulation in a way that will have decisive consequences for the description of certain phenomena which are taking place on the cosmological scale. But, once there, I will go the opposite way and show how those original insights regarding cosmology shall affect our understanding of *quantum* physics and open up the way to a more pragmatic approach toward a quantum theory of gravitation. The level of this discussion is clearly philosophical, but remains very precise in its reference to quantitative aspects and concepts, unlike most philosophical essays concerning physics. Mathematical developments will be kept to a bare minimum, however, and will be introduced only when absolutely necessary and of utmost significance. This is obviously in contrast with the current tendency observed in the physical sciences to focus on technical aspects and to relegate epistemology to the backseat.

Concerning the methodology which is reflected in the style of this treatise, I must emphasize that I have been introduced to quantitative methods very early on, but I later came to realize that in the context where all the really useful mathematical developments that could be carried out in the field of fundamental theoretical physics have already been performed over and over again by competent people, real progress can only arise at the level of interpretation. Indeed, a fully consistent interpretation of the existing frameworks is currently missing, perhaps because the vast majority of competent researchers prefer to dedicate their efforts to more technical aspects, and this is restraining our ability to distinguish between what are viable developments and what is logically and empirically inappropriate. But as I do believe that the objective of a philosophy of science should be to elaborate and to justify, through logical arguments constrained by observational data, a globally consistent world-view, and as I'm convinced that it is only when this goal is successfully achieved that we are allowed to consider this particular vision of the world to be a valid representation of it, then this is the objective toward which I directed my efforts.

Furthermore, it is important to note that, if mathematical developments do not dominate the content of this report, this is also simply a consequence of the fact that, while I have achieved a crucial revision of the mathematical framework of relativity theory and a necessary improvement of the interpretation of quantum theory, I nevertheless ended up confirming the validity of the basic mathematical structures of both theories within a certain limit, so that practically no further mathematical developments were required. The

reader must be warned, however, that the density of information that is to be found in the text of this document is very high. In some cases, it took me years of dedicated reflection and careful investigation to gain confidence in the validity and inevitability of certain specific results which may be discussed only once in the main portion of the report, as otherwise the length of the treatise would be excessive. Therefore, you must pay attention to every detail of the discussion and be careful not to miss out some important information that may become necessary, later on, for understanding and appreciating the value of other elements of the discussion (this is also true concerning the footnotes, which provide essential, complementary information rather than just references). This, however, does not mean that the present essay is difficult to read, to the contrary. In fact, I tend to follow a rather educational approach and I do not avoid making statements and providing explanations that may appear obvious to some or even most readers, because I think that it is better to make too many unnecessary statements than to more or less willingly avoid making some which would have been useful. This approach should not be considered as condescending or as an indication that this work is intended mainly for a beginner audience.

Now, I must mention that I do recognize that the approach I followed in order to achieve the valuable results that will be described and justified here *is* different from that which is usually followed in theoretical physics. Indeed, very early on in my career, I was led to concentrate my efforts on questions of an epistemological nature and to rely on the expertise of specialists concerning certain technical aspects which are not essential to an accurate understanding of the issues on which I was concentrating my efforts. Thus, instead of assimilating all the complex machinery that allows to solve specific problems in various fields of theoretical physics, I was satisfied with studying problems of a more general nature that still required careful reasoning and analysis, but that were not considered serious work by most conventional researchers. I'm convinced that, if I had insisted on following a more conventional approach, I simply would never have been able to derive all of the important results that figure in this report. Indeed, achieving such a comprehensive understanding of the interpretative issues of so many different fields of fundamental theoretical physics while keeping in touch with the latest experimental advances in cosmology, particle physics, and quantum theory was a full-time occupation that required dedicated efforts sustained over a long period of time. But even more demanding was the task of actually reflecting on those issues and of exploring the effectiveness of various poten-

tial solutions to the multiple problems encountered, whether it was those which already existed or those which developed as a result of the tentative solutions I was myself proposing.

What I'm trying to say is that the kind of work I have accomplished requires specialization, but while most researchers develop very elaborate technical skills in one specific field of study, my specialization consisted in developing skills in analyzing certain general aspects common to several different fields of fundamental theoretical physics which all have to do with time directionality. If I had not focused my attention on questions of interpretation and had rather tried to develop all of the elaborate skills required to solve more specific problems in every field I studied, as I thought to be necessary when I began studying physics in a traditional academic environment, I would certainly have failed to contribute to our understanding of the physical world. The truth is that a certain level of technical expertise was required to achieve those results, but I was lucky enough that, when I first began to work at a more qualitative level, I had already developed most of that mathematical proficiency.

But the very fact that, for many researchers, the preceding comments will merely reflect incompetence indicates that, at the present epoch, theoretical physics has reached a point in its evolution which is similar in many regards to that in which natural philosophy ended up when it began deviating into mathematical idealism during antiquity. Indeed, it has recently been emphasized that the absence of philosophical underpinning that characterizes some currently favored approaches and the excessive recourse to mathematics in formulating physical theories (which is often achieved even at the expense of clarity or usefulness), has driven the field of fundamental theoretical physics into a state of stagnation. But this overly technical strategy is not a requirement of the scientific method and there is no need to use complex mathematics at *every* level of discussion and under *all* circumstances, especially when language allows sufficient or better clarity and contains clear references to precise quantitative constructs which have already been developed. In fact, I believe that there is a trend in the evolution of scientific research, from the first theoreticians who invented their own mathematics, to later physicists who made use of existing mathematical developments to build their models, and on to some present day physicist using already existing mathematical physics frameworks to produce further original insights, still building on what had previously been achieved.

It must be clear, however, that I'm not trying to deny the effectiveness,

or the usefulness, and certainly not the necessity of a quantitative approach to physics, but simply to emphasize that, in order to develop a globally consistent understanding of so many different aspects of fundamental theoretical physics, I had no choice but to follow an unconventional approach and to adjoin to mathematical reasoning the benefits, nowadays somewhat forgotten, of rigorous philosophical analysis. But, even though I would not myself have believed that one could achieve significant results by concentrating on interpretative issues when I started studying physics, which I did the usual way by learning about the mathematics of quantum theory, statistical mechanics and general relativity, it is through experience and by force of circumstance (although not as a result of mere inability), after having slowly and partly unwillingly deviated from the traditional path, that I began to understand that there is real value in such an approach, which I developed by making systematic a learning process that initially appeared to merely be a faithful, but irresponsible time-wasting improvisation. If the reader has enough courage to immerse herself in a similar experience and to loosen her grip on more traditional ways of achieving deep understanding, while nevertheless being ready to spend a minimum amount of effort to follow simple logical arguments, I can assure her that she will not be deceived and will learn useful physics, which is not so bad already by today's standards.

It must be noted, however, that due to the unusually large number of disciplines affected by the developments which will be introduced in this report, it may be difficult, at first, to gain a proper appreciation of the value of some of the most radical ideas that will be discussed, because a good portion of the arguments that motivate results which are discussed in the first sections will only become fully understandable after reading the latter sections of the report. Therefore, it is important to keep in mind that, if some conclusion or hypothesis may at first appear to be unjustified, it is more likely to be a consequence of the fact that not all of the arguments that will make it a valid proposal have yet been discussed. You can trust me for having spent a considerable amount of time verifying the validity of my claims in order to create the simplest and yet also the most universally valid explanation of the facts considered here, sometimes by rejecting my own earlier conclusions. I do believe that anyone who carefully reads the entire document will be able to recognize that, in the end, there is little choice, if one wishes to obtain a globally consistent picture of physical reality, but to accept the validity of some ideas that may at first perhaps seem a little extravagant, or to the contrary just plain reactionary.

In fact, I must admit that I have myself gone through phases when I realized that I had to revise my understanding of certain concepts associated with a certain field of theoretical physics based on new developments I had achieved in another, not totally unrelated field and sometimes this revision itself had an impact on the validity of other results in other fields, because, at this level, nothing can really be conceived independently of anything else. But it takes time to get a proper understanding of the whole picture in which everything agrees with everything else. There may, however, remain aspects which have not yet been fully integrated into the global picture I developed, simply due to the fact that I did not had the chance to rethink their significance in the context of all the other advances. This is perhaps unavoidable given the considerably large scope of the subject of time directionality, which is relevant to so many fundamental aspects of physical theory. But I have done my best at providing the most exhaustive account of the progress achieved and at identifying the various relationships between the many insights that form the substance of this report. Yet, so many shifts in understanding as has occurred during the process of developing this more consistent picture of so large a portion of physical reality may have left some consequences of even the most decisive insights on various other aspects of the global picture not fully assimilated. Given that I wanted to publish the results of my inquiry within my lifetime, I had no choice but to eventually let the outcome which I believe to be the best achievable account of my research go out for others to benefit, but also to criticize for what it may still contain of imperfections.

1.3 Historical context

There are many similarities between the current state in which science finds itself and those through which it went at other crucial turning points in its history. Indeed, the situation we have now arrived at is characterized by an accumulation of unanswered questions which creates an impasse that prevents further progress from being achieved. It is my belief that answering just a few key questions among those will greatly facilitate future theoretical research. When we examine the present situation in physics it becomes clear, in effect, that if there are questions which we are justified in not being able to answer right now, because they are related to what may be said of reality under conditions which we cannot yet reproduce in experiments (think of

trying to explain the origin of the free parameters of the standard model of particle physics), there are also questions which have to do with known difficulties which we have puzzled about for a long time and which we have no reason to believe further experiments may be particularly useful in helping resolve. But those are problems whose existence is often simply unknown to most people or which are sometimes considered to have already been solved, while careful examination shows that this is not always entirely the case. Most current programs in fundamental theoretical physics are focused on trying to solve the problems raised by questions of the first type and this is unfortunate, because here is precisely the domain in which progress is limited by technological constraints of a practical nature and the cost of achieving the required experiments. Very early on I recognized that, if I was going to enable progress to be made in physics, I had to concentrate on questions of the latter type, where progress could occur not only in my lifetime, but also as a consequence of the success or failure of my own enterprise.

Among the questions we may have hope to answer using our current knowledge is the question I mentioned earlier on as having being that which launched the reflection process from which this report emerged. It is, in effect, one of those unsolved questions whose very existence is usually unrecognized or which is considered to have already been solved, while this is clearly not the case, as I will explain later. You will not see it mentioned in most accounts as being one of today's open questions in physics, but it is one of the most important categories of question regarding classical physics and a field most people currently consider to be free of major difficulties. This problem of negative energy states could actually be called the 'classical gravitation theory problem' or the 'general relativity problem', because properly answering that question requires introducing slight modifications to that theory, which actually consist in a generalization of its own founding principles. This is the first question I will address in this report and satisfactory answers will be provided to the mostly unrecognized issues it currently raises. Doing so will require reconsidering the significance of certain aspects of the problem of vacuum energy and gaining a new understanding of the gravitational effects of homogeneous and inhomogeneous matter distributions that can be extended to our description of the physical vacuum.

An additional category of questions, which is also related to classical gravitation theory, can be collectively described as the 'cosmology problem'. It asks what is the origin of the constants of the standard model of cosmology, what is dark matter and what is dark energy, how are we to resolve the

flatness and horizon problems, and what explains the existence of certain visible large-scale structures which appear to have developed at too early an epoch or on too large a scale to be explainable by conventional theories? It also asks why it is that the energy which is contained in zero-point vacuum fluctuations is so low in comparison with the very large value that is provided by most theoretical estimates? Finally, it asks whether there was a beginning to time in the past and why it is that matter was present in the first instants of the Big Bang? While it is often considered that some of those questions have already been answered by developments like inflation theory, I will explain that there remain important unresolved issues in this context and that we are justified in seeking alternative answers, which I will show do exist. In fact, even though the objectives I had in mind when I started this research project were quite humble, in the end I was able to provide original solutions to nearly all aspects of the cosmology problem.

But I will also address a further category of questions that is usually considered to regard classical physics, but which actually sits right at the interface between the classical theory of gravitation and quantum theory. This is the traditional question of the origin of the statistical properties of matter which are reflected in the unidirectionality of the evolution in time of systems with a large number of microscopic degrees of freedom. Given that this problem of the origin of the thermodynamic arrow of time can be traced back to the peculiar properties of the distribution of matter energy which existed during the first instants of the Big Bang (as I will explain), it follows that the question of the origin of the unidirectionality of thermodynamic processes is, in effect, also a question for cosmology and as such, it will hugely benefit from the insights I have gained while solving other aspects of the cosmology problem. What was somewhat unexpected to me, though, was the realization that answering those questions actually constitutes an essential condition for addressing an additional and apparently unrelated category of questions. Indeed, as I have mentioned above, the solution I will propose to the problem of the origin of time irreversibility turned out to be essential for developing a proper understanding of quantum theory and in order to provide a satisfactory explanation for the emergence of a quasiclassical world and this is why I will discuss the problem of the interpretation of quantum theory as part of my analysis of the question of time directionality.

Richard Feynman has emphasized the fact that, acquiring knowledge about one physical law, or getting insight into one important problem, and being confident in the validity of those developments, often allows us to find

other laws. I have been able to experience the validity of this remark while my understanding of physics progressed. Indeed, by carefully applying the knowledge I had gained by examining the problems which I initially tried to address, I was allowed to achieve further improvements to our understanding of distinct and apparently unrelated issues, always based on an unassailable confidence in the validity of what I had already been able to understand. This report, therefore, provides a complete account of the subject it covers and anybody with a basic knowledge of relativity theory, cosmology, quantum theory, and statistical mechanics will be able to benefit from the revised understanding that it brings to this entire domain of scientific research. It is my hope that by reading about what I have found, some young and not yet indoctrinated mind will be inspired to explore even more remote territories and bring forth a significant shift in our understanding of reality that will prove, again, that it is only by wandering far from the beaten track that one can gain the perspective necessary to see the vast landscape that goes unnoticed to those who do not dare to deviate from the normal course of research imposed by the practices which are of common use at a given epoch.

1.4 Organizing principles

Every successful venture into unknown territory requires relying on the appropriate beacons and guidelines and this is particularly true when the voyage takes you to the boundaries of traditional certainties and brings you to question some essential aspects of what had previously appeared to constitute a fixed background for scientific exploration. I would, therefore, like to briefly describe what were the essential principles that guided me on developing the revisions of classical gravitation theory and quantum field theory which are described in this report. It must, first of all, be understood that those principles were not given as preconditions imposed on any vision of the world, but actually developed alongside improvements in my and other people's knowledge and understanding of that portion of physical reality we actually experience and through the possibility that this probing allowed of inferring the regularities present in an even larger and more encompassing domain of the same reality.

My awareness of the importance of the first of those principles developed mostly in conjunction with my appreciation of the requirements imposed by the classical theory of gravitation. Indeed, it is while tackling the problem

of negative energy that I realized the importance (emphasized by others in a different context) of a relational view of the physical attributes of objects and that I understood the real significance of the requirement of relativistic invariance. This allowed me to perceive the true meaning of Einstein's insistence that the objects of physics must be conceived of only in relation to the spacetime structure to which they belong, because I saw that the metric properties of space and time must be understood to depend on the sign of energy of an object (as will be explained later), in contrast with what one would expect traditionally. Thus, if a determination of the relationships between physical objects in different spatial locations or states of motion is possible only when we determine the common spacetime structure shared by those objects, then the fact that the spacetime structure itself is dependent on the nature of the objects means that the relationships between them are dependent on their nature and in particular, their energy signs. It therefore appeared to me that it is not merely the position and state of motion of an object which require a relational description, but that any physical quantity must always be defined or characterized only in relation to similar quantities of other objects present in the same universe (the physical attributes of a system enable to characterize it merely in relation to the similar attributes of other systems and those relationships are determined through the use of reference systems).

When I tried to understand what could logically impose such a requirement, I slowly came to realize that it is the very fact that it would be meaningless to relate some physical quantity, in order to define its value, to a reference point not part of the same physical universe. Indeed, in the absence of a well-defined, continuous network of causal relationships that would extend to those immaterial reference systems, there can be no meaningful definition of the physical quantity involved, because physical relationships are material relationships and an object cannot be put into relation with something that is not part of the same causally related ensemble (the universe) to which it belongs. This requirement of a relational definition of physical quantities will have enormously important consequences on many aspects of the developments to be discussed in the following chapters. It is important to understand, however, that the necessity to define the value of physical attributes in a relational way does not imply, as some authors have suggested, that nothing can exist other than the physical reality we observe in our universe. Indeed, it must be clear that what I have found is that there can be no reference, by observers in a given universe, to physical attributes not

related to one another by the network of causal relationships belonging to their own universe. But this does not mean that other such ensembles, or universes, cannot exist as logical possibilities, with similar, purely relational, and mutually referring properties, objectively distinct from those existing in the observed universe.

This remark illustrates the importance of another broad requirement that slowly emerged as being unavoidable for a solution to the problems which will be discussed in the fifth chapter of this report. There is, in effect, a tendency, nowadays, to designate as metaphysical every aspect of reality which may be impossible to probe through direct observation and to conclude that such aspects are not worth the attention of the scientific community. What I have come to understand is that the self-imposed requirement of systematically characterizing as metaphysical any notion that refers to aspects of physical reality which may not be directly accessible to observation is actually a mild form of solipsism and constitutes one of the most serious obstacle on the way to developing more accurate models in fundamental theoretical physics. In fact, I think that the greatest challenge with which science is currently faced may well be that of surmounting the obstinate refusal to accept as a legitimate object of scientific inquiry what cannot be *directly* observed by the means of measuring instruments and as physically meaningful what lies outside the limits of observation of a given observer (think of the reality behind event horizons for example). In this particular sense, the success of science might, in the end, depend on our willingness to adopt a position analogue to scientific realism and opposite to instrumentalism, concerning ultimately the idea that something really exists outside our immediate domain of perception of reality.

This requirement is not so different from the original condition of objective reality which was advocated by Einstein and which he proposed in an attempt to demonstrate the validity of an approach based on the hypothesis that reality actually exists, even when it is not subjected to direct observation. But given that, in the physical sciences, objectivity has rather come to characterize any conception of reality that is derived solely from empirical knowledge and observation, then it would not be appropriate to use the term ‘objective reality’ in order to refer specifically to a reality that is not directly observable under all conditions, even if the nature of this reality was still *derived* from experimental facts. Thus, I cannot avoid having to speak about a *realist* conception of reality as being essential to a consistent interpretation of quantum theory, even if that may appear tautological, as there does not

exist a more appropriate term to denote this kind of approach. It must be clear, however, that it cannot be required of such a reality that it be classical in nature, despite the fact that it would be characterized as objectively real (in the philosophical sense). Anyhow, I think that this scientific ‘realism’ must be considered a necessary ingredient for the elaboration of an accurate understanding of the nature of reality at a fundamental level and this is what motivates my position with respect to certain unresolved issues regarding the problem of the interpretation of quantum theory.

Such a conviction, however, should not be confused with a belief in the validity of theoretical constructs that have no experimental justification, which does not constitute a desirable position to hold on to and which would actually consist in the exact opposite of the viewpoint I’m defending here. What I’m suggesting, in effect, is that it may sometimes be appropriate to extend the validity of what we know to be true with absolute certainty to a larger domain of reality where this validity may not be *directly* assessed and not that it would be right to try to extend the domain of validity of a description for which there does not yet exist any empirical evidence. In other words, if we are justified in extrapolating beyond the domain of direct observation, as may be found necessary, principles and notions which we have good reasons to believe are indeed valid, it would be wrong to take advantage of the absence of observational data to try to justify hypotheses which cannot yet be independently corroborated and which may therefore have no validity whatsoever from a scientific viewpoint. Those considerations will have decisive consequences for the formulation of an interpretation of quantum theory that contains no contradiction when considered in the broader context of the representation of reality that emerges from the progress which will be achieved, in the first portion of this report, in solving other long-standing problems in the fields of gravitational physics, cosmology, and statistical mechanics.

Chapter 2

Negative Energy and Gravitation

2.1 The negative-energy problem

Regarding the question of negative energy, the current situation has much in common with that in which we were at the turn of the previous century with regard to the quantization hypothesis. There was, in effect, some reluctance, initially, to recognize the validity of the original suggestion by Max Planck that energy is quantized, despite the fact that this proposal would have solved the problem of black body radiation. The trouble was, of course, that recognizing the validity of the quantization hypothesis required abandoning classical physics. There is a similar dilemma with negative energy today because, as I will show, the hypothesis that matter can be found in such a state has the potential to solve many important problems facing theoretical physics, but those benefits come at a price which may, at first, appear to be too high. Indeed, the introduction of negative-energy matter as a concept somewhat distinct from that which is currently favored (which I believe is required in order to allow it to be consistent from a basic theoretical viewpoint) seems to imply that general relativity has to be abandoned. But rejecting a theory so well established and so beautifully simple as general relativity is not something that most people would do without very good motives. Yet, while the current assumptions concerning the rules governing negative-energy matter (if it was to actually exist) may appear to better agree with relativity, they actually contradict some of the basic principles on

which this theory is founded, therefore making it just as untenable. We must then either abandon the idea that negative-energy matter can exist, or else provide a better interpretation of negative energy states which may force a reformulation of relativity theory itself. But I will show that the conclusion that the latter alternative is the only viable one is not necessarily as dramatic in terms of its consequences as one would expect, because what is required, in this context, is mainly a reinterpretation of the equivalence principle and not a rejection of the whole mathematical framework of relativity theory.

There is, however, an additional problem regarding the hypothesis that negative energy states cannot be rejected, which is that there appears to be no observational evidence for matter in such a state. But here also there is an analogy which should teach us a lesson. This is the case of the neutrino as a massive particle. For a long time, when I was reading physics papers or any book on the subject of particle physics, I could see that it was nearly *always* assumed, more or less implicitly, that the neutrino is massless, as if this was a fact, while there was absolutely no evidence that this is actually the case and it was merely the difficulty to prove that the hypothesis is wrong that justified that everyone just assumed that the neutrino is massless. But just as for the idea that negative-energy matter does not exist, I thought that it was incorrect to simply assume that the neutrino is massless when this could not yet be considered a fact. Thus, I always kept an open mind about those issues, because I saw that there were strong arguments (usually not recognized) for rejecting those commonly held assumptions and in the case of the neutrino at least it appears that this position was justified. In fact, I will later explain that there are very good reasons to expect that it won't be easy to confirm the existence of negative-energy matter, because, as I have come to understand, even the portion of it that may still be present in the universe today should not be directly observable, just as the more common, hypothetical dark matter. Thus, if I'm right, the implicit assumption that negative energy is forbidden would be just one of those 'reasonable' assumptions which we should be careful not taking too seriously.

The problem of negative energy also has a parallel in a distinct but not entirely unrelated problem which is that of the origin of the arrow of time. Indeed, it was suggested by some eminent figures that the problem of irreversibility could be solved by integrating some fundamental element of irreversibility into the formalism of even the most elementary physical theories. This would seem to be justified by the fact that the problem of time asymmetry has been known to exist for a long time and no acceptable solu-

tion to it that would be based on boundary conditions imposed on otherwise time-symmetric evolution has ever been found. But again, I think that the difficulty to prove a hypothesis (that time asymmetry can arise from time-symmetric physical laws) should not be taken as evidence that what may perhaps be its only alternative (that time asymmetry is fundamental) is right. In the case of negative energy, we are also in a situation where we have built into the very formalism of our most fundamental theory of matter (which currently is quantum field theory) the apparently necessary, but clearly unjustified (from a theoretical viewpoint) hypothesis that only positive frequencies (associated with positive energies) are allowed to propagate in the future (the constraint on negative frequencies being merely that they must propagate toward the past).

However, I think that the fact that this artificial restriction appears to be valid does not imply that positive frequencies cannot propagate backward in time or that negative frequencies cannot propagate forward in time, but merely that if there exist two kinds of matter related by their opposite energy signs (the frequency signs relative to the direction of propagation in time) then, for some reason, they can only interact with matter having the same energy sign (I will eventually explain why such a limitation naturally occurs). This absence of interaction or interference (in the classical sense) is what really justifies the observation that quantum field theory only deals with matter of one energy sign under most circumstances (when gravitation is not involved). But given that I'm suggesting that energy sign is a relatively defined physical property, so that there is no absolute (non-relational) distinction between positive- and negative-energy matter, then it must, in effect, be concluded that there cannot exist a constraint that would impose that negative-energy matter, and only matter with such an energy sign, cannot exist under any circumstances, if positive-energy matter itself is allowed to exist, as required, because it is not even possible to identify the distinguishing property, specific to negative-energy matter, that would justify that its existence be ruled out in such a way. Thus, I'm allowed to conclude that any attempt at getting rid of the apparently intractable problem of negative energy states by simply imposing a constraint to be applied on the formalism itself is misguided and unnecessary, because, indeed, once an appropriate understanding of the true nature of negative-energy matter is available it becomes apparent that a restriction on allowed frequencies is no longer necessary. In fact I believe that the same can be said of the problem of irreversibility, because in chapter 4 I will show that the thermodynamic

arrow of time is not an intrinsic feature of fundamental physical laws, but instead originates from an unavoidable constraint that applies on the boundary conditions at the Big Bang.

In the context where we must recognize that there is no theoretical motive to reject the possibility that negative-energy matter may be present in our universe it becomes apparent that one often mentioned argument that must definitely be rejected concerning the nature of the gravitational interaction is the idea that the strength of gravitation on the largest scales is a consequence of the ‘fact’ that this interaction is always attractive. This is a conclusion which is usually assumed to follow from the observation that there does not exist negative gravitational charges (negative-energy matter is assumed not to exist). Yet, what actually explains the fact that gravitation is a dominant force on larger scales (in addition to its long-range property) is not the absence of matter in a negative energy state, but the simple fact that gravity is attractive between objects with the same positive gravitational charge, that is, between objects with a positive sign of energy. Thus, if gravitation dominates over electrical forces on astronomical scales, this is really a consequence of the fact that while identical electric charges tend to disperse under mutual electrostatic repulsion, positive energies have a tendency to coalesce and to accumulate under mutual gravitational attraction and the fact that electromagnetism is already known to have both positive and negative charges has nothing to do with the fact that those charges do not so readily accumulate, because even if there were only positive electric charges they would not cluster, because identical electric charges mutually repel one another and the possibility for such opposite charges to cancel out actually facilitates an *accumulation* of those charges, but only in neutral configurations and under the influence of gravitation.

It must therefore be understood that there is no requirement for gravitation to always be attractive merely on the basis of the fact that its existence can be felt despite its extreme weakness, as is sometimes suggested. Indeed, if it was found that there actually exist negative-energy particles, the possibility for energy to cancel out would not necessarily prevent the accumulation of matter with one or another energy sign, because negative-energy matter may also be gravitationally attracted to itself (despite what is usually assumed) and could therefore also be subject to accumulation. To summarize, what makes electrical forces negligible on the large scale is the fact that identical electric charges do not attract one another and therefore do not accumulate as may identical gravitational charges. Instead electric charges of opposite

signs are attracted to each other and immediately cancel out, therefore preventing further accumulation, at least under the influence of electric forces. But this does not mean that gravitation would be submitted to the same fate if negative-energy particles were found to exist, because it may well be the case that gravitational charges with the same sign always attract one another, given that this is already known to be true for positive-energy matter, and this would not even forbid opposite-energy bodies from gravitationally repelling one another. The frequently encountered remark that gravitation is attractive for all particles should therefore be understood to mean only that it is attractive for all currently known forms of matter.

Thus, again, the observation of large accumulations of positive-energy matter is not an argument against the existence of negative-energy matter. But it is also true that the apparent absence of large accumulations of *negative-energy* matter would not necessarily mean that such matter cannot exist, even if we were to assume that this matter gravitationally attracts matter of the same kind. Indeed, it may turn out that this matter is dark and given that it may also be repelled by positive-energy matter (even if this is not what we usually assume), then we might be justified to expect that it should be located mainly in regions of the universe where the density of positive-energy matter is the lowest. Therefore, negative-energy matter would be virtually absent from regions where positive-energy matter is more abundant, like that in which we are located, and this would help explain that we have never noticed its existence. What's more, there may be other reasons to recognize that baryonic negative-energy matter, even if it is allowed as a stable form of matter, and even if it has had decisive consequences on cosmic evolution, may no longer be present in large quantities in the universe today. I will explain later why the assumptions discussed here concerning the nature and the abundance of negative-energy matter should, in effect, be those which are retained, thus confirming the validity of the above explanations as to why it is that negative-energy matter *appears* to be absent from our universe (even though it is not). It will then be clear that, theoretically, it is to be expected that if negative-energy matter exists it should have the properties which are responsible for our very ignorance of its existence.

I think that what must be recognized above all is that the commonly held view that the occurrence of negative energy in a theory is necessarily always indicative of a problem is not rationally motivated and that it is not true that all traces of negative energy *must* be eradicated, at all costs, whenever they

are encountered. Dirac, at least, understood that the prediction of negative energy states could not be ignored and tried to provide an explanation for the absence of transitions to such states [1]. His solution, based on the idea that negative energy states are already all occupied, was not satisfactory, but at least he did not simply reject the possibility that negative-energy matter might have to be considered real. There is no motive to argue, as people often do, that negative energy is totally unacceptable, other than the difficulty to find an appropriate interpretation that would be compatible with empirical facts for this logically unavoidable counterpart to positive energy. In the absence of a theoretical justification for the absence of negative-energy matter I think that the only appropriate approach would be to seek to find out why it is that we never observe matter in such states, rather than try to build that assumption into a then necessarily incomplete theory of quantum fields. In this particular sense, it is significant that the prediction of antiparticles was a by-product of Dirac's original interpretation of negative energy states, because this contributed to the belief that the discovery of antiparticles constitutes a definitive solution to the negative-energy problem. But, given that Dirac's interpretation was later found to be inappropriate, I think that we need to recognize that, in fact, antiparticles can only be one particular aspect of a complete solution to the problem of negative energy, which therefore remains unsolved.

In any case, it must be understood that, even if we were to succeed in justifying that it should be imposed that there cannot be transitions from a positive energy state to a negative energy state, we would not have solved the problem of negative energy. This is because such a restriction would merely impose that no positive-energy particle can turn into a negative-energy particle (and vice versa maybe), but there would be nothing in that constraint to forbid a particle to already be in a negative energy state, in which case we would still need to provide a consistent description of the properties of matter in such a state and to justify that we do not observe those negative-energy particles under most conditions. In fact, I will later provide arguments to the effect that just such a restriction on energy-sign-shifting transitions is to be expected to occur very naturally, even if negative-energy matter must indeed be allowed to exist. Anyhow, the fact is that, if there is no reason to assume that some restriction applying to energy sign reversal would forbid *positive-energy* matter from existing, then there cannot be more justification in assuming that such a restriction forbids negative-energy matter from being present in the same way under certain circumstances. I must insist

again that there is no reason to assume that the concept of negative energy is problematic all by itself and that negative energy must be avoided systematically, because the only requirement, regarding negative energy states, may be that there cannot be transition to such states by a particle in a positive energy state and this only when the transition would be to a state of negative energy propagating forward in time. Such a requirement is necessary to keep positive- and negative-energy matter virtually isolated at the elementary-particle level, so that the experimental constraint of an absence of interference from negative-energy matter into the theoretical predictions involving positive-energy matter can be satisfied.

I do understand, of course, that there are a number of issues associated with the possibility that matter may occupy negative energy states. Of particular concern would be the issue of ‘vacuum decay’ or the apparent problem that all positive-energy particles should fall within a very short interval of time into the available negative energy states by releasing a compensating amount of positive-energy radiation, if those states are not assumed to be forbidden. In fact, this problem would seem to affect negative-energy matter itself, even if transitions to negative energy states by positive-energy particles were found to be impossible. This is of course the difficulty that motivated Dirac’s problematic proposal that those energy states should already be nearly completely filled, so that no further decay should occur. But I will show in later portions of this chapter that this problem and also some others which may seem to arise in relation to the possibility for negative-energy matter to exist in a stable form are merely a consequence of the inappropriateness of the current interpretation of the concept of negative energy. In fact, it will be shown that it is not even necessary to assume (in order to prevent a decay of the vacuum) that negative energy states cannot be reached by matter in a positive energy state, because even matter already in a negative energy state cannot be assumed to fall to even ‘lower’ energy states.

I also recognize that the tentative interpretation of negative energy states that came to replace Dirac’s solution does, in effect, provide some level of relief in that it at least allows to take into account those negative energy states that cannot be ignored as they actually interfere with processes involving ordinary matter. This is because we are indeed allowed to consider that antiparticles are negative-energy particles propagating backward in time. But even under that particular interpretation, antiparticles can still be conceived as ordinary particles (submitted to normal gravitational interactions) from

the forward-in-time perspective relative to which their energy is positive and therefore they cannot be considered to provide an interpretation of negative energy states of the kind that would be truly significant from a physical viewpoint. Again, the exclusion of true negative energy states may appear to be justified from an observational viewpoint, but it still constitutes an arbitrary rule which would at least require an explanation, as there is no consistency principle behind it. It is therefore quite amazing that so many otherwise well-informed authors suggest that no negative-energy, or negative-mass particle can exist, as if this was an obvious and unavoidable conclusion. It must be clear that I'm not complaining about this situation, I merely want it to be recognized for what it is, because I will take a different course and it should be understood that I'm not doing this without good motives or out of a fondness for hopeless, exotic or eccentric ideas.

I must therefore mention that I'm aware that the originators of the steady state theory of cosmology once also criticized (based on distinct motives) the traditional position according to which the existence of negative-energy matter is forbidden. But if I do find this criticism to be valid and appropriate, I do not, however, find suitable the whole concept of negative-energy matter (which is actually very traditional) proposed by these authors, nor do I agree with the objectives they unsuccessfully (given the failure of steady state cosmology) sought to achieve by using this otherwise interesting idea. I think that the fact that the hypothesis that negative-energy matter may exist was historically associated with such failed theoretical models and was also developed into many different inconsistent formulations lacking any epistemological support is more than anything else responsible for the state of suspicion and confusion that currently surrounds the whole idea of negative-energy matter. The objective I will try to achieve in this chapter will therefore be to clarify the situation regarding what should be expected regarding the properties of matter in a negative energy state and to demonstrate the validity of the concept itself, in the context where it is properly defined and justified.

2.2 The time-direction degree of freedom

What emerges from my re-examination of the assumptions behind our current understanding regarding the possibility that particles may occupy negative energy states, is that we must first recognize that, for any elementary parti-

cle, there exists a fundamental degree of freedom related to the direction of propagation in time of its charges, including the gravitational charge, that is to say, including energy. The existence of such a degree of freedom means that a positive charge can, in effect, be positive either in relation to the positive direction of time, if such a charge propagates in the positive direction of time, or in relation to the negative direction of time, if the same positive charge propagates in the negative direction of time. But the particles so characterized would be physically different from one another. It is not possible therefore to completely specify the physical properties of a particle at a given instant by simply providing the sign of its charges independent from their direction of propagation in time. But given that a particle can actually be identified by the charges (including energy) it carries (it has no other physical properties except for its momentum, position, and spin at a given time) this means that the apparent nature of a particle may depend on whether it propagates its charges in the positive or the negative direction of time, that is, it may depend on whether it is itself propagating forward or backward in time¹. The physical attributes of a particle can only be unambiguously defined in relation to the direction of time in which this particle propagates and this is true also for energy.

This is what the insights gained by considering the consequences of the relativity of simultaneity for the quantum description of particle interactions should be understood to imply. Indeed, it is the fact that some processes involving the exchange of a virtual particle of interaction cannot be assigned a unique definite order of occurrence in time that renders the notion of particles propagating backward in time unavoidable. This is because the emission and absorption events of such an exchange process are space-like separated, so that their order of occurrence in time is dependent on the state of motion of the observer. Thus, what would appear, for one observer, to be a conse-

¹I'm here considering a particle in a semi-classical way, as if we could always associate with it a definite position and momentum, even though it is clear that actual knowledge of those conjugate attributes cannot be obtained at the same time. This idealization simply allows to gain insight into what would be the properties of an elementary particle if it could be observed at the energy scale of an actual macroscopic body, while still carrying a mere unit of its other charges. We may alternatively consider a real macroscopic body and assume that it has physical properties that evolve in a perfectly coordinated fashion, with all its charges necessarily propagating in the same direction of time at all times (therefore acting as one 'macroscopic' charge), but such a viewpoint is actually even less realistic than the former idealization (for reasons that will appear more clearly later on) and would change nothing to the following conclusions.

quence of the emission of some particle carrying a negative charge, would, for another observer, appear to be a consequence of the absorption of a similar particle carrying a positive charge, which certainly requires the sign of charge to be dependent on the perceived direction of propagation in time. Given the undeniable validity of this viewpoint, the only argument that could still allow one to reject the reality of a degree of freedom associated with the direction of propagation in time would be one based on the second law of thermodynamics and the apparent impossibility for a macroscopic body to ‘travel’ backward in time. It appears, however, that this argument is not valid, because the thermodynamic constraint only applies to the flow of information as it occurs through the formation of records and in no way forbids individual particles from propagating backward in time as long as they are not involved in processes which (collectively) would allow information to be transferred from the future to the past (I will better explain what motivates this distinction in the first portion of chapter 5). It is therefore merely this limitation on the flow of information that explains the fact that our experience of reality has made us suspicious of the possibility that objects themselves (or particles) can propagate backward in time and not the actual impossibility of such an occurrence.

In such a context, the possibility to distinguish the sign of a charge, including energy, would depend on the possibility to determine the direction of propagation in time of this charge. Thus, even independently from the argument based on the relativity of simultaneity, we may consider that the sign of charges and in particular the sign of energy is defined only in relation to the state of motion of the particle carrying those charges, where ‘motion’ is here relative to time instead of space. But if we may also assume that the attribution of a direction of propagation in time is merely a matter of convention, because all that can be asserted is whether any two particles are propagating in the same direction of time or in opposite directions, as I will suggest later, then it would appear that the sign of energy itself would become a relative notion, dependent on which direction of time is chosen as that in which a given particle propagates. In this particular sense we would have to recognize that associated with the relativity of ‘motion’ in time there is also a relativity of the sign of energy.

Acknowledgment that the sign of energy is a relative property actually allows one to reject the validity of the constraint usually imposed that all energy must be positive, because it means that, even what appears to be positive energy according to one particular convention for the direction of

propagation in time, is actually negative energy according to an alternative choice for the same time-direction parameter. The possibility for particles to propagate backward in time, which is made unavoidable by the fact that backward-in-time motion is actually required for a consistent understanding of the constraints imposed by a relativistic treatment of quantum processes, as mentioned above, therefore actually implies that negative energies must also be allowed in physical theory, because even what we usually describe as a positive-energy particle could be redefined as a negative-energy particle if we were to also assume, as a matter of convention, that the direction of propagation in time of the particle is opposite that which is usually (more or less implicitly) assumed. Negative energies must be considered to be possible states of matter, even if only for particles propagating in the backward direction of time. This dependence of energy sign on the assumed direction of propagation in time is what allows antiparticles to actually be described as particles propagating backward in time with negative energies and unchanged non-gravitational charges, as Feynman once suggested [2], even if we are also allowed to consider those particles as positive-energy particles with reversed non-gravitational charges propagating in the usual forward-in-time direction.

What is essential to understand, here, is the dependence of the value of any charge, including energy, on the direction of time in which this charge is assumed to be propagating. Thus, simply saying that a particle has positive electrical charge or positive energy doesn't make sense. We must also always specify the direction of propagation of this energy with respect to the time parameter. What appears to be a positive charge or a positive energy relative to the positive direction of time would be a negative charge or a negative energy relative to the negative direction of time. Thus, all those energy signs are merely established on the basis of practical conventions and can never be asserted in an absolute fashion. It must be recognized, however, that if the energy of an electron is by convention considered positive relative to the future direction of time in which it is, again by convention, assumed to propagate, then the energy of an anti-electron must *necessarily* be considered negative relative to the past direction of time in which it must, under the same convention, be assumed to propagate. It is merely because we ignore the requirement to describe the positron as propagating backward in time that we can attribute to it a positive energy (and a positive electric charge). As a consequence, it would seem that even on the basis of current observations we would not be allowed to assume that particles are forbidden from occupying properly defined negative energy states.

Yet despite the unavoidable character of this conclusion and even in the face of the enormous simplification of our world-view that is made possible by the hypothesis of the existence of a fundamental degree of freedom related to time direction, it is still often suggested that the interpretation of antiparticles as particles propagating backward in time with negative energy is merely a mathematical artifact and corresponds to nothing real. But I think that this attitude is similar to that of nineteenth century philosophers and scientists rejecting the hypothesis of the existence of atoms, even in face of the overwhelming evidence in favor of this concept, supposedly because the atoms could not be seen directly, but actually because of an unjustified prejudice in favor of a continuous, macroscopic description of matter. Given the above discussion concerning the relative nature of energy sign, I think that it is clear that there is no basis for assuming, as is often done, that the negative energy of antiparticles as particles propagating backward in time is not real and that those particles are merely ‘ordinary’ particles which happen to be carrying opposite non-gravitational charges. If we are allowed to describe antiparticles as particles propagating backward in time, then we must recognize the existence of negative energy states.

It must, in this context, be understood that the commonly met suggestion that all physical properties are simply reversed for an antiparticle (by comparison with those of the associated particle) is wrong, because the signs of all physical quantities are dependent on the direction of propagation in time and we would at least have to specify with respect to which direction of time the various quantities are to be assumed reversed. Indeed, even from the viewpoint where antiparticles are assumed to propagate in the same direction of time as do regular particles we would have to admit that energy is not reversed for an antiparticle, otherwise a pair-annihilation process should release few or even no energy in the form of radiation, contrarily to what is routinely observed. Also, if we do consider instead the viewpoint of an antiparticle’s true (when ordinary particles are assumed to propagate forward in time) direction of propagation in time, then energy would indeed be reversed as I already mentioned, but all non-gravitational charges far from being reversed would have to be considered rigorously unchanged given that from the forward-in-time viewpoint they actually appear to be reversed, while from my perspective the sign of charge is a relative notion, dependent on the assumption that is made regarding the direction of propagation in time of a particle.

Thus, what appears to be a positively-charged particle in relation to an-

other particle propagating forward in time would actually appear to be a negatively-charged particle in relation to yet another particle propagating backward in time and the same would be true of energy sign. Those relative alterations of the sign of charges occurring as a consequence of a reversal of time are manifested merely in the fact that what is found to be a repulsive non-gravitational interaction between two identical particles propagating in the same direction of time, would upon a reversal of the direction of propagation in time of one of the particles become an attractive interaction, or vice versa, as a result of the equivalent reversal of the sign of charge that occurs when a particle reverses its direction of propagation in time without actually reversing its charge. This is an unavoidable consequence of the fact that the departure of a positively-charged particle from a region of space would from a reversed-time viewpoint necessarily appear as the arrival of a particle of opposite (negative) charge, therefore implying that there is a relationship between the relative direction of propagation in time and the relative sign of any conserved physical quantity. We do not even need to know what an electric charge is or what energy is, from an exact mathematical viewpoint, to draw that conclusion. The reversal of charges associated with a reversal of time simply illustrates the subtlety of the relational definition of the sign of conserved (time-invariant) physical quantities in the context where there is a fundamental degree of freedom associated with time direction.

It must be remarked that in the context where there is, in effect, a dependence of the sign of charges on the direction of propagation in time, it follows that there no longer needs to be a mystery regarding why all charges come in two varieties, each having the exact same magnitude, but a polarity opposite that of the other. This is because, even if there were only, say, positive electrical charges, the fact that particles are free to propagate either forward or backward in time (under appropriate conditions) means that, from a practical viewpoint, there would still occur phenomena involving negatively-charged, but otherwise identical particles and it would not be possible to say whether it is the positive or the negative charges which constitute the ‘true’ charges. In such a context it seems possible that the requirement imposed by Grand Unified Theories that the sum of charges of all elementary particles cancel out, so that the overall symmetry is preserved in the context where it is not spontaneously broken, could ultimately be understood to be made possible (if the current elementary particles are actually composed of more fundamental building blocks) by the relativity of the sign of charges with respect to the direction of time, which not only allows, but actually requires the

existence of opposite charges. What I'm now suggesting is that we would in fact be justified to consider that the same requirement also applies to energy, which would therefore come in two varieties with opposite signs, not only for particles propagating in opposite directions of time, but even relative to the conventional, forward direction of time.

In any case, it should be clear that it is no longer possible to consider the sign of charges, including that of energy, independently from their direction of propagation in time. The traditional viewpoint according to which it seems possible to define charge without reference to some direction of time is valid merely because we implicitly always consider the sign of charge with respect to the positive direction of time (conventionally assumed to be the future). The positive-definite value of energy under all circumstances is thus an artifact of this implicit choice of the positive direction of time as the direction relative to which energy is measured. It is true, though, that if it was not for the non-gravitational charges carried by a particle it would, in effect, be impossible to distinguish between the case of a positive energy propagating forward in time and that of a negative energy propagating backward in time, just as it would be impossible to distinguish between the case of a negative energy propagating forward in time and that of a positive energy propagating backward in time. But there is no reason to assume that there would be no distinction between positive and negative energies propagating in the same direction of time and therefore the truly significant measure concerning energy is the sign of action, which is obtained by multiplying the sign of energy by the sign of time intervals. If the hypothesis that energy must necessarily be positive has always appeared valid it is merely as a consequence of the fact that we always measure energy relative to the positive or forward direction of time and for all known particles action remains positive. As I suggested above, however, this does not mean that energy really is always positive, but merely that action, or the sign of energy relative to the sign of time intervals, is, in effect, always positive for all currently known particles, independently from the true sign of energy of those particles.

What must be understood is that, ultimately, it is not only the sign of energy that is to be viewed as a relative quantity, but that the sign of action itself is purely relative, in the sense that there could never exist a generally agreed, absolutely defined, positive or negative value for the sign of action of a particle. In this context not only would the sign of energy be dependent on the direction of time in which a particle is assumed to propagate, but the sign of action would itself depend on the choice of what direction in time is to be

that in which what are assumed to be positive-energy particles propagate, or what is the sign of energy of those particles which are considered to propagate forward in time. Here all that matters is that once you define one particle as having positive action, because you assume that it is this particle that propagates positive energy forward in time, then the particles that you must assume to be carrying negative energies forward in time or positive energies backward in time, as a consequence of this *choice*, are those which will have negative action. But it must be clear that you are always free to describe the first particle as propagating negative energy forward in time and therefore as having negative action, as all by itself this choice is arbitrary, but in this case the other particles would then necessarily have to be assumed to carry positive action instead of negative action, because their *relationships* of time directionality and energy sign with the first particle (the difference or the identity of the signs of time intervals and energy) would remain unchanged.

It must also be remarked that the fact that what we would currently define as negative-action particles are related to ordinary matter through a simple convention regarding the direction of propagation in time means that the motive for rejecting the possibility that negative-action matter may actually exist is no stronger than that which would consist in arguing that ordinary matter itself is not allowed to exist. There is absolutely no rational motive for rejecting the viewpoint described here and many reasons to recognize its validity. In any case, the fact that the sign of action is a purely relative concept, which can vary as a consequence of assumptions regarding the direction of propagation in time, means that if the direction of the gravitational acceleration produced by a local matter distribution depends on the sign of action of its source, then it should also vary as a function of the assumptions made concerning the direction of propagation in time of the objects submitted to it (which determine their own action signs in relation to that of the source) and therefore the gravitational field must itself be considered a relative concept dependent on the conventions used by an observer.

Regarding the relation between the sign of charges in general and the direction of propagation in time it must be noted that energy actually distinguishes itself from non-gravitational charges by the fact that it is naturally reversed when a particle reverses its direction of propagation in time. Indeed, in the context where a particle-antiparticle annihilation process must be considered as an event during which a particle bifurcates in time to be-

gin propagating the same non-gravitational charges backward in time (which would effect the same kind of change as reversing the charges and keeping the direction of propagation in time unchanged), it must be assumed that the energy of the particle is reversed, along with the direction of time intervals, when the bifurcation occurs, given that the particle now propagates backward in time while its energy remains positive from the conventional forward-in-time viewpoint. In fact, we have no choice but to consider that only non-gravitational charges are left unchanged (relative to the true direction of propagation in time) when the particle begins propagating backward in time during what appears to be a particle-antiparticle annihilation process, because energy is always released by such a process and if the sign of energy had remained unchanged along with that of non-gravitational charges when the direction of propagation in time of the particle reversed, then an antiparticle's energy would be opposite that of its particle counterpart with respect to the forward direction of time and therefore the annihilation of such a pair could occur without any energy at all being released, as I previously mentioned. Thus, energy must actually reverse along the 'true' direction of propagation in time of a particle, when the particle reverses its direction of propagation in time during a pair-annihilation process, just like momentum naturally reverses when a particle changes its direction of motion in space. The negative energy of an antiparticle simply propagates backward in time so that relative to the positive, or forward direction of time it is left unchanged and from a mathematical viewpoint this interpretation fully agrees with the traditional description.

If this relational interpretation of the energy signs of particles involved in pair-annihilation processes is valid, then, based on the fact that we also have many reasons to believe that the gravitational properties of antiparticles are the same as those of particles, I can deduce that, from a gravitational viewpoint, the sign of energy is physically significant merely in relation to the direction in which a particle with that sign of energy is propagating in time. In other words, to produce an anomalous gravitational field, or to respond anomalously to a gravitational field, a particle would have to propagate its negative energy forward in time rather than backward, as does an ordinary antiparticle. This is a simple, but very significant result whose consequences will be developed in the following sections. What must be understood is the fundamental character of the degree of freedom associated with time direction, which, in a general-relativistic context, simply embodies the sum of all relationships of time directionality between a given particle

and all the other particles in the universe. This physical property must be considered distinct from any property of time directionality which is merely statistically significant and which is associated with the flow of information, like that which characterizes the irreversible processes obeying the second law of thermodynamics.

Concerning the gravitational properties of antimatter, it appears that it is actually unnecessary to appeal to any independent constraint, like the equivalence principle (which seems to require all matter to have the same acceleration in a gravitational field), to justify that antimatter should not ‘fall’ up in the gravitational field of a positive-energy planet like the Earth, as was often proposed before experiments began to rule out such a possibility. Indeed, any of the arguments traditionally provided to rule out the possibility of an anomalous gravitational behavior of antimatter become unnecessary once it is understood that it is actually only matter propagating its negative energy forward in time that could experience gravitation distinctively from normal matter, while it is already known that if negative energy is to be associated with antiparticles then this energy would in fact propagate backward in time. There is, thus, a very good reason to assume that antimatter falls down in the gravitational field of the Earth, but this is not an argument that we could use to rule out the possibility that some matter that would not be antimatter could perhaps be subject to anomalous gravitational interaction with ordinary matter, because there is no *a priori* motive for assuming that there cannot exist particles propagating negative energy forward in time. In fact, I will later explain that even the argument against anomalously gravitating matter which arises from the necessary application of the equivalence principle is not really unavoidable, because it is possible to better define this principle in a way that allows for the existence of anomalously gravitating matter of the appropriate type, while retaining the general form of the mathematical framework of relativity theory which can accommodate such a generalization.

In any case, it must be recognized that all those properties of fundamental time directionality discussed above are a reflection of the fact that the sign of charges (including energy) is not only defined in relation to the direction of propagation in time of the particle carrying those charges, but is actually determined completely arbitrarily as being merely significant in relation to the similar physical properties of other particles. From a relational viewpoint it would be incorrect to assume that the direction of propagation in time

of a given type of particle, carrying a unit of electric charge with a given, arbitrarily assigned positive or negative sign, is definitely the future direction, say, while the direction of propagation of the antiparticle of the same type is definitely the past, or even that there exists an absolutely defined character of being an antiparticle by opposition to being a particle. The only physical property that can be objectively defined without referring to quantitative attributes of objects that are not part of our universe is the relative direction of propagation in time of two particles. Two particles with the same type of charge may be both propagating in the same direction of time or they may be propagating in opposite directions of time and this is all we can ascertain through physical means.

What must be understood is that, while the relationship between the direction of propagation in time and the sign of a given charge, including energy, is a matter of coordinative definition (a definition that must be applied similarly to all processes in the whole universe on the basis of their relationships to one particular process for which an arbitrary choice of properties is assumed), once such a definition is applied, the difference between the sign of time intervals and the sign of charges is an objective physical property that is not dependent on a particular viewpoint. But it is not just the relationship between the sign of charge and the direction of propagation in time of a particle which can be given clear meaning through the use of a coordinative definition, because once we define one kind of particle as actually propagating a positive charge forward in time, then it should also be possible to differentiate such a particle from an otherwise identical particle propagating a negative charge in the opposite direction of time.

It must be clear, therefore, that once we assume an ordinary electron to be propagating its negative charge forward in time, it is not possible to consider another *ordinary* electron as perhaps propagating backward in time while carrying a positive electric charge in this direction of time (so that the electron would still appear to be propagating a negative charge relative to the forward direction of time). Indeed, if a certain condition of continuity of the flow of time on which I will elaborate in section 4.3 is assumed to apply, such a backward-in-time-propagating ordinary electron could only annihilate with an anti-electron which would be propagating the same positive charge forward in time (instead of propagating a negative charge backward in time). But this would actually mean that certain positrons cannot annihilate with certain electrons, while no constraint of this kind is observed to apply, as all known electrons have the same unique probability of annihilating with any

positron. Thus, if a constraint of continuity of the flow of time does indeed apply along an elementary particle's world-line, then an ordinary electron must be assumed to propagate in one and only one direction of time, while its antimatter counterpart must similarly be assumed to always be propagating in the opposite direction of time. Perhaps that this restriction is a consequence of the fact that there actually exists only one electron or that all electrons are 'the same particle' propagating forward and backward in spacetime, as John Wheeler once argued, but the condition of continuity of the flow of time does not specifically require the validity of this hypothesis.

On the basis of those considerations and given the previously reached conclusion that only the sign of energy with respect to a given direction of time has physical significance, it must, in effect, be recognized that only a particle propagating either negative energy forward in time or positive energy backward in time (in the context where ordinary matter is considered to propagate positive energy forward in time) could potentially respond in an anomalous way to the gravitational interaction. What is important to know about such a particle, which we may call a negative-action particle² to distinguish it from a particle merely propagating negative energy backward in time like an antiparticle, is that the preceding considerations, regarding the relational definition of physical quantities, would also mean that the particle cannot possibly be considered to have physical properties that would qualify it as responding to the gravitational field of a positive-action body in an anomalous fashion that would not also be shared by an ordinary matter particle (propagating positive energy forward in time) submitted to the gravitational field of a negative-action body. This must be considered an unavoidable conclusion in the context where one can physically distinguish only a difference or an equality in the signs of action of any two particles and cannot attribute objective meaning to the sign of action itself. That does not mean that there would actually be no anomalous response, only that, in a configuration where all 'anomalously' gravitating matter is replaced by ordinary matter and all ordinary matter is replaced by anomalously gravitating matter, we should observe no difference (attributable merely to the gravitational interaction). Thus, a particle defined as having negative energy relative to the positive direction of time and which would be located in the

²Despite the ambiguity, I still use the term 'negative energy' in place of 'negative action' to identify such anomalously gravitating matter when the context clearly indicates that I mean negative energy propagating forward in time, or equivalently positive energy propagating backward in time.

gravitational field of a planet having opposite energy relative to the positive direction of time should behave in the same way as a positive-energy particle in the gravitational field of a negative-energy planet and similarly for any combination of energy signs of particle and planet, because only the relative difference in forward-propagated energy signs can be considered significant. Given the preceding discussion, this should be crystal-clear. But that is not what is usually assumed to occur by people discussing negative energy or making quantitative predictions involving matter in such an energy state.

What is usually assumed is that a positive-energy or positive-mass body would attract all bodies, regardless of whether those bodies have positive or negative energy or mass, while a negative-mass body would repel all bodies, again regardless of whether those bodies have positive or negative mass. It is currently believed that this is the consequence of taking inertial mass to be reversed along with gravitational mass, as would appear to be required by the equivalence principle. It must be clear, however, that those are not results which are ‘derived’ from relativity theory, as is sometimes suggested, but merely the consequence of a choice that is implicitly made regarding what properties should be associated with negative inertial mass, while trying to be as accommodating as possible with the traditional conception of the principle of equivalence. But if I find it appropriate and indeed necessary to consider, as most people do, that inertial mass is reversed along with gravitational mass when we are considering an object with negative energy (which would normally allow the equivalence principle to apply), I cannot agree with the conclusion that is usually drawn from such an assumption. Indeed, for the response of various masses to the presence of a negative mass to be in line with common expectations, it must be possible to determine the sign of mass, or the sign of action of particles in an absolute, non-relational manner, because we are assigning the attractive or repulsive character of the gravitational field in precisely such an absolute manner (the field is either repulsive for everything or attractive for everything), which I believe could never be justified.

I think that it cannot be assumed that a negative mass is repulsive in an absolute, invariant way, because it would not be possible to tell relative to what reference point the distinctiveness of this character is defined, given that positive mass cannot be used as a reference if its gravitationally attractive nature is itself absolutely defined (does not vary merely in relation to a variation of the sign of mass of the object with which it is interacting). I will explain, in a later section of this chapter, why it is that the assumption

that a negative inertial mass is associated with a reversal of the sign of action, far from having the undesirable consequence of allowing absolutely defined physical properties into physical theory (if there could ever be such a theory), actually gives rise to a description of the gravitational interaction between positive and negative-mass bodies that is in perfect agreement with the requirement of relational definition of the sign of mass or energy (once the inertial properties of negative-mass matter are well understood). All that would then remain to understand is how the equivalence principle can still be satisfied by such a description. For that purpose, I will provide arguments to the effect that a simple reconsideration of the true significance of the principle of equivalence, and a better understanding of its motivation in the principle of relativity of accelerated motion, allows its foundations to be preserved while enabling the more consistent, relational viewpoint on the sign of mass to be retained and to actually be integrated into the core mathematical framework of relativity theory by introducing a slight modification to this classical theory of gravitation that is actually a simple generalization of it. In order to further justify this approach, I will first try to identify what should be the true properties of negative-action matter and why we should not expect such matter to behave in ways that would make it undesirable, not only from the viewpoint of the requirement of a relational description of physical quantities, but with respect to other constraints and other physical principles which we can be confident must also be obeyed.

2.3 Our current understanding

Before addressing the question of how a negative-energy particle would actually behave, we may first want to explore what the current situation is regarding the notion, or indeed the problem of negative energy. For this purpose, it should first of all be noted that for many reasons no one seems to like the idea that there could exist negative-energy particles. Thus, it is no surprise that one of the most basic and often implicit assumption that enters our description of physical reality is that energy must always be positive. There are many different mathematical formulations of that requirement which impose various degrees of conformity to the hypothesis that matter cannot find itself in a negative energy state (for a technical review of those conditions see Ref. [3]). In its least restrictive form this condition is called the weak energy condition and merely constitutes a statement about the positivity of

the components of the stress-energy tensor (the most general representation of the energy content of matter). More constraining conditions have also been proposed, among which is the appropriately named strong energy condition which, if obeyed under all circumstances, would mean that gravity must always be attractive (between all forms of matter which would then be allowed to exist). Those conditions are used as rigorously defined hypotheses in various theorems dealing with the behavior of matter under the influence of the gravitational interaction.

The problem is that it was found, at some point, that configurations involving negative energy densities are actually allowed to occur in quantum field theory [4]. This does not mean that negative-energy particles are explicitly allowed by current theories, but merely that unlike what we would expect from a classical viewpoint, where the vacuum is described as a total absence of matter, quantum field theory allows for the local density of energy to not always be positive definite, even in the context where only positive-energy matter is present. A well-known experiment illustrates the kind of phenomena involved. It requires placing two parallel mirrors a very small distance apart in a vacuum, so as to forbid some states, which would normally exist in the vacuum, from being present in the space between the mirrors, as a consequence of the incompatibility of their characteristic wavelengths with the spatial constraints imposed by the presence of the mirrors. The predicted result, which is actually observed, is that there should arise a small pressure pulling the mirrors together as a consequence of the comparatively larger pressure exerted from the outside, which is actually caused by a decrease in pressure from between the mirrors that can be attributed to the restriction imposed on which virtual particles can be present in this volume. This is of course the phenomenon known as the Casimir effect [5]. It is clear though that we are not directly measuring a negative energy density in such an experiment, but merely the indirect effects of an absence of some positive contribution to vacuum energy, which is then assumed to imply that the energy density in the small volume between the mirrors is smaller than that which exists even in the absence of any matter and which would traditionally be considered null. But even this particular occurrence of negative energy is assumed to be so serious a problem by some theorists that they suggested that the description of the vacuum as involving virtual particles coming in and out of existence is actually only a mathematical trick and does not reflect what is really going on in the absence of ‘real’ matter.

However, this aversion for whatever is negative of energy is not shared

by all authors and some more open-minded specialists have tried to address the issue of negative energies as they occur in quantum field theory and in so doing gained some significant insights into what exactly is allowed by a quantized description of the vacuum. A modified version of the weak energy condition was thus proposed that allows to take into account the fluctuations of energy which arise in the quantum realm. This condition, which is appropriately called the averaged, weak energy condition, involves only quantum expectation values of the stress-energy tensor averaged over some period of time during which the observations are assumed to occur, rather than idealized measurements at a spacetime point. A feature of the constraint provided by this condition is that it allows for the presence of large negative energies over relatively large regions of space if there is a compensation by the presence of a larger amount of positive energy during the time period over which the observations are made. It was indeed found out [6] [7] [8] [9] that quantum field theory places strong limits on the values of negative energy density that can be observed over finite periods of time under various conditions. What emerges from those developments is that there appears to be a constraint on the magnitude of negative energy that can be observed and it indicates that negative energy can be merely as large as the time interval during which it is measured is short. I believe that this is indicative of the fact that while negative energy states cannot be ruled out as strictly forbidden, they should also clearly not be expected to materialize in stable form in the context where we are dealing with ordinary matter configurations, for which the particles are already predominantly in positive energy states.

A similar limitation can also be observed to restrain another form of negative energy that occurs in the presence of an attractive force field, even in a classical context. Indeed, the energy contained in the force field between two particles submitted to an attractive interaction must be considered negative. This is because work and positive energy must be provided to separate two particles attracted to one another in such a way and given that it must be assumed that the attractive field responsible for this interaction would contain no energy at all when the particles are separated by a distance that tends to infinity (in the context where the strength of the field associated with a long-range interaction decreases in proportion with the square of the distance, so that it must be null when this distance is infinite), then we must conclude that the energy initially contained in the same attractive force field when the particles were near one another was actually negative (so that adding pos-

itive energy can produce a null final value). This conclusion is undeniable, given that it is actually observed that the energy of a bound system, formed of many interacting particles, is lower than the sum of the energies of those particles when they are free.

Thus, the energy contained in an attractive force field must definitely be considered negative, as this energy is required to provide the negative contribution that reduces the energy of the whole bound system. The additional energy that was present before the formation of a bound system is in fact released (through the emission of radiation for example) when the system is created, but except for the additional negative energy contained in the attractive force field, the system is identical, in terms of its matter particle content, to what it was initially and therefore we definitely need the negative energy. This is made more obvious when we consider larger systems like those bound by the gravitational interaction. It was shown, in effect, that even a system as large as the Earth-Moon system has an asymptotically-defined total mass (providing a measure of its total energy) which is smaller than that of its constituent planets (when it is possible to neglect any contribution which would normally be attributed to the presence of dark matter) and observations confirm this prediction. Therefore, it is clear that the energy contained in the gravitational field maintaining the two planets together must be negative.

What is crucial to understand regarding the situation described above, however, is that even if we must acknowledge the existence of a well-defined negative contribution to the energy of some physical systems that diminishes their total energy, it is again impossible to measure that energy directly and it can merely be deduced to occur from the behavior of the positive-energy subsystems which are submitted to the attractive interaction. Here also, the negative energy must be associated with virtual particles, namely the unobserved bosons that mediate the interaction, and cannot be measured independently from the total energy of the bound systems, which usually remains positive. It is simply not possible to isolate the attractive force field of a bound system from its positive-energy sources and this is true for systems of any size. It would, nevertheless, certainly be a concern if the negative binding energy of a system made of positive-energy components could become so negative as to make the total energy of the bound system itself negative. Once again, however, it was shown that there are unavoidable theoretical constraints on the values that observable total energy can take. It was shown [10] [11] [12] [13] [14] [15] [16] [17] [18], concerning the gravita-

tional interaction in particular, that the energy of matter (everything except gravitation) plus that of gravitation is always positive when the dominant energy condition is assumed to be valid, which actually amounts to assume that the energy of the component particles is itself positive. If we compress a positive-energy body too tightly, it simply collapses into a black hole before its surface area is allowed to become so small and its energy density so large that the magnitude of its negative gravitational potential energy would be larger than the positive energy of the matter. Thus, positive-energy matter cannot turn into negative-energy matter through an increase of negative gravitational potential energy.

What must be retained from the preceding considerations, therefore, is that even though it is often present, negative energy seems to never be measurable. But this conclusion is valid merely under the condition that we are dealing with situations where matter was already in a positive-energy configuration to begin with. It must be clear, however, that we still have no argument to rule out the possibility that there may exist configurations where the component particles themselves would have negative energies and for which there would exist constraints, similar to those unveiled here, enforcing the *negativity* of energy.

In a previous section of this chapter I mentioned that it is desirable from a certain viewpoint to consider antiparticles as propagating negative energy backward in time. Indeed, if antiparticles are propagating backward in time, as the reversal of their non-gravitational charges clearly suggests, then they *must* have negative energy relative to the direction of time in which they are propagating (which is the past), so that relative to the opposite direction of time (which is the future) they would still appear to have positive energy, as required. In fact, it was discovered a long time ago by Paul Dirac (when he achieved his unification of special relativity and quantum theory) that there is a mathematical requirement for the existence of negative energy states. Indeed, it turned out that in order to obtain Lorentz-invariant equations for the wave function one had to sacrifice the positivity of energy. After having considered various possible interpretations for what in nature could possibly correspond to those negative energy states Dirac concluded that it required the existence of a new category of particles, the antimatter particles, which would consist of holes in a filled distribution of such negative-energy matter. But despite the fact that it was later found that antiparticles do exist, as he predicted, Dirac's solution to the problem of negative energy states was

never considered fully satisfactory.

Antiparticles were eventually described by Feynman (following Ernst Stückelberg) as particles propagating negative energies backward in time, which allowed to fulfill the mathematical requirements imposed by the existence of the negative energy states (by providing an interpretation for those transitions which were predicted to involve a reversal of energy) without requiring the presence of the filled, negative energy continuum. But in the process, it seems that the discovery that particles could actually occupy negative energy states, which appeared to be implied by the original developments, was somehow forgotten and lost in the details of the proposed solution. This indifference was probably justified by the fact that antiparticles could still be considered to have positive energy, for all practical purpose. But what is usually unrecognized is that while attributing a positive energy to antiparticles may appear more ‘reasonable’ than assuming that those particles propagate negative energy backward in time, such a choice would actually imply that it is the particles themselves (by opposition to antiparticles) which must then be considered to carry negative energy backward in time, because it must be either that or the opposite. This is what the subtleties of the quantum-mechanical definition of energy seems to require that was not apparent classically.

The reluctance to recognize the true physical significance of negative energy states is probably also in part a consequence of the apparently insurmountable difficulties which would be associated with the possibility for particles to occupy those physically allowed states. First of all, it is certainly not desirable from a theoretical viewpoint to assume that antiparticles would be submitted to anomalous gravitational interaction as a consequence of propagating negative energy backward in time, because it was demonstrated some time ago [19] that if, for any reason, antimatter was to be found experiencing repulsive gravitational interactions with ordinary matter, we would run into a number of problems ranging from violations of the conservation of energy and up to the undesirable and unlikely (from a theoretical perspective) possibility of producing perpetual motion machines. But an analysis of the arguments presented against the possibility of anomalously gravitating antimatter has led me to conclude (for reasons which will be explained later) that the problem really has to do merely with the possibility for antimatter ‘as we know it’ to experience gravitational repulsion. It cannot be considered to mean that matter in a true, negative energy state (propagating negative energy relative to future-directed time intervals) could not exist and experience

anomalous gravitational interactions with ordinary matter without violating the principle of conservation of energy or the second law of thermodynamics, because matter in such a negative energy state may also, by necessity, have properties different from those which are known to characterize antimatter, in particular with what regards non-gravitational interactions.

Nevertheless, most people today seem to consider that the developments that followed the introduction of the early theory of relativistic quantum mechanics and which gave rise to modern quantum field theory have eliminated the problem of negative energy states, which can now be considered a mere artifact of the former single-particle theory. Thus, the predicted negative energy states would simply be nonphysical solutions that must be discarded as irrelevant to physical reality. But it must be clear that this is indeed what we are doing here. We are rejecting the possibility that a particle could be found in a whole set of states that are allowed by the most basic equations without providing any justification as to why those states should be forbidden. Indeed, upon closer examination it becomes clear that if ‘true’, negative energy states do not explicitly arise in quantum field theory it is not because the structure of the theory forbids them, but simply because we *choose* to ignore those solutions to start with and then integrate that choice into the formalism. More specifically, it turns out that what prevents negative-action particles from showing up in quantum field theory is merely a choice of boundary conditions for the path integrals that provide the probability amplitude for transitions involving particle trajectories in spacetime. There are several possible choices for expanding those integrals which all constitute valid solutions of the equations of the theory, but only those solutions propagating positive frequencies forward in time and negative frequencies backward in time are usually considered to be physically significant, while the solutions propagating negative frequencies forward in time and positive frequencies backward in time, which are also valid from a mathematical viewpoint, are systematically rejected. But this actually amounts to retain only the positive-action portion of the theory, while ignoring all transitions involving negative-action (although not negative-energy) particles. There is no other origin for the often-mentioned conclusion that quantum field theory does not involve negative-energy matter. It is our own arbitrary decision to reject all transitions involving negative-action particles.

In order to make the choice of boundary conditions responsible for the absence of negative-action particles in quantum field theory more acceptable it is sometimes suggested that the negative energies predicted by the single-

particle, relativistic equations are simply transition energies, or differences between two positive energy states and there is obviously no reason why those variations could not be negative if they can be positive. But no explanation has ever been provided for why the same reasoning could not be applied to the energy states themselves, which are also energy differences, given that the energy of a particle is always defined in relation to the zero level of energy associated with the vacuum in which it propagates. There is no justification for this arbitrary distinction between transition energies and particle energies, except for the satisfaction that is obtained by the physicist in having easily disposed of an embarrassing problem. It may of course be argued that there is nothing wrong with those methods, given that they appear to be validated by experimental results. Indeed, we have never observed interferences by negative-action particles into the outcome of any experiment conducted at any level of energy and to any degree of precision. But I would like to emphasize that this still doesn't constitute an explanation for the absence of negative-action particles.

Thus, the problem I have with the modern approach to quantum field theory is that the formalism is generally introduced in a way that encourages us to believe that, after all, no particle is actually propagating backward in time with negative energy and that a positron is really just another particle, identical to the electron, but with an opposite electrical charge. However, this viewpoint does not only complicate things unnecessarily as a consequence of rejecting the possibility for electrons and positrons and all other particles and their related antiparticles to actually consist in the same particles observed from different perspectives, it is also completely ignorant of the requirement of a relational definition of any physical attribute dependent on the fundamental time-direction degree of freedom. But if we choose to recognize the validity and the greater value of the viewpoint defended here and according to which antiparticles are really just ordinary particles propagating backward in time, then we must accept that there definitely exist in nature particles which are known as carrying negative energies and if the arguments provided above concerning the arbitrariness of the current restrictions imposed on the propagation of those negative energy states are valid, then we would have to conclude that there should necessarily also exist particles propagating such energies forward in time and which could be submitted to anomalous gravitational interactions in the presence of ordinary matter.

2.4 The negative-mass concept

When discussing the issue of negative mass, what must first of all be understood is that, if the physical property of mass is to have any polarity associated with it, such that we could attribute to mass either a positive or a negative sign, then this polarity must be directly related to the sign of action, that is, to the sign of energy relative to the positive direction of time. This is because, as I previously emphasized, the sign of action is the only physical property from which the attractive or repulsive character of the gravitational interaction between two bodies could depend. We may, thus, attribute positive mass to a positive-action particle and negative mass to a negative-action particle. Mass being a Newtonian concept, its polarity must be determined in relation to a particular Newtonian gravitational field. From this viewpoint the sign of mass of a given particle could, in effect, be understood as determining the response to the gravitational field of a given source, in the sense that it would determine the *direction* of the gravitational force exerted on such a particle. If we may consider the gravitational field of the source (represented by a vector in Newtonian mechanics) to be uniform, then only its own direction or polarity (which we may assume to be dependent merely on the sign of mass of the source, when its position is assumed to be fixed) would be decisive in determining the kind of response experienced by a given type of mass submitted to it. Equipped with such a definition, we can meaningfully discuss the problem of the gravitational interaction of negative-action particles with positive-action particles and with themselves as the problem of the gravitational interaction of positive and negative masses. This will allow us to better grasp the significance of the assumptions that will form the basis of the new interpretation of negative-energy matter which I shall propose and therefore, also, to gain better confidence in their validity, even in the more appropriate context of a general-relativistic theory.

If we may agree on those requirements, then I think that what must emerge is that, if it is indeed important to have a well-defined concept of negative mass, then it also seems that such a negative mass must be negative in all respects. That there could be a difference between the sign of gravitational mass and the sign of inertial mass is usually considered to be forbidden merely by the general theory of relativity which is, in effect, founded on the principle of equivalence which requires the equality of gravitational and inertial masses. However, I think that if this hypothesis is justified, it is not because our negative-mass concept must comply with some per-

ceived requirement from general relativity theory, but because it would not be acceptable to attribute mutually exclusive values to one single physical property. Thus, I do believe that the mass of any particle or bound object should be either definitely positive, or definitely negative (but still in a relational way), regardless of whether we are considering gravitational mass or inertial mass, if the concept itself is to have any consistent physical meaning. But unlike most theorists, I do not consider that this requirement must be assumed to imply the kind of behavior that is usually attributed to negative-mass matter, where gravitational repulsion is an intrinsic property of this type of matter itself, independently from the sign of mass of the matter with which it is interacting. This is indeed the conclusion I was able to draw based on the outcome of the previously discussed analysis of the constraints imposed by a relational definition of the sign of energy, for reasons I will now explain.

The difficulty I originally met when I first began to explore the possibility that inertial mass could be reversed along with gravitational mass when we are dealing with negative-mass matter is that, if both the gravitational mass and the inertial mass are to be negative at once, then it seems that there could occur situations where the principle of inertia would be violated (I will explain what motivates this conclusion below). I was able to understand, however, that this is merely a consequence of the inappropriateness of current assumptions regarding what we should expect to be the behavior of matter with both a negative gravitational mass and a negative inertial mass. Actually, despite the fact that it is usually taken for granted that we know for sure at least what the behavior of matter with positive mass is, because we routinely observe gravitational phenomena involving this kind of matter and there can be no mistake here, I will explain that this is not entirely the case and that there is still much confusion as to even what we should expect concerning the response of positive-mass matter to a concentration of negative mass. Currently, it is assumed that given that positive-mass matter gravitationally attracts all matter and resist the action of any force exerted on it, then this must be an intrinsic property of such positive masses. On the other hand, it is usually assumed that two choices exist for what could possibly characterize the behavior of matter with a negative mass. The situation we have right now is thus the following.

First of all, we must assume that *gravitational* mass is indeed negative when mass is reversed. This would give rise to gravitational repulsion when only the mass of the source (the active gravitational mass) is negative, be-

cause it reverses the polarity of the Newtonian gravitational field to which any passive gravitational mass is submitted and therefore should at least reverse the force exerted on positive-mass bodies. But once this is recognized it is usually considered that two possibilities actually exist for a negative-mass particle submitted to a given gravitational field, depending on whether *inertial* mass is assumed to remain positive or is itself also negative. Here, the inertial mass of a particle is assumed to determine the response of that particle (actually the direction of its acceleration) to any force, including a gravitational force, while the gravitational mass of the same particle is assumed to determine both the polarity of the gravitational field it produces and the response of the particle to a given gravitational force field. If we were to agree with those assumptions, then we would have to conclude that a negative-gravitational-mass particle with a negative inertial mass, should actually respond normally to any gravitational force field (because the nature of its response is changed twice, once by the reversal of its inertial mass and once by that of its gravitational mass) while its response to non-gravitational forces would be reversed (same force, opposite acceleration), as current assumptions concerning the effects of a reversal of inertial mass would imply. But we must also keep in mind that the fact that this kind of matter would respond normally to gravitational force fields would, under current assumptions, still mean that it is repelled by matter of the same type, because the gravitational field produced by such matter is also assumed to be reversed. Thus, such negative masses would repel masses of all signs, be repelled by other negative masses and be attracted to positive masses, still under the hypothesis that the above stated commonly accepted assumptions are valid. Given that it is usually considered that, in a general-relativistic context, all mass (gravitational and inertial) must be negative, this is the choice that is usually retained as defining the behavior of negative-mass matter if it could exist.

But despite the support that is usually granted to such a conception of negative-mass or negative-energy matter, I think that enormous problems would arise if it was retained as a valid proposal. Some of those problems, involving black holes and the second law of thermodynamics, will be discussed later, but even if we remain at the level of classical Newtonian dynamics we can readily identify one very serious problem, which is that the existence of such matter would allow violations of the principle of inertia (considered as a generalization of Newton's first law) or the very idea that no physical system can accelerate without work being done on it by an external force. This is

because indeed, as stated above, from the current viewpoint, a negative-mass body would both repel positive-mass bodies and be attracted to them. Such a combination of features could then give rise to unlikely phenomena like pairs of opposite-mass bodies chasing one another and in the process accelerating to arbitrarily large velocities, still without any external energy input, as described in Ref. [20]. The fact that energy would in principle be conserved under such conditions (because the energy gained by one of the bodies would be opposite that of the other) is no consolation as it actually makes the problem worse, given that it would allow both positive and negative kinetic energies to be produced out of nothing, while consistency requirements require kinetic energy to be conserved as a positive-definite quantity from the viewpoint of positive-energy observers (or as a negative-definite quantity from the viewpoint of negative-energy observers), as I will explain in sections 2.11 and 2.13).

In fact, an even more basic violation would occur if such phenomena were made possible by the existence of negative-energy matter. The problem I see is that there would be no equal and opposite force to that applied on a given body that could be attributable to its assumed interaction with the other body and this would be a violation of the principle of action and reaction (Newton's third law, which is a particular aspect of the principle of local causality), while this is one requirement that in all fairness we should recognize as being as essential as that of conservation of energy, because if it does not rigorously apply then absolutely anything could occur and under such conditions we could not give much of even the principle of conservation of energy. However, I think that what those observations show is not the nonphysical nature of negative mass, but merely the ineffectiveness of the traditional approach to describe the behavior of this kind of matter. It is important to mention, by the way, that even though this hypothetical situation of accelerating opposite-mass pairs has been described by other authors in the past, none of them has ever recognized that what it actually demonstrates is the inconsistency of the currently favored notion of negative mass, which I believe is illustrative of the state of denial in which most people remain concerning the possibility that there could actually exist negative-mass matter.

What is also significant concerning the unlikely phenomenon described above is that it would necessarily be the positive-mass bodies that would be chased in this way, while the negative-mass bodies would inevitably be those trailing them. But isn't it strange indeed that there should be such a clear

and decisive distinction between what constitutes the role of positive masses and what constitutes that of negative masses? Doesn't it seem like there is something wrong with such a hypothetical phenomenon? Shouldn't we only be allowed to define the property of gravitational attraction and repulsion in such a way that we could not observe such mass-sign-distinguishing behavior? What I have understood is that the unease we may experience in face of the strangeness of such phenomena is in fact justified. Indeed, it does not just seem like there is something wrong here, because what we have just described is actually the perfect example of an attempt to distinguish a physical property (the positivity of mass or the attractiveness of gravitation) despite the absence of any reference in the physical universe to which that arbitrary distinction could be related, which violates the very basic requirement of relational determination of physical attributes discussed above. The mistake which is made by assuming the validity of the traditional viewpoint is that we suppose that we can define attraction and repulsion in an absolute (non-relative) manner such that one kind of mass always attracts all kinds of masses regardless of their polarities and another always repels all masses, still regardless of their polarities, as if attractiveness and repulsiveness were intrinsic aspects of one and the other type of mass.

However, if the sign of mass is to be considered a meaningful physical property of elementary particles, then it must be taken to indicate that there can be a reversed or opposite value to a given mass and this reversed value can be considered to be reversed merely in relation to a non-reversed mass and to nothing else. A mass cannot be considered to be reversed with respect to an absolute point of reference lacking any counterpart in the physical universe. Therefore, if a gravitational field is to be assumed repulsive as a consequence of the reversed (negative) sign of the mass of the matter that is the source of the field, then this gravitational field should be repulsive only for an unchanged (positive) mass particle and not with respect to other negative masses. It would be incorrect to assume that the attractive or repulsive nature of gravitational fields depends solely on the sign of mass of the source itself, because no distinction exists for the sign of a mass other than its sameness or oppositeness compared to that of another mass. That does not mean that the field itself must be assumed to change as a consequence of the reversal of the sign of mass of the particle experiencing it (even though that may be one way to describe things if other conventions are adopted for the sign of mass itself as we will see later), but merely that the response of a negative-mass particle to a given gravitational field must be reversed in

comparison to the response we would expect from a positive-mass particle submitted to the same field, despite the associated reversal of the inertial mass of such a particle. If that was not the case, then I think that we would have to conclude that negative mass is, in effect, forbidden.

If the incorrect hypothesis on which the traditional approach is based, regarding the effect of a reversal of inertial mass, nevertheless allows to successfully (from my viewpoint) predict that a positive mass would be repelled in the gravitational field of a negative mass, it is simply because we assume the right inertial properties for the positive-mass matter submitted to the gravitational force of the negative mass. Thus, the positive mass responds in the appropriate way to the gravitational force exerted by the negative mass, which is correctly assumed to be a repulsive force, given that the gravitational field produced by the negative mass is necessarily opposite that which would be produced by a positive mass of similar magnitude located in the same position. The problem is that, given that it seems that we cannot expect the same kind of behavior from a negative mass submitted to the gravitational field of a positive mass, then it would appear that the behavior of both positive and negative masses is the consequence of some predetermined property of absolute attractiveness and repulsiveness (that cannot be related to any property of the source defined with respect to a property of the matter with which it interacts) associated with the gravitational fields emanating from positive and negative masses respectively.

The difficulty to which the traditional interpretation gives rise is also made apparent when we consider the case of a negative mass in the gravitational field of another negative mass, given that now the negative mass would be repelled by the same negative-mass matter (because the gravitational force is unchanged, but the response to this force would be reversed), while, on the basis of the relational definition of mass sign, there should be no difference between this case and that of a positive mass in the gravitational field of another positive mass (which is symmetric to the other case under exchange of mass signs). The appropriate outcome could only be obtained if, in addition to the assumption regarding the nature of the gravitational force between two negative-mass bodies, it was also assumed that the reversal of the inertial mass of the negative-mass body submitted to this force actually changes nothing to the response of that body to the force that is exerted by the other negative-mass body. Thus, the problem of the absoluteness of the attractive or repulsive nature of the gravitational field arises as a direct consequence of current assumptions regarding the effect of a reversal of inertial

mass. It is only in this context that the direction of the Newtonian gravitational force associated with a concentration of matter of positive or negative mass sign acquires an absolute meaning and is not merely dependent on the identity or the difference between the sign of mass of the matter submitted to the gravitational field and that of the matter that is the source of this field.

Even if merely as a consequence of the previously discussed considerations regarding the relative nature of the sign of energy (as dependent on the direction of propagation in time of a particle) and the purely conventional (subject to an arbitrary coordinative definition) significance of the sign of action, it would appear that a consistent notion of negative mass would require that it is the relative difference or absence of difference between the mass signs of two gravitationally interacting bodies that determines the attractive or repulsive character of this interaction, so that two negative-mass bodies should be submitted to the same mutual gravitational attraction that is experienced by two positive-mass bodies, while the same negative-mass bodies would also repel ordinary positive-mass bodies and be repelled by them, unlike is usually assumed. But the fact that it is often not even fully understood that negative mass should, in effect, be associated with negative action is illustrative of the confusion that surrounds the whole question of negative energy and gravitational repulsion, because there should be no doubt that, if it is possible for the sign of mass of a given body to be negative in some way, then this would necessarily have to occur as a consequence of the fact that this body has negative energy, or more precisely negative action. In any case, if the traditional viewpoint allows predictions that violate the expectations of a relational definition of mass sign, it is precisely because it allows to assume that there can be an absolute character of attractiveness or repulsiveness associated with a given sign of mass. To be fair, I must acknowledge that some authors did suggest in the past that the gravitational interaction should perhaps be repulsive between two bodies with opposite mass signs, while it would be attractive between two negative-mass bodies (just as it is between two positive-mass bodies), but simply on the basis of the fact that the sign of the gravitational force that is obtained by reversing the sign of one of the masses in Newton's equation for universal gravitation would itself be reversed, while it would be unchanged if the signs of the two masses were together reversed.

But even though it is not necessarily wrong to suggest that the repul-

sive or attractive nature of the gravitational interaction is determined by the signs of mass in Newton's equation for universal gravitation, it is only when we realize that the sign of mass must be related to the sign of action that we can begin to understand why it is that there should be a symmetry under exchange of positive and negative masses. This is because, as I previously mentioned, positive action states are related to negative action states by a simple convention regarding the sign of energy and that of time intervals, so that the sign of action is itself a purely relative notion. There must consequently be a symmetry under exchange of positive and negative-action matter, which would then require the behavior of positive masses in relation to themselves and in relation to negative masses to be similar to that of negative masses in relation to themselves and in relation to positive masses. I may add that, in such a context, it appears that the suggestion that, if negative mass bodies have never been observed it is perhaps simply because they do not assemble themselves into larger masses (as a consequence of their assumed absolute, gravitationally-repulsive nature), cannot be valid and if negative-mass matter exists, then alternative arguments would have to be proposed to explain this absence of observational evidence. Later on in this and the following chapters I will explain how it is possible, in effect, to reconcile the apparent absence of concentrations of gravitationally-repulsive matter on stellar and galactic scales with a more consistent notion of negative-mass matter.

The contradictions of the traditional conception of negative-mass matter can be illustrated by using a rarely discussed thought experiment. It has, in effect, been proposed that the sign of energy of a negative-mass particle could be determined by measuring the energy lost or gained while raising or lowering the particle in the gravitational field of some large object. Now, according to the traditional conception, if we were to raise a negative-mass body in the gravitational field of a positive-mass object like a planet, we would have to produce work and exert a force directed downward, because the inertial mass of the body is negative, which according to the traditional viewpoint means that it responds perversely to the applied force. But then, it is also the case, according to this same viewpoint, that the gravitational force exerted by the planet on the body should be attractive, because the planet has positive mass. Thus, we would be in the situation where we would have to exert a force downward to raise a negative-mass body in the gravitational field of a planet that exerts an *attractive* force on that body. I do not know to what extent people actually believe in the validity of such a conclusion, but I think that, faced with such absurdities, one has to come

to realize that the contradictions involved are a clear indication that the traditional assumptions regarding the behavior of negative-mass or negative-action matter are incorrect and that a better interpretation of what such a state of matter may involve is required.

Despite the fact that the question of the validity of the traditional conception of negative-mass matter had never been clearly analyzed before, it is no doubt the general feeling that there is something wrong with the possibility of observing phenomena of the type described above (including that where pairs of opposite-mass bodies accelerate without any external force being applied on them) which is responsible for having transformed the idea of negative-energy or negative-mass matter into the synonym of nonsense it has become in the minds of so many researchers. But, is negative mass really to blame here, or could it be that we are not attributing to it the right physical properties? There is, of course, even under the conventional assumptions regarding the response of negative-mass particles to applied forces, another possibility, which is that when gravitational mass is negative, inertial mass may remain positive for some reason. Of course that would not only appear to contradict the equivalence principle, as is already understood, it would also, if I'm right, itself be nonsense, as we would have to assume that one single physical quantity related to one single particle (the mass of that particle) is at once both positive and negative, for the same observer. The latter problem has never been discussed, but I think that it is actually the strongest argument one can make against this second possibility. We may nevertheless begin by exploring the consequences of such a choice.

Under the same commonly held assumption to the effect that the response of a particle to any force is dependent on the sign of its inertial mass, we would have to conclude that a negative-gravitational-mass body, to which a positive inertial mass would be attributed, would respond anomalously (in comparison to the response expected of a positive mass) to any gravitational force field (because the nature of the response is changed only once by the reversal of its gravitational mass), while its response to non-gravitational forces would be unchanged (same force, same acceleration), because the inertial mass remains positive or unchanged in comparison with that of positive-mass bodies. Therefore, if material bodies were to exist that would be made of such negative-mass matter they should, from the traditional viewpoint, gravitationally attract one another (as do positive masses), repel positive-mass bodies and also be repelled by those same positive-mass

bodies. As a consequence, we would observe no violation of the principle of inertia in this case and also no acceleration without work. If this behavior was to be observed, it would in fact be possible to exchange all positive-mass bodies with negative-mass bodies and vice versa and no apparent change in the phenomenology of the gravitational interaction would be detectable, because gravitational repulsion would only occur when there is a difference in the signs of the *gravitational* masses which are interacting. Thus, from a purely phenomenological viewpoint there would be equivalence between positive and negative-mass bodies.

Given the previous discussion regarding the necessity of a relational determination of the sign of energy, which would here be a requirement for the relational determination of the sign of mass, this situation would appear more appropriate, because, indeed, it would be impossible in principle to differentiate any intrinsic property of gravitational attraction or repulsion and only the difference or the equality of the signs of *gravitational* mass of two particles would be physically significant. The problem that most people would have with this possibility, however, is that it would explicitly violate the equivalence principle, because positive and negative gravitational masses would respond differently to a given gravitational field, produced by a given matter distribution, even if they are located in the same local inertial reference system.

But I think that, even before we consider the issue of the apparent incompatibility with the principle of equivalence, we must first of all ask how could it be determined which of the two types of matter would indeed have the *inertial* mass opposite its gravitational mass? And then it is obvious that this question could never be settled (because we could never decide which type of matter actually has a negative gravitational mass) and yet, in such a context, this would be a highly pertinent question, as we do assume a physical difference, analogous in this respect to an absolute distinction between positive- and negative-mass bodies. Indeed, why would the inertial mass remain positive when the gravitational mass is reversed? It is only confusion to pretend that there are multiple aspects of mass and that each of those independent mass properties can have a different sign. An electric charge is either positive or negative and mass, appropriately defined as the charge associated with the gravitational interaction, must also be either positive or negative. I think that we would be right to object trying to save the principle of inertia by assuming that some masses could be at once positive and negative, not because this would forbid all masses from always having the

same acceleration in a gravitational field, thereby allowing violations of the principle of equivalence, but simply because such a hypothesis would involve a contradiction. Clearly there is still something wrong, even with the second possibility that is traditionally considered for assigning physical properties to negative-mass bodies.

The preceding discussion should then have made clear the fact that there are two unsolved issues regarding negative mass. First, if we accept the requirement for a relational definition of the attractive and repulsive character of a gravitational field, then we must conclude that the currently favored assumption for what would be the behavior of negative-mass bodies, having at once negative gravitational mass and negative inertial mass, is incorrect, because, as I explained, it would involve absolutely defined properties of attractiveness and repulsiveness that would not depend merely on the difference or equality of the signs of the interacting masses. But if we consider the other traditionally considered (but not favored) possibility for the definition of negative *gravitational* mass, we may obtain the required relational definition of gravitational attraction and repulsion, but as I have explained a distinct problem would arise.

Indeed, under such conditions the behavior expected of negative-mass matter would have to be that which we currently assume to be shared by particles with a contradictory definition of their mass sign, which is not only objectionable on the basis of logical consistency, but which still involves a certain violation of the constraint of relational definition of physical attributes, by requiring one and only one type of gravitational mass to have an opposite inertial mass. Arguing that the problem here is with the notion that there exists only one single property of mass, while the difficulty can be avoided when the appropriate distinction is made between what we would call the inertial mass, which always remains invariant, and what constitutes the ‘real mass’, which we would call the gravitational mass and which may alone be reversed, would in my opinion not just be confused, it would be nonsense. What is positive cannot also at the same time be negative, if this polarity is to have any meaningful physical significance. Mass is not an abstruse, complicated property, with multiple independent and yet interrelated aspects, it is the gravitational charge and even though the stress-energy tensor replaces mass as the source of gravitational fields in a general-relativistic context, the lessons learned here are still valid and significant even in the context of the modern theory of gravitation.

It took me some time to realize that the problems we are dealing with

here (if we are willing to recognize that the whole question of identifying the properties of negative-energy matter is not itself insignificant) originate from what is usually assumed concerning the response to any force field in the case of a body with negative inertial mass. It is only after a rather long process of getting to understand the meaning of the phenomenon of inertia that I was finally able to gain the insight required to solve the problem of identifying the actual properties of negative-mass matter, in the context where we consider it a consistency requirement to impose on such matter that it should have both a negative gravitational mass and a negative inertial mass. Keep in mind that this explanation will be easier to grasp when the consequences of the integration of such a concept of negative-energy matter to the modern theory of gravitation will have been more thoroughly explored. Basically, what must be understood is that the direction of the equivalent gravitational field experienced by a given mass, in a reference system in which it is accelerating, even in the absence of nearby matter inhomogeneities, is in fact dependent on the sign of the mass that is accelerating. As a consequence, the inertial *force* associated with a given acceleration is left invariant, even if the sign of inertial mass is itself reversed along with the gravitational mass for a negative-energy particle.

In order to appreciate the following discussion at its true value, it is essential to remember that relativity theory does imply, in effect, that there exists a Newtonian gravitational field exerting a gravitational force on a positive-mass body which is accelerating relative to a local inertial reference system, even far from any large mass. The existence of the inertial force associated with this equivalent gravitational field is what allows a dynamic (by opposition to static) equilibrium to occur when an external force is applied on a body, which gives rise to an acceleration. Indeed, in the accelerated reference system relative to which a positive-mass body submitted to an external force does not accelerate, a gravitational force is present which balances the applied external force and this is what explains that there is no acceleration of the body relative to this particular (accelerated) reference system. In fact, the equivalent gravitational field is a general feature of acceleration and is present in any accelerated reference system, but in the absence of an external force to balance the associated inertial force the equivalent gravitational field only serves to determine the local inertial reference system associated with free-fall motion.

Indeed, given that the force associated with the equivalent gravitational field is a gravitational force, we must conclude that when the force respon-

sible for the acceleration is itself gravitational, we are actually in a situation where there would appear to be no force at all. It is therefore possible to assume that what determines the local inertial reference systems relative to which a positive mass experiences no gravitational force is the local matter distribution which is the source of the applied gravitational forces which are balanced by the inertial force which would otherwise be present relative to those reference systems (this is the essence of the insight that led to relativity theory). In any case, it is clear that the inertial force attributable to an equivalent gravitational field is always directed opposite the direction of the external force which gives rise to the corresponding acceleration, for a positive-mass body, and this means that the direction of the equivalent gravitational field experienced by a positive-mass body is opposite the direction of its acceleration, that is, opposite the direction of acceleration of the reference system relative to which this equivalent gravitational field exists. But what would occur if we had a negative-mass body in place of a positive-mass body?

First of all, it must be clear that the gravitational force $\mathbf{F}_g = m\mathbf{g}$ on a particle of mass m attributable to a given matter distribution would be reversed if the mass of the particle was reversed, because the Newtonian gravitational field vector \mathbf{g} at the particle's position would be left unchanged (because the matter distribution that is the source of the field does not change), while the sign of mass of the particle experiencing the field would be reversed. Now the problem usually is that when we want to determine the response of a particle to some gravitational force \mathbf{F} using Newton's second law $\mathbf{F} = m\mathbf{a}$, if the mass of the particle is reversed (negative), then the resulting acceleration \mathbf{a} would appear to have to be opposite that experienced by a positive mass submitted to the same force (the acceleration would be in the direction opposite that of the applied force). This is the traditional conception regarding negative mass. But if we consider things in a more general context, where Newton's second law would be an equation expressing the equilibrium between external forces \mathbf{F}_{ext} and the inertial force $\mathbf{F}_i = m\mathbf{g}_{eq}$ produced by the equivalent gravitational field \mathbf{g}_{eq} associated with a given acceleration, then we may write $\mathbf{F}_{ext} + \mathbf{F}_i = 0$ or $\mathbf{F}_{ext} = -\mathbf{F}_i$, so that for example if the external force is gravitational $\mathbf{F}_{ext} = \mathbf{F}_g = m\mathbf{g}$ then we would have $m\mathbf{g} = -m\mathbf{g}_{eq}$ and this means that the equivalent gravitational field \mathbf{g}_{eq} is usually opposite both the applied gravitational field and the acceleration, because in the present case we also have $\mathbf{F}_{ext} = m\mathbf{a}$, which means that $m\mathbf{g}_{eq} = -m\mathbf{a}$ for the considered positive mass m at least.

But would the equivalent gravitational field experienced by a negative-mass particle really be directed opposite the direction of its acceleration, as is the case for a positive-mass particle? To that question I think that, contrarily to what is usually assumed implicitly, we would have to answer that this cannot be the case. I will explain that, in fact, the equivalent gravitational field \mathbf{g}_{eq}^- that would be experienced by a negative-mass particle accelerating in a given direction away from any local matter inhomogeneity is the opposite of the equivalent gravitational field \mathbf{g}_{eq}^+ that would be experienced by a similar positive-mass particle with the same acceleration under the same conditions, so that we have $\mathbf{g}_{eq}^- = -\mathbf{g}_{eq}^+ = -(-\mathbf{a}) = \mathbf{a}$ for a negative-mass particle and given that we still have $\mathbf{F}_{ext} = -\mathbf{F}_i = -m\mathbf{g}_{eq}^-$ it means that $\mathbf{F}_{ext} = -m\mathbf{a}$ when the mass m is negative. If this is correct, then it would mean that the acceleration which a negative-mass particle would experience as a result of the action of a given force would actually be the same as that which would be experienced by a positive-mass particle submitted to the same force (not the same force field but really the same force), even if the mass, including the inertial mass, is indeed negative. The validity of this conclusion depends on only two assumptions. First, the proposed generalized Newton's second law (explicitly involving inertial forces instead of accelerations) must be considered more fundamental than the original formulation involving accelerations, so that the equilibrium it describes is really between forces and not merely between a force and an acceleration. Secondly, it must be assumed that the equivalent gravitational field associated with a given acceleration is reversed when the mass is reversed.

If the preceding conclusions are accurate it would appear that the fact that Newton's second law was always observed to work in its original form, that is, when the equivalent gravitational field is implicitly considered to be opposite the acceleration, is merely a consequence of the fact that it has only ever been verified to apply using positive-mass matter. But what is it indeed that might allow one to assume that the equivalent gravitational field would be reversed (would be directed in the same sense as the acceleration) for an accelerating negative-mass particle in comparison to what it would be for a similarly accelerating positive-mass particle? To understand what is going on we may consider the example of Einstein's elevator experiment. Indeed, we are allowed by the equivalence principle to assume that the effects observed inside an elevator accelerated in the vacuum, away from any local matter inhomogeneity, could also be explained by assuming that the elevator is not accelerating in the same vacuum (relative to the local inertial reference

system which would exist in the absence of any local matter inhomogeneity), but that it is instead maintained in place in the gravitational field of a large mass (located beneath the elevator) by the same external force which was originally causing it to accelerate. Thus, it seems that acceleration relative to a local inertial reference system always gives rise to an equivalent gravitational field similar to that which we would normally attribute to the presence of a local concentration of matter. We may then define an *equivalent source* to be the matter distribution which would give rise to the equivalent gravitational field experienced by an accelerated body if the presence of this field was not merely the consequence of acceleration.

Now, if we are allowed to assume that the equivalent gravitational *field* associated with the inertial gravitational force is actually reversed when the mass of the accelerated body is itself reversed (even without speculating about what the phenomenon of inertia might actually involve), it is simply because we can expect that the sign of mass of the equivalent source associated with the equivalent gravitational field experienced by a negative-mass body should itself be reversed. There should be no question, in effect, that if an accelerating positive-mass observer is allowed to assume that the equivalent gravitational field she experiences is actually attributable to the presence of an equivalent source with *positive* mass located in the direction opposite her acceleration, then a similarly accelerating negative-mass observer should himself be allowed to attribute the equivalent gravitational field that he would experience to the presence of some equivalent source with *negative* mass also located in the direction opposite his acceleration, otherwise we would have a way to determine in an absolute fashion, the positivity of mass.

Indeed, if it was always an equivalent source with positive mass (located in an invariant position relative to the accelerating body) that gave rise to the equivalent gravitational field, we could simply accelerate an observer of any mass sign and measure the equivalent gravitational field experienced by this observer, which could then be identified as the gravitational field attributable to a positive mass in the assumed position. Therefore, any gravitational field exerting on a given body a force such as that which was observed could be identified as the gravitational field of a positive mass, independently from the mere difference or equality between the polarity of the mass producing the field and that of the particle experiencing it. But this is a violation of the above discussed requirement of relational definition of the sign of mass. Thus, the problem with the traditional conception of negative inertial mass is that it would again allow to differentiate between positive and negative mass in an

absolute (non-relative) way, this time by referring to the predefined positive mass of the equivalent source whose gravitational field should invariably be observed under otherwise arbitrary motions of acceleration.

As it turns out, an additional difficulty arises when we try to assess the response of negative-mass matter to applied forces if we insist on assuming that the equivalent gravitational field associated with acceleration is an invariant property of the acceleration itself. Indeed, it is not only in the presence of an external force that the inertial force on a negative-mass body would have to be in the direction of its presumed acceleration when it is assumed that the equivalent gravitational field is opposite this acceleration (as is the case for a positive-mass body). The truth is that, when one recognizes the validity of the generalized form of Newton's second law, then under the inappropriate assumption that it is an equivalent source with positive mass that gives rise to the inertial force experienced by a negative-mass body in an accelerated reference system, it follows that even in the absence of external forces the inertial force would have the same direction as the acceleration, which means that the negative-mass body would actually accelerate in the same direction as the *accelerated* reference system itself. As a consequence, there would no longer be an equilibrium between the applied forces and the inertial force that is experienced by a negative-mass body due to its acceleration, which is certainly not a desirable outcome. Thus, even if the equivalent gravitational field experienced by an accelerating negative-mass body was the same as that experienced by a similarly accelerating positive-mass body, this would not give rise to the kind of motion which is traditionally expected from a negative-mass body.

What is important to understand, in effect, is that, in the context of a generalized formulation of Newton's second law, it must actually be imposed that there is always an equilibrium between the applied forces and the inertial force and under such conditions, the acceleration to which a body with a given mass sign is submitted is determined solely by the requirement that the inertial force it experiences actually balances the applied forces. Thus, once the direction of an applied force is known the acceleration of the body submitted to this force is determined only by the condition that it does, in effect, give rise to an inertial force which balances the applied force. But if the equivalent gravitational field which gives rise to the inertial force is dependent on both the direction of acceleration and the sign of mass of the accelerated body then the fact that the sign of mass would be reversed would not affect the direction of the acceleration, because the equivalent gravitational field

would also be reversed, which allows the inertial force associated with this acceleration to remain invariant under a reversal of mass.

Under such conditions, it would not be appropriate to assume that it is the sign of mass itself which determines the direction of the acceleration, because in fact the acceleration of a body submitted to a given force is determined merely by the requirement that the inertial force experienced by such an object balances the applied force in the accelerated reference system relative to which this inertial force is present. There is no *a priori* justification for considering that a negative-mass body with negative inertial mass should experience an acceleration opposite the applied force. This would be an incorrect interpretation of the classical equation between force and acceleration, which must be assumed to be valid only when the mass is positive. What the preceding argument shows, in effect, is that it would be a mistake to assume that the traditional formulation of Newton's second law also applies when the mass is negative. This equation does not apply when the mass is negative simply because the formula was not derived under the assumption that mass can be negative and was never intended to apply under such circumstances. But in the context of a generalized formulation of Newton's law and when the mass of the equivalent source responsible for the equivalent gravitational field is appropriately reversed for an accelerating negative-mass body, it follows that the equivalent gravitational field experienced by such an object must itself be opposite that experienced by a positive-mass body, which means that the inertial force remains unchanged, as does the body's acceleration.

If we are willing to recognize that it would be a serious inconsistency to allow for the same equivalent source (with the same mass sign) to give rise to both the equivalent gravitational field experienced by positive-mass particles and that experienced by negative-mass particles, then we must also recognize that similarly accelerating positive and negative-mass bodies would experience opposite equivalent gravitational fields, because those gravitational fields would arise from equivalent sources with opposite mass signs. But given that a negative mass must experience a force opposite that experienced by a positive mass of similar magnitude in response to any gravitational field, it follows that the inertial *force* actually has the same direction for both positive- and negative-mass bodies accelerating in the same direction, as a consequence of being submitted to the same external force (which is more constraining than requiring the same applied force *field*), even if we consider inertial mass to be reversed along with gravitational mass, as I previously

argued to be necessary.

In the present context, we would actually be allowed to assume that the requirement to consider that the equivalent gravitational field is reversed for a negative-mass body (in comparison with the equivalent gravitational field experienced by a positive-mass body with the same acceleration) is justified by the fact that it allows the dynamic equilibrium of forces on such an object to be maintained in the accelerated reference system relative to which this equivalent gravitational field is experienced, because if, in order to meet this constraint, we must consider the same inertial gravitational force to arise from the same acceleration, then it means that a negative-mass body would necessarily have to experience a reversed equivalent gravitational field, given that its mass is indeed reversed. No circular reasoning is involved here, because those results actually follow from the mere requirement of relational definition of the sign of mass applied to the equivalent source that gives rise to the equivalent gravitational field experienced by an accelerating negative-mass body.

For this argument to be valid, what must be recognized is that the negativity of the inertial mass of a negative ‘gravitational’ mass is an independent consistency requirement, which actually amounts to assume that mass is mass and that it cannot be both negative and positive at the same time and once this is acknowledged we are allowed to also and independently conclude that, just as there is not a unique sign of mass, there is not a unique equivalent gravitational field for bodies with opposite mass signs in the same accelerated reference system. In such a context we have no choice but to recognize that the response of a negative-mass body to any applied force would be that which we ordinarily (but inappropriately) attribute to a negative gravitational mass whose inertial mass would remain positive.

It is now possible to understand why it is that the inappropriate choice of a positive inertial mass in association with a negative gravitational mass would seem to agree, from a purely phenomenological viewpoint, with the independently motivated requirement of a relational definition of mass sign (given that it would allow gravitational attraction and repulsion to themselves be features dependent merely on the *difference* between the signs of gravitational mass of any two bodies). It is simply because, in such a case, instead of appropriately reversing the equivalent gravitational field for a negative mass accelerating in a given direction, we would reverse the sign of inertial mass (which must be negative for a negative-mass particle) a second time, from negative to positive again (while keeping the gravitational mass

negative), which, superficially, would be equivalent to simply reversing the direction of the equivalent gravitational field while keeping the mass negative as required. But I must emphasize again that, if that was the only possible approach to obtain consistent behavior from negative-mass bodies, then we would have to conclude that negative mass is not an appropriate concept in physical theory, because we would have to assume that a single unique physical property (what we may call the gravitational ‘charge’) is required to have at once and *from the exact same viewpoint* (for an observer with a given mass sign) two opposite values and this is clearly unacceptable.

It must nevertheless be mentioned that, as later developments will illustrate, it appears that, in fact, the reversal of the equivalent gravitational field is the trade-off we have to accept for keeping the value of the gravitational field attributable to a local matter inhomogeneity generally invariant while assuming that it is actually the mass experiencing it that can be reversed. But if, instead, we considered that the motion of a body must always be determined using the measure of gravitational field experienced by an observer made of matter with an invariant sign of energy, then it would be natural to assume that the sign of mass of the body (both inertial and gravitational) is positive definite, while it is the gravitational field attributable to a given matter inhomogeneity that is an observer dependent property.

From this viewpoint, the equivalent gravitational field due to acceleration far from any local matter inhomogeneities would no longer be dependent on the sign of mass of the accelerating body (because the mass itself would not change), while the gravitational field due to the presence of a local matter inhomogeneity would depend on the perceived sign of energy of its sources, which would become an observer-dependent property (again because the mass or energy of the body experiencing the fields would actually be considered positive definite). In this context there would then still be a practical (although not fundamental) distinction between an equivalent gravitational field due to acceleration far from any local mass concentration (which wouldn’t depend on the nature of the accelerating body) and the gravitational field due to the presence of a local matter inhomogeneity (which would depend on the nature of the object submitted to it). I will explain below what is the profound origin of this distinction and why it does not constitute an insurmountable difficulty for a consistent general-relativistic theory of gravitation based on the equivalence principle.

What must be retained here is that we can still consider the direction of the gravitational field attributable to the presence of a local matter in-

homogeneity to be an observer-independent property, while it is the mass experiencing it and therefore also the equivalent gravitational field experienced by this mass which may be reversed, but only at the price of changing the equations of motion which will be shown to otherwise describe the trajectories of particles submitted only to the gravitational interaction in a way that is equivalent to considering that the mass experiencing the gravitational field (attributable to this local matter inhomogeneity) is invariant, while it is the field itself which is reversed (in comparison to what it would be if we had considered its effect on a body with reversed mass). Now, if we do consider the mass (both gravitational and inertial) of the particle experiencing a gravitational field to always be positive definite, so that that it is the direction of the gravitational field itself which varies as a function of the *relative* difference between the observer-dependent sign of mass of the source (which can still be either positive or negative) and that of the particle experiencing the field (which would always be assumed to be the positive one) then we obtain a framework that can be more easily generalized to a relativistic theory. But it must be clear that the two approaches discussed here are equivalent in the Newtonian context and still require all mass (gravitational and inertial) to be either positive or negative and when the direction of the gravitational field due to a local matter inhomogeneity is not considered to be an observer-dependent property we must indeed consider the equivalent gravitational field to itself be dependent on the sign of the accelerated mass (which is no longer positive definite), otherwise the equivalence between the two viewpoints breaks down.

From the viewpoint where the mass experiencing a gravitational field is considered positive definite, a Newtonian gravitational field experienced by a particle we would normally consider to have positive mass, if it is not the result of an accelerated motion far from any matter inhomogeneity (in which case we would be dealing with an equivalent gravitational field), would be experienced by a particle we would normally consider to have negative mass as an oppositely directed Newtonian gravitational field, while the mass of the particle experiencing this relatively defined gravitational field would not even show up in the equations used to determine its motion. But if the gravitational mass experiencing this reversed gravitational field is kept positive, then it must be assumed that the inertial mass is also kept positive and under such conditions the equivalent gravitational field would appear not to be reversed. It is because we do not appropriately keep the sign of the mass experiencing the equivalent gravitational field invariant when we try to determine the

motion of what we currently describe as a negative-mass particle in an accelerated reference system that we need to consider this gravitational field to be reversed. But when the external force applied on what we would currently describe as a negative-mass particle is gravitation itself, it is possible to assume that this force is reversed (from that which would be experienced by what we currently describe as a positive-mass particle), not because the mass of the particle is reversed, but because the local gravitational field itself is reversed. In such a case the inertial *force* would not be reversed, because the mass (both gravitational and inertial) that is experiencing the field is not reversed and it must also be assumed that the equivalent gravitational field is left unchanged (is identical to that which is experienced by what we already consider to be a positive-mass particle). Therefore, acceleration still doesn't take place in the direction opposite the applied force and this is all a consequence of the fact that even though the local gravitational field appears to be reversed from such a perspective, the equivalent gravitational field, in contrast, is left invariant along with the sign of mass of the particle.

It should be clear, then, that in the context of an approach according to which the particles experiencing a gravitational field are always assumed to have a positive mass, the crucial assumption is that while the gravitational fields attributable to local matter concentrations are dependent on the nature of the body experiencing their effects, the equivalent gravitational field associated with acceleration away from local masses would, for its part, remain invariant, regardless of how the body experiencing it perceives the gravitational fields attributable to local matter inhomogeneities. This hypothesis can be considered to be equivalent to that which in the above described approach consists in assuming that the equivalent gravitational field must actually be reversed for a negative mass, because this is indeed what allows the inertial properties of an object to be independent from its mass sign. I believe that this observation clearly shows that I'm justified in analyzing the problem of negative mass from a conventional perspective, according to which the mass experiencing a gravitational field is explicitly assumed to be reversed, because in such a context the underlying assumptions are made more apparent and it is also easier to explain what I'm referring to when discussing the case of anomalously-gravitating matter. In a Newtonian context I will therefore continue to use the first viewpoint, according to which it is possible for the mass experiencing a gravitational field to be negative.

Now, we may want to dig a little deeper and ask why it is exactly that we are allowed to assume that the direction of the equivalent gravitational

field is dependent on the sign of mass of the object experiencing it? I have tried very hard to develop a better understanding of the whole phenomenon of inertia and what I have learned has actually helped me to derive the above discussed results. Indeed, this investigation has enabled me to realize that the assumption that the equivalent gravitational field is reversed, when the mass which is subject to acceleration is itself reversed, is not just a requirement of the necessary relational definition of the sign of mass, but must be imposed in order to allow a relational description of the phenomenon of inertia itself, in the sense that inertia should be conceived as arising from purely relative motions between matter particles, as suggested by Ernst Mach a long time ago. In this context, I have become convinced that the inertial forces acting on a particle can be understood to arise as a consequence of an imbalance, caused by acceleration relative to the global inertial reference system (associated with the distribution of matter on the largest scale), in the sum of forces attributable to the interaction of the accelerating particle with each and every other particle in the universe.

What happens, in effect, is that there must be a similar imbalance of the gravitational forces exerted on similarly accelerating positive- and negative-mass bodies arising from their interaction with the rest of the matter in the universe, because the imbalance responsible for the existence of the inertial gravitational force is similar to a skewed mass distribution and if the actual large-scale matter distribution responsible for those effects is roughly the same from the viewpoint of both positive and negative masses, in the absence of local matter inhomogeneities, then the imbalance should develop in a similar way for both positive and negative masses from the viewpoint of their own mass sign. Thus, what must be retained of this investigation is that the equivalent gravitational field which applies on a negative-mass body should in fact be the opposite of that which would be experienced by a positive-mass body with the same acceleration that is located within the same matter distribution, even if simply as a consequence of the fact that for a reversed mass the same motion relative to the same matter distribution should give rise to a similar imbalance in the sum of *forces* attributable to interaction with all the matter in the (visible) universe.

Indeed, given that the mass itself is reversed, the invariance of this imbalance would mean that the equivalent gravitational field responsible for the inertial force must also be reversed in the accelerated reference system, so that the force existing relative to it can itself be left invariant. But if the equivalent gravitational field associated with the acceleration of a negative-

mass body is the opposite of that associated with the same acceleration of a positive-mass body, it follows that the reaction to any applied force is indeed the same for opposite-mass particles, despite the fact that there is no distinction between inertial and gravitational mass signs (even for negative-mass particles). This may be considered to actually explain why it is appropriate to assume that it is the inertial force itself, instead of merely the product of mass and acceleration, that would be opposite the direction of the applied external force for a negative-mass body, as the generalization of Newton's second law that I proposed allows to express.

But it must be clear that if there is a requirement for inertial mass to be reversed, along with gravitational mass, it does not follow from imposing the validity of the equivalence principle as a condition that all matter should have the same acceleration in the absence of any interaction other than gravitation, as is usually considered. Indeed, as the previous analysis allows to understand, even a negative-mass body for which both the gravitational and the inertial masses are negative cannot be expected to follow the same trajectory as a positive-mass body in the presence of a local positive or negative mass concentration (despite what is usually assumed). What I have tried to explain is precisely that, even when inertial mass is assumed to be reversed along with gravitational mass, it is not possible to preserve the validity of the equivalence principle integrally. Thus, a local inertial reference system cannot be defined independently from the sign of mass of the body experiencing it, given that the direction of the gravitational force resulting from a particular matter distribution depends on the sign of mass of this body. What I will explain in the following section is that the requirement that all matter with the *same mass sign*, in the same location, experiences the same acceleration is in fact restrictive enough for the equivalence between gravitation and acceleration to apply in a certain way, that allows a metric theory of the gravitational field to emerge which merely relativizes the curvature of spacetime by making it an observer-dependent aspect of reality.

2.5 The equivalence principle with negative mass

It is not usually recognized that the general theory of relativity is actually based on two postulates, because only the first postulate, which concerns the

equivalence between acceleration and a Newtonian gravitational field, is well-known and is explicitly taken into account. But actually, a second postulate is required to obtain the current formulation of the theory and is implicitly assumed to be valid without justification. It is the hypothesis of absolute significance of the sign of energy. This second assumption appears to be necessary in order to preserve the validity of the first postulate under conditions where the presence of negative-energy matter would, in effect, need to be taken into account. But even though the postulate of the absolute definiteness of the sign of energy may be considered problematic in the context of the preceding analysis, it remains to be shown whether it is possible to provide a consistent classical theory of the gravitational field in which only this second postulate would be rejected. Thus, I will try to show, in this section and later on, when discussing the mathematical aspects of a generalized theory of gravitation, that it is perfectly possible and indeed actually necessary to maintain the validity of the equivalence principle in a certain form, while nevertheless rejecting the assumption of an absolute significance of the sign of mass or energy.

First of all, it must be emphasized that the true motivation behind the equivalence principle is to be found in a requirement which we may call the principle of relativity and which is actually one particular expression of the requirement of relational definition of all physical quantities. This relativity principle imposes that the state of motion of an object, and in particular its rate of acceleration, is to be determined merely in relation to the state of motion of other physical systems, so that there is no absolute state of acceleration relative to an arbitrarily-chosen, unique, metaphysical reference system. The principle that there is an equivalence between a Newtonian gravitational field and an acceleration enables this requirement to be fulfilled, because it allows what might have otherwise appeared to be an acceleration relative to absolute space to merely be a state of rest in the vicinity of a local mass concentration not accelerating relative to the same ‘absolute’ space, as Einstein understood, but as we tend to ignore nowadays in favor of the mere mathematical requirement of general covariance of the field equations. I think that it must be recognized that, in fact, the only essential implication of the equivalence principle is that there is no longer any motive for arguing that because acceleration is felt (unlike velocity) it must be absolute. Thus, it may appear problematic that even if we can find generally covariant equations for the gravitational field in the presence of negative-energy matter, the fact that according to the previous analysis such matter would not share the same

accelerated motion as positive-energy matter in the presence of a local matter inhomogeneity (while it should in the absence of such a perturbation, for reasons I explained before) would appear to allow the effects of acceleration relative to matter at large to be distinguished from those attributable to the gravitational field of a local mass.

There is indeed a tension between the principle of relativity and the previously discussed requirements concerning negative-mass matter which we may illustrate by once again using Einstein's elevator experiment. Under circumstances where what I have identified as appropriately behaving negative-energy matter would be present it may seem, in effect, that we could differentiate an acceleration of the elevator occurring far from any local mass from an acceleration of the elevator occurring while it is at rest near such a large mass. This is because, near a planet or another large matter inhomogeneity, positive- and negative-mass bodies would accelerate in opposite directions, one toward the local mass and the other away from it (one upward, the other downward), while in the elevator which is simply accelerating far from any large mass, positive- and negative-energy bodies would share the same acceleration, apparently betraying the fact that the acceleration is 'real'. We may, therefore, assume that an observer in the elevator would be able to tell when it is that she is simply standing still in the gravitational field of a planet and when it is that she is actually accelerating far from any big mass. The 'true' acceleration would have been revealed to the occupants of the elevator as that for which both the positive- and the negative-mass bodies have the same acceleration. Consequently, we would seem to be justified to conclude that the notion that the effects of acceleration are totally equivalent to those of a gravitational field (which is the essence of the principle of equivalence) is no longer valid when we introduce negative-mass matter with properties otherwise required to make it a consistent concept (according to the preceding analysis).

Indeed, I made it clear before that it is not possible to abandon the principle of inertia or Newton's third law (action and reaction) in order to accommodate the existence of negative-mass matter, because if those rules were not strictly obeyed under all conditions then not much else would remain valid. We cannot even tell what a world devoid from this constraint would look like and there is no reason to assume in particular that the equivalence principle itself would still be obeyed, as is usually assumed, because, after all, this principle is a reflection of the phenomenon of inertia. Trying to save the principle of equivalence by simply allowing negative-mass matter

to react anomalously to applied forces (as if that was required when inertial mass is negative), so that it can accelerate in the same way positive-mass matter does in the presence of local matter inhomogeneities, would not make sense, because this would mean that the principle of inertia no longer applies in general and again, in such a case there is no guarantee that even the alternative situation we expect to observe under those conditions would really occur. I believe that there are reasons why no violations of the principle of inertia have ever been observed despite the fact that the techniques required to reveal such transgressions have long been available. It would not be clever to think that it is by rejecting this principle that we can maintain the requirement of the equivalence between a gravitational field and acceleration. Clearly, there must be something wrong with certain assumptions we take for granted concerning the equivalence principle itself. The fact that this is the principle upon which relativity theory and our modern concept of gravitation is founded should not prevent us from reexamining some of the implicit assumptions surrounding it. Failing to do so would mean that we have to give up on the idea that negative-energy matter can exist, because only so could we then avoid being faced with the annoying and unpredictable consequences of an alternative choice concerning the properties of this matter.

It is important to note, at this point, that it would be inappropriate to suggest that it may be possible to accommodate the requirement that the principle of equivalence also applies in the presence of negative-mass matter by assuming that opposite-mass bodies always share *opposite* accelerations instead of always sharing the same acceleration, as is traditionally believed. It is certainly true that, under such circumstances, it would still be impossible to distinguish a true acceleration given that opposite-mass bodies would always accelerate in opposite directions, whether those accelerations are the result of the presence of a local concentration of matter or the result of the presence of an equivalent gravitational field far from any large mass. But this situation could only occur, in the context of an appropriate conception of the phenomenon of inertia based on the previously discussed generalized formulation of Newton's second law, if it was assumed that the equivalent gravitational field associated with acceleration is not reversed despite the reversal of the mass of the accelerated body experiencing it.

From that viewpoint we should actually expect that one of two opposite-mass bodies would fall down while the other would fall up in the accelerating Einstein elevator far from any local mass, even when no force is applied on any of the two masses independently. However, this kind of behavior would

constitute an even more severe violation of the principle of inertia than that which would occur in the case of the chasing pair of opposite-mass bodies described before, given that, in this case, there wouldn't even exist any identifiable cause for the upward acceleration of one of the two bodies, because the elevator does not even interact with any of the masses and merely constitutes a reference system. In fact, this situation is so devoid of plausibility that it clearly means that it is not possible to try to salvage the equivalence principle by assuming that the equivalent gravitational field is not reversed for an accelerating negative-mass body. The fact that the kind of uniqueness of the equivalent gravitational field that is involved here would also violate the requirement of relational definition of the sign of mass, as I explained in the previous section, only contributes to confirm the validity of this conclusion. We must therefore accept that while the local inertial reference systems can differ for positive- and negative-mass bodies near some local matter inhomogeneities, they must nevertheless be identical for opposite-mass bodies far from local mass concentrations.

I will soon explain why it is exactly that we are allowed to consider that the principle of relativity of motion (concerning acceleration in particular) is not threatened by the conclusion that the free-fall state of motion of a negative-mass body can be different from that of a positive-mass body in the presence of local matter inhomogeneities. But it is important to first point out that in the case of the elevator suspended in the gravitational field of a local mass we are, in effect, considering an inhomogeneous matter distribution for which positive- and negative-energy matter concentrations are *not* superposed in space (in the classical sense) and therefore do not produce mutually compensating local gravitational fields. If such compensations between the effects of *local* matter inhomogeneities were to occur, as would be the case for example in the presence of two superposed gas clouds of opposite energy signs with the same overall motion or rotation, then the acceleration of positive- and negative-energy bodies located near or within those matter distributions would be the same despite the presence of local inhomogeneities in the configuration of positive- and negative-energy matter. This actually means that there couldn't be any effect from the motion relative to such a matter distribution, because whatever gravitational effect positive-energy matter would have, would be compensated by the opposite effect of the similarly distributed negative-energy matter present around the body. This is true also of rotation, which according to Einstein's theory induces a frame dragging effect which we may assume to be dependent on the sign of mass

like any other gravitational phenomenon.

Now, you may recall this earlier discussion (from the preceding section) in which I suggested that it should be possible to attribute the inertial gravitational forces experienced by positive- and negative-mass bodies in the accelerating elevator away from local masses to some imbalance in the sum of gravitational forces attributable to interaction with all the matter in the universe, arising as a consequence of acceleration relative to the reference system associated with the average state of motion of this large-scale matter distribution. However, given what I just mentioned regarding the compensating effects of superposed matter distributions with opposite masses and identical motions, it seems that one would have to assume that no imbalance could arise from the gravitational interaction with positive- and negative-energy matter if they are similarly distributed in space on the largest scale. Thus, one must conclude that if the positive- and negative-energy matter distributions are indeed mostly homogeneous and are at rest with respect to one another on such a scale (as appears necessary if the cosmological principle applies equally to both matter distributions), then there should be no effect on both positive- and negative-mass bodies from the presence of matter on the largest scale.

What this means is that there could not be any imbalance in the equilibrium of gravitational forces attributable to the large-scale matter distribution that would give rise to inertial forces or the equivalent gravitational fields, because one imbalance, attributable to motion relative to positive-energy matter, would be compensated by a similar, but opposite one arising from the same motion relative to negative-energy matter (all masses would experience two opposite, equivalent gravitational fields all at once). It thus appears that there is something wrong with one or more of the implicit assumptions entering this deduction, because inertia does exist and indeed, if there was no inertia, the world would not be anything even remotely similar to what we observe. Of course, the idea that there simply never was any negative-energy matter in the universe (so that the imbalance due to acceleration relative to the positive-energy matter distribution is not compensated by an imbalance due to acceleration relative to the superposed negative-energy matter distribution) may be tempting, because after all we do not observe any such matter. But keep in mind that it will later be explained that this hypothesis is not required and that, in any case, it would again amount to simply reject the possibility that such matter may exist, without providing any justification for this very convenient hypothesis.

We may summarize the situation by noting that what we know for sure is that if the expected identical accelerations of the opposite-energy bodies relative to the elevator far from any local mass are due to a similar imbalance in the gravitational forces attributable to the interaction of those bodies with matter on the largest scale, then this imbalance must be attributed to a motion that takes place relative to opposite-energy matter distributions which share the same motion (or absence of motion) and the same rotation and which *should* therefore have mostly compensating effects on positive- and negative-energy bodies with the same motion *relative* to this homogeneous matter distribution. If this is recognized, then we have to admit that in the context where negative-energy matter actually exists it would be difficult to see how a local inertial reference system could be determined by the large-scale matter distribution through the gravitational interaction. In such a case, it would then seem that we have to conclude that there may need to exist something like absolute acceleration relative to an arbitrarily-chosen, unique reference system lacking any physical underpinning. What I have understood though (for reasons that will be discussed later) is that the hypothesis that both the large-scale positive- and negative-energy matter distributions have an effect on positive- or negative-energy bodies, considered independently, constitutes the incorrect assumption which appears to invalidate the hypothesis that all motion (including accelerated motion) is relative, even in the presence of negative-energy matter.

If we drop the assumption that a *negative-energy* matter distribution that is uniform on the cosmological scale can exert a force on *positive-energy* matter (and vice versa for the effects of positive-energy matter on negative-energy matter), then it seems that we can explain the imbalance responsible for the force of inertia as being the consequence of an acceleration with respect to the one particular, but relatively defined, reference system which is that relative to which most of the matter in the universe is at rest, because, in such a case, there would be no canceling of the effects attributable to the positive-energy matter distribution by those of the negative-energy matter distribution (and vice versa) on the largest scale. Therefore, what I suggest we have to recognize, even if only by necessity, is that there is no compensation, for a positive-mass body accelerating relative to the average matter distribution on the cosmological scale, between the equivalent gravitational field attributable to positive-energy matter and that which we could have attributed to negative-energy matter. Similarly, there should be no equivalent gravitational field attributable to acceleration relative to the average

distribution of positive-energy matter to compensate the equivalent gravitational field attributable to acceleration relative to negative-energy matter for a negative-mass body. I believe that this is due merely to the fact that, on the cosmological scale, particles of one energy sign interact only with the matter distribution that has the same energy sign. I'm particularly confident in the validity of this conclusion, given that I had actually understood the requirement of absence of interaction between a positive-energy body and the uniform, large-scale distribution of negative-energy matter before I even realized that it was required to solve the problem of the relativity of motion, in the context where negative-energy matter is allowed to exist. I will explain what independently justifies this conclusion in sections 2.6 and 2.8.

What happens, therefore, is that only the very-large-scale distribution of *positive-energy* matter determines the local inertial reference system that is experienced by *positive-energy* bodies in the absence of local matter inhomogeneities, while only the overall distribution of *negative-energy* matter determines the local inertial reference system experienced by *negative-energy* bodies in the absence of local matter inhomogeneities (this language would also be appropriate from a general-relativistic viewpoint). Thus, what differentiates the situation of the elevator near a large mass of positive *or* negative sign and the situation we have in the elevator accelerating far from any such local mass is that, in the first case, the force responsible for the observed acceleration is the result of an imbalance that is caused by unequally distributed inhomogeneities in the positive- and negative-energy matter distributions and this imbalance is dependent on the sign of energy of the body experiencing it (as there are two possibilities for both the sign of mass of the source and that of the accelerated body), while, in the latter case, the observed force responsible for the acceleration is the result of an imbalance that is always caused by the motion of a body of given mass sign relative to a uniform matter distribution with the same mass sign (necessarily and invariably), so that it is not dependent on the sign of energy or mass of the body experiencing it (positive- and negative-energy bodies react in the same way to acceleration relative to matter on the largest scales), as long as the distributions of positive- and negative-energy matter are both homogeneous and are not accelerating or rotating relative to one another on the largest scale.

All accelerations are therefore relative accelerations between well-defined physical points of reference within the universe and no absolute state of rest (more exactly of absence of acceleration) can be identified. This is true even

if there does exist a unique particular reference system (actually two unique, but corresponding reference systems) which is singled out as that relative to which the motion (state of acceleration) of positive- and negative-mass bodies is the same in the absence of local disturbances, as a result of the correspondence of the average state of motion of the positive- and negative-energy matter distributions on the largest scales. But this conclusion applies merely in the context where, *globally*, any particle is gravitationally influenced only by its interaction with matter of the same energy sign, whose state of motion relative to the particle, therefore, alone determines the local inertial reference system in which the particle evolves. Thus, despite the expected correspondence of the states of motion of the uniform positive- and negative-energy matter distributions on the largest scale (which may seem to imply an absence of resulting effect on both positive- and negative-mass bodies), there nevertheless exists a resulting effect from the presence of this matter on a local mass of any sign that allows to determine a unique reference system and this is what explains that there appears to be a difference between acceleration far from any local mass and the acceleration attributable to the gravitational force of local matter inhomogeneities, while, actually, the difference observed is merely the consequence of the fact that a body with a given mass sign interacts only with the large-scale matter distribution with the same sign of mass, so that no compensation can exist in this case.

In light of those developments, it appears that what the previously discussed insight concerning the nature of the equilibrium involved in determining local inertial reference systems should be understood to mean is that free-fall motion, instead of involving a total absence of forces, as is usually assumed in a general-relativistic context, must be considered to be the consequence of an acceleration-dependent equilibrium in the sum of gravitational forces attributable to interaction with both local masses and the large-scale matter distribution. This interpretation appears to be required in the context where negative-energy matter must be recognized to exist, given that, in such a case, there cannot even be a unique inertial, or free-fall reference system dictated by the geometry of spacetime, so that we are forced to consider the reality of the general-relativistic gravitational field as being associated with such a physical interaction. Indeed, it is only when we are dealing with a universal force, defined precisely as a force that affects all bodies in the same way, that we can *choose* (as a mere convention) to include this force in our definition of the metric properties of space and time, given that, in principle, geometry must be shared by all objects present in the related space.

What remains to decide is whether this convenient choice is still appropriate for gravitation, in the context where the force in question can no longer be assumed to affect all bodies similarly (therefore betraying its material nature).

Einstein himself insisted that once we recognize the validity of a principle of general relativity of motion, then the speed of light can no longer be assumed to be constant (even though it is left invariant locally, along a geodesic), given that, in the elevator experiment, light rays may follow curved paths. But, from this viewpoint, the curvature of spacetime should naturally be expected to arise as the consequence of a local perturbation in the equilibrium of gravitational forces attributable to the interaction of the bodies experiencing it with all the matter in the universe (except the large-scale matter distribution with opposite mass sign), otherwise it would be impossible to determine what affects the trajectory of light in an accelerated reference system far from any local matter inhomogeneity. Indeed, even in a flat space, far from any local matter concentration, the motion of light in a straight line, which is usually considered to be a consequence of geometry itself, would, from my viewpoint, be a consequence of the equilibrium of forces arising from the gravitational interaction with the rest of matter in the universe. This does not mean, however, that the geometrical interpretation of gravitation is incorrect, but merely that the geometrical properties of space must definitely be conceived as arising from those interactions and more precisely, from some sort of equilibrium in the sum of gravitational forces that can be altered by the presence of local matter inhomogeneities. As I will explain in section 5.13, such a viewpoint has the added benefit of being more easy to generalize to a theory where the gravitational interaction must not only be described as an interaction mediated by quantum particles, as is already recognized to be necessary, but must really be integrated into the quantum framework in the manner I shall propose.

In any case, I think that it is clear that statements to the effect that relativity theory has made the concept of gravitational interaction obsolete and replaced it with that of spacetime curvature (so that gravitation is merely a manifestation of the geometry of spacetime) can no longer be assumed meaningful, if curvature is itself a relatively defined property which arises as a consequence of an equilibrium of local and inertial gravitational forces which depend on the sign of energy of the objects involved. I think that the situation we have here is similar to that in which electromagnetic theory was before the quantization of energy and the photon concept were proposed,

because spacetime is now viewed as a continuous medium that directly takes part in determining the motion of objects, just like the electromagnetic field was originally considered to be a continuous wavelike phenomenon, directly influencing the motion of charged bodies. When it was shown that light is a corpuscular phenomenon, the whole notion of electromagnetic wave was not abandoned, of course, because there was something real about the wavelike character of electromagnetic phenomena and this is the element which came to be integrated into quantum mechanics. Similarly, I think that the concept of spacetime curvature cannot and need not be abandoned when gravitation is described as an interaction which, ultimately, would need to be mediated by the exchange of quantum particles in a way that would allow mass-sign-dependent local inertial reference systems to emerge, only, the curvature of spacetime can no longer be considered as actually *being* gravitation itself.

As Hans Reichenbach once emphasized [21] (p. 256), if we choose to integrate the gravitational force into our definition of spacetime we may no longer need to explicitly take the force into consideration to explain the motion of bodies, but we must still invoke a force as the cause of the geometry itself. Thus, it is not gravitation which was replaced by curved geometry, but all of geometry that became a manifestation of the universality of the gravitational interaction, and I think that this is particularly relevant in the context of a theory of gravitation that allows to take into account the possibility of the existence of negative-energy matter. Actually, the commonly made remark to the effect that relativity allowed to eliminate gravitation as a real force appears to be motivated by the fact that the gravitational force arising from local mass concentrations was given the status of inertial force (similar in kind to the Coriolis force) by relativity and given that inertial forces were never seen as real forces, then it is believed that gravitation can now be considered a fictitious force under all circumstances. But I believe that it is rather the contrary that is true and that it is the inertial forces which can be considered as real gravitational forces in a general relativistic context. The fact that inertial forces are involved in giving rise to the dynamic equilibrium which determines the mass-sign-dependent local inertial reference systems would then be a further indication that the geometry of spacetime is the product of an equilibrium of real gravitational forces arising from the interaction of local masses with the rest of matter in the universe.

Having properly identified the origin of the identical response of positive- and negative-mass bodies to acceleration, I do not want to immediately en-

ter into a discussion as to what are the true elements of justification behind the assumption that particles with a given mass sign are not affected, from a gravitational viewpoint, by the presence of a uniform distribution of matter of opposite mass sign on a cosmological scale. But it may nevertheless already be noted that the fact that one particular reference system appears to be singled out as having unique status among all possible states of acceleration is not a unique feature of the approach described here. Actually, in a general-relativistic context, even in the absence of negative-energy matter, this feature of our description of the motion of objects should appear all the more natural given that all inertial reference systems are an outcome of the gravitational interaction and are therefore determined by the surrounding matter distribution. There exists, in effect, one very particular reference system in our universe, which we may call the global inertial reference system and which is that which is determined by the average motion of all masses together and relative to which most masses in the universe do not accelerate in the absence of a local force. That there may be such a unique point of reference does not mean that it is not relationally defined. Relativity theory allows to explain the existence of this particular reference system as being a result of the combined gravitational interactions of a local body in any state of motion with all the other masses in the universe (with the same mass sign) and therefore in relation to the average motion of those masses. Indeed, even far from any big mass, there remains the gravitational effect of the universe as a whole, which can never be ignored. Thus, the situation we usually refer to as corresponding to an absence of gravitational field and which we expect to be experienced far from any local mass concentration, is not different, in fact, from that occurring in the presence of such a local mass, only it is characterized by the fact that the gravitational field is then attributable to the average distribution of either positive- or negative-mass matter present on the cosmic scale and cannot be compensated by the presence of matter with an opposite mass sign, as long as all matter is uniformly distributed on such a scale.

The fact that inertial reference systems are always determined by the average state of motion of matter in the universe becomes particularly obvious when we consider the reference system associated with a felt motion of rotation which, as experiments have revealed, must be one that takes place relative to the most distant galaxies and therefore relative to the largest ensemble of matter in the universe. The reference system relative to which a positive-mass observer feels no rotation must then be determined simply

by the gravitational field attributable to all matter particles with the same mass sign present in the visible universe, in a way that is dependent on the average state of motion of those particles. Such a reference system, therefore, is definitely unique, even though its description involves only relationally defined properties. We may still consider the average matter distribution on the largest scale to be rotating, but then its gravitational field would give rise to a rotating inertial reference system which, through relativistic frame dragging, would put the whole matter content of the universe in rotation with it³. Since Einstein, there is no longer any mystery with the existence of such a preferred reference system and what I'm trying to explain is that there is also no problem with the fact that there is a unique reference system relative to which both positive- and negative-mass bodies have no acceleration when free from external non-gravitational forces. We are not faced here with a metaphysical reference system associated with absolute acceleration, but merely with an ordinary reference system relative to which the sum of allowed gravitational interactions of local masses with the ensemble of matter present on the largest scale imposes an absence of acceleration that is shared by positive- and negative-mass bodies.

Again, it must be stressed that even when it may seem that we are dealing with empty space, what the objects actually experience are the effects of the whole surrounding matter distribution conveyed by the gravitational field as an intermediary material entity, which, in a general-relativistic context, actually determines the possibly distinct local inertial reference systems affecting positive- and negative-energy bodies. This aspect of the general-relativistic (or physical) space is what allows to conceive of rotation as being purely relative, even when the distance of some objects to the rotation axis of a rotating observer becomes large enough that the objects would actually have to move at faster-than-light velocities in the reference system tied to the observer. Indeed, it is the rotation of the whole gravitational field, as a material entity (which would also occur in a universe totally devoid of 'real' matter), that explains that this motion of the remote objects is possible as a true

³It has been mentioned that a (positive-mass) observer uniformly rotating with respect to the distant stars and which would choose to consider himself motionless would observe a gravitational field which from a Newtonian viewpoint could not exist, therefore weakening the equivalence principle. But it is interesting to observe that this difficulty would no longer exist in the context where a repulsive gravitational field that grows in proportion to the distance from an axis could be produced by an appropriately configured inhomogeneous, static distribution of negative-energy matter.

motion, because locally the objects are not moving (accelerating) relative to the gravitational field (or the local inertial reference systems), which is then itself rotating, and this is what makes their large velocities and accelerations possible, as is already well understood.

But if acceleration occurs merely relative to the inertial reference systems determined by the gravitational field, it must not be forgotten that the state of motion of matter also contributes to determine the gravitational field and therefore it should naturally be expected that there is no acceleration of matter as a whole relative to the global inertial reference system determined by the gravitational field produced by this large-scale matter distribution. It may also be remarked that the situation we are dealing with here, concerning the relativity of acceleration in the presence of negative-energy matter, is similar to that regarding the relativity of velocity, because there also exists a preferred reference system relative to which the temperature of the cosmic microwave background is mostly uniform and which may appear to define a state of absolute rest, but this unique reference system is merely that which is not moving relative to the average state of motion (not acceleration) of matter on the largest scale. If there is no conflict with the principle of relativity in such a case, then there need not be a problem in the case of the global inertial reference system singled out as being that relative to which there is no difference between the states of acceleration of freely falling positive- and negative-mass bodies.

There would then be no substance to the argument that the distinction between acceleration and gravitation, which appears to be revealed by the distinct accelerations of positive- and negative-energy bodies in the standing still elevator near a local mass (in the context where negative-energy matter does not respond perversely to applied forces), allows absolute acceleration (or absolute absence of acceleration) to be determined. Indeed, the local gravitational fields and the associated local inertial reference systems are always determined in a relative fashion as dependent on the presence of the local masses which are the source of the fields, while the reference system where the states of acceleration of positive- and negative-energy bodies are identical is determined as that relative to which the large-scale matter distribution (which we may assume to be in the same average state of motion for positive- and negative-energy matter) is itself not accelerating. This all follows from the fact that positive- and negative-energy bodies interact only with the homogeneous matter distribution with the same sign of energy as

their own on the cosmological scale⁴, so that motions relative to those matter distributions must be treated differently from motions relative to local matter inhomogeneities, although they are still relative motions.

It must be noted, however, that if the distributions of positive- and negative-energy matter were in motion relative to one another on the largest scale, there would then actually be two different global inertial reference systems associated with the two types of mass (positive and negative) experiencing them, even away from any local mass. In such a case it would be more difficult to differentiate between the situation of the elevator far from any large mass and that in which unequally distributed concentrations of positive- and negative-mass matter are present locally. It remains, though, that if a certain condition of zero energy and momentum, which will be discussed in section 4.5, must be imposed on the universe as a whole, then, in the absence of very-large-scale inhomogeneities in the two matter distributions, we should not expect negative-energy matter to be accelerating or even only moving, on the average (on the largest scale), relative to positive-energy matter, because negative-energy bodies have momentum pointing in the direction opposite their motion, which means that the global inertial reference systems associated with positive- and negative-energy matter should be indistinguishable, especially in the context where there exists a constraint on the magnitude of density fluctuations in the initial Big Bang state (as I will explain in section 4.9).

Based on the above discussed considerations, I have come to the conclusion that, after all, the principle of relativity is not really threatened by the introduction of negative-energy matter obeying the requirement of relational definition of its mass sign. But clearly the equivalence principle itself (which allows accelerated motion to be treated relativistically) is no longer to be considered valid in the sense it was traditionally believed to be and if it need not and indeed cannot be abandoned it must, however, be generalized or somewhat relativized. In fact, we already know for sure that the equivalence principle always applies only in local reference systems whose states of motion can be different in various locations. We can tell, in effect, that a gravita-

⁴In fact, as I will later explain, the large-scale distribution of negative-energy matter may exert an influence on positive-energy bodies, but only when inhomogeneities are present in this matter distribution. The nature of those interactions is such, however, that there is necessarily a cancellation in the sum of the effects involved on the largest scale, so that there can be no overall effect and the same is true for the effects of positive-energy matter on negative-energy bodies.

tional field is attributable to the presence of local masses instead of being the consequence of an acceleration, even in the total absence of negative-energy matter, when we consider a portion of space that is sufficiently large. For example, if we consider two elevators suspended on opposite sides of a planet, instead of a single elevator, it is obvious that even though observers in each of those elevators could assume that they are accelerating far from any local mass, from the global viewpoint, where we would be observing oppositely directed gravitational fields and an absence of relative motion of the elevators, we would have to conclude that those fields are due to the presence of a local mass and not to acceleration relative to the homogeneous large-scale matter distribution, even in the absence of negative-mass bodies in the elevators. In fact, even in a single elevator standing still on the surface of a small planet, freely falling positive-mass particles would have a tendency to slightly converge toward one another, therefore betraying the fact that the observed acceleration is an effect of the presence of a nearby mass attracting the particles toward its center. Yet we do not consider the equivalence principle to be violated under such conditions.

What I'm suggesting, therefore, is that, instead of assuming that the equivalence of gravitation and acceleration applies only locally, we have to recognize that it really applies only for a single elementary particle, which would be the most localized physical system we may consider. If we assume that no two such particles can be exactly superposed in an elementary volume of space (which ultimately may be true for bosons just as for fermions if there is a maximum value of energy associated with the Planck scale) we could say that the hypothesis that the equivalence of acceleration and gravitation applies merely within a local free-fall reference system is equivalent to the assumption that the equivalence principle always applies only for one single elementary particle. But then such a particle could have either positive or negative mass and the equivalence principle could be considered to apply not merely to one particle at once, but to one particle with one mass or energy sign at once, which would be a simple generalization of the discussed hypothesis and as such, should not raise any further issue (of the kind I have considered so far). For one elementary particle, with one energy sign, there would never be a difference between acceleration and a gravitational field. It is only when we consider two or more particles of *any* mass sign together, or more precisely in relation to one another, in the presence of a gravitational field attributable to a local matter inhomogeneity (when there is no compensation between the gravitational fields attributable to the local positive- and

negative-energy matter distributions) that we can tell the difference between acceleration relative to the large-scale matter distribution and such a gravitational field, but this may be assumed irrelevant when we are considering that no two particles (especially two opposite-mass particles) can actually be found in the exact same position at the same time.

It is generally recognized, however, that what makes gravitation different from other interactions is the fact that the motion of bodies in a gravitational field does not depend on the physical properties of those bodies (when no other force field is present). But even though this characteristic would appear to be violated in the presence of negative-energy matter obeying the consistency conditions I have identified, this does not make gravitation any less distinct. Indeed, in the context of the previously discussed viewpoint where it is the direction of the gravitational field attributable to a given matter distribution which varies upon a reversal of the mass of the particle submitted to it (which would actually be considered positive definite), the equivalence principle would merely be relativized by the presence of such negative-energy matter, because the difference between the motion of positive-energy bodies and that of negative-energy bodies would actually be a consequence of the different measures of spacetime curvature which (as I will explain later) can be associated with those two measures of the Newtonian gravitational field. But in such a situation it appears natural to expect that opposite-mass bodies should not be restricted to share the same local inertial reference systems, because, in fact, they do not even evolve in the same space, but in spaces characterized by different metric properties.

Thus, the fact that the gravitational field can be conceived in such an observer-dependent way means that, in the case of gravitation, it is not the reaction that varies when the ‘charge’ is reversed, but the field itself (to which is associated a given spacetime curvature). It is still true, therefore, that, in any given situation, all bodies (sharing the same measure of gravitational field) follow the same motion (acceleration does not depend on the detailed characteristics of the bodies experiencing the same gravitational field). The equivalence principle can thus be assumed to still be valid in the presence of negative-energy matter, only it would apply separately for positive- and negative-energy bodies (just as it applies separately for separate portions of space), because each of those two kinds of matter particle is to be attributed its own free-fall reference system defined in relation to its mass sign. Therefore, all particles with the same energy sign still share the same local inertial reference system and this is all that is truly required for a general-relativistic

gravitational field theory to apply.

2.6 An effect of voids in the matter distribution

It is sometimes recognized that there is a kind of equivalence between the presence of a void in an otherwise uniform matter distribution and what would be the presumed effect of the presence of gravitationally repelling matter present in a quantity and with a distribution equivalent to that of the missing matter. In the context of an expanding universe, we would, in effect, observe underdense regions of the cosmos to be producing a local acceleration of the rate of expansion, while overdense regions would produce a local deceleration of it. The acceleration observed in the case of underdense regions would have all the characteristics of a gravitational repulsion originating from those regions, which would force the matter still remaining inside their volume to migrate to the periphery of what would become the observed voids in the matter distribution [22]. The same effect would also cause nearby underdense regions to merge into even larger spherical voids, as if they were attracted to one another by the force of gravity. This is what all authors who have considered the issue agree must occur when underdense regions form in an expanding universe. Thus, in this particular case, it seems that the gravitationally repelling matter formations would actually be submitted to mutual gravitational attraction with similar formations, even while they would repel oppositely configured formations consisting of overdense regions and would presumably also be repelled by them.

But it is usually considered that there is nothing more than an accidental analogy between the case of those matter formations and any gravitationally-repulsive matter, because if the effect occurs as described above then, according to the traditional understanding, such gravitationally-repulsive voids would need to have not only negative gravitational mass, but also positive inertial mass [23] and as everyone ‘knows’, this kind of negative mass is forbidden by the equivalence principle and relativity theory, which require the equality of gravitational and inertial masses. Thus, what we would *observe* to be happening is not what most people would consider should occur if we were actually dealing with gravitationally-repulsive matter. Indeed, as I previously explained, what is usually assumed is that gravitational repulsion is a kind

of definite and invariable property of matter of some type and that this kind of matter would therefore itself also be repelled by matter of the same type. This is usually assumed to be the unavoidable consequence of attributing a negative inertial mass to negative-energy matter. But, given the previous discussion and the insights I provided concerning what should be a consistent concept of negative-mass or negative-energy matter, it should be clear that we would not be justified to argue that the observed phenomenon involving voids in a uniform matter distribution does not replicate the behavior we should expect of negative-mass matter. In fact, from my viewpoint it rather seems that the described interaction between overdense and underdense regions of an expanding universe would be exactly that which we should expect to occur if positive and negative masses were actually involved. Therefore, we cannot so easily reject the possibility that the discussed phenomenon is actually telling us something important about the nature of negative-energy matter.

I do believe that there is actually more than a valid analogy between voids in a uniform positive-energy matter distribution and gravitationally-repulsive matter and that there is something very profound which we need to understand concerning the phenomenon described here. Indeed, I think that the discussed equivalence should not be restricted to the case of expanding matter, but must be considered valid even in a local context, where the rate of universal expansion is a negligible factor. But if the gravitational dynamics of voids in a homogeneous positive-energy matter distribution actually reflects that which we should expect of a phenomenon involving gravitationally-repulsive negative-energy matter, then it may suggest an interpretation of negative-energy matter which would have to do with an absence of positive energy of some kind. It must first be explained, however, why it is that we may actually be allowed to consider that the equivalence discussed above is valid exactly and constitutes a very general feature of the gravitational interaction, despite the objections which might be raised against that possibility.

Basically, what we may object, concerning the idea that the presence of a void in a uniform positive-energy matter distribution could be equivalent to the presence of an excess of negative-energy matter, is that it is usually assumed that there can be no net gravitational force inside a spherical void in a uniform matter distribution that would be attributable to matter outside the void; a conclusion that seems to be supported by Birkhoff's theorem [24]. What Birkhoff's theorem implies is that there can be no net gravitational

force on matter located inside any spherically symmetric region in a globally uniform matter distribution from matter located outside that region. This is usually assumed to imply that there cannot be any net gravitational force inside a spherical void in a uniform matter distribution, given precisely that there is no matter inside the void, while it would appear that the surrounding matter itself cannot exert such a gravitational force. Thus, it seems that in the absence of any matter inside a spherical region, there can be no local gravitational field on the boundary of that region, as any gravitational acceleration could only be attributed to matter located inside the region considered, while there would be no matter inside that region.

The influence of voids on the local rate of acceleration of cosmic expansion which was discussed above would thus merely be a result of the fact that the rate of growth of the distance between two galaxies located on the boundary of such a void actually depends on the density of matter *inside* the void and given that this density would be lower than the average, then the rate of growth of the distance, or the local rate of expansion would be larger in proportion with the amount of matter missing inside the void. But that does not mean that it is usually assumed that there would actually be a repulsive gravitational field on the surface of the void. In fact there appears to be some confusion surrounding the issue discussed here, as some authors recognize that there cannot be an equilibrium of gravitational forces in the presence of a void in the cosmic matter distribution and yet they fail to recognize that this may actually give rise to repulsive gravitational fields for the surrounding positive-energy matter, probably because they assume that the effect of the noted disequilibrium would be that which is observed to affect the local rate of expansion, while actually this is a distinct (but not entirely unrelated) effect associated merely with cosmic expansion. But what I believe must be recognized is that there would, in effect, be gravitational repulsion in the presence of an underdensity in an otherwise uniform matter distribution, not only at the boundary of the surface, but everywhere inside the void, with a net force that would decrease as we approach the center of the void, where it would have a null value. This situation would then clearly be different from that we would have in the case of a hollow sphere of finite size, inside of which the Newtonian gravitational field should indeed be zero everywhere.

It must, in effect, be understood that contrarily to what is usually believed, Birkhoff's theorem does not forbid this conclusion, because the decisive condition entering this theorem is that of spherical symmetry, which

would actually be obeyed if we were considering a hollow sphere or a universe that was spherically symmetric around *any* point on any scale, but which, I suggest, would fail, locally, for a universe with an actual void in its matter distribution. Indeed, the case of a homogeneous and isotropic universe is equivalent to that of a sphere of finite size only when the universe is considered on the scale at which its matter is uniformly distributed and no significant void is present, which explains why Birkhoff's theorem (which is a necessary element of current cosmological models) is observed to apply on a cosmological scale. But I think that it would only be in the case of a spherical region centered on an actual sphere of matter of finite size, located within an otherwise empty universe, that the theorem discussed here would actually remain valid regardless of the distribution of matter inside the spherical region, because only in such a case would we be dealing with a spherical symmetry that is not dependent on the position of the observer. What we usually fail to recognize is that, the fact that the matter distribution in the universe would be symmetric around *any* location in the absence of a void in its homogeneous and isotropic matter distribution, means that the presence of a void would necessarily alter the equilibrium of forces around that void.

It is clear, indeed, that in the presence of a uniform matter distribution extending throughout the universe, an equilibrium exists locally between the sum of forces attributable to the interaction of a freely falling body with all the matter in the universe and therefore the removal of a certain quantity of matter in a region of finite volume must have an effect that would be the opposite of that which we would otherwise attribute to the matter that is missing in this region of the universe (in the absence of a void). This should be expected to occur due to the fact that the removal of a certain amount of positive-energy matter, to create a void, would eliminate the attractive gravitational force which would otherwise be exerted on positive-energy matter by the matter in the void and given that there was no net force before the creation of the void, then the other forces which are still present would give rise to a gravitational acceleration directed away from the void and of similar magnitude to that which would have been produced by the matter that filled the void. Thus, for positive-energy matter, there would appear to be a repulsive gravitational force originating from the presence of a void in such a uniform matter distribution, which would actually be the consequence of an uncompensated gravitational attraction attributable to the positive-energy matter outside the void. But this is a valid conclusion only when we recognize that Birkhoff's theorem is not valid in the sense it is usually assumed to

be and that the case of a spherical distribution of matter of finite size, with a central cavity, is *not* equivalent to the case of a void in a uniform cosmic matter distribution.

What must be understood is that if, in the case of a hollow sphere of finite size, the subtraction of matter to create the cavity does not result in a net force originating from the matter surrounding the cavity *that is part of the sphere* this does not mean that it would also be the case that there would be no acceleration inside the cavity resulting from the gravitational interaction with *all* the matter that is present in the universe (unless it was actually assumed that the universe is empty except for the presence of the sphere). What is wrong, therefore, is the idea that when we are considering a spherical region of the universe, the rest of the universe surrounding that region can be considered as a hollow sphere simply on the basis of the fact that, according to the cosmological principle, matter is distributed uniformly in all directions. In fact, such a spherical region in a uniform matter distribution would be free of uncompensated external forces only if it was itself filled with matter as uniformly distributed as the matter found outside the region (which is actually verified on a cosmological scale in our universe), because it is only in such a case that the spherical symmetry would apply to any point inside the spherical region. Again, it must be noted that, in this context, the fact that the concept of the hollow sphere is nevertheless appropriate to describe the dynamics of the universe on the largest scale is due merely to the fact that we do not actually consider the case where spherical voids are present in the matter distribution, but really the case of a uniformly filled matter distribution for which no spherical regions devoid of matter are present on the particular scale that is considered (as a requirement of the cosmological principle).

It must be clear that I'm not suggesting that there would be uncompensated gravitational forces in the case of the finite-size hollow sphere itself (if it was located in an empty universe for example). In fact, the problem here has to do again with the fact that we fail to apply the requirement of relational definition of physical properties when we are dealing with the resultant effect of the gravitational forces attributable to the universe as a whole. Indeed, from the traditional viewpoint, when we are dealing with a chosen spherical region of the universe, we are implicitly assuming that the surrounding matter which may influence the particles located inside that region (through the gravitational interaction, even if there is no net force) is spherically distributed around the center of the spherical region consid-

ered, as if the location of the center of mass of the universe was an intrinsic invariable feature of the whole configuration. But the center of a matter distribution, in a physical universe without boundary, is not an absolute feature (as would be the case for a hollow sphere), it must rather be defined in a relational manner as any other property, if we are to be able to determine the consequences on a given object of being located in such a position. When we are dealing with the matter distribution in a universe without spatial boundary, in which the local inertial reference systems are determined by the entire matter distribution (following Mach's principle), the true center of mass, defined in terms of the influences exerted on a given body, is *always* located right at the position where that body is to be found, wherever this position may be in the matter distribution.

Thus, a particle located *at the center* of a void in a uniform matter distribution could actually be considered to be in the situation of a particle in a hollow sphere, because for this particle the whole sphere of influence of the universe is centered on the void (in this situation the surrounding matter actually is a hollow sphere centered on the particle's position). Therefore, such a particle would feel no uncompensated gravitational force from the whole universe, as required. But if this particle moves to one side or another in the void, the matter distribution influencing the particle in its new position would be centered on the new position and this means that the void in the previous hollow sphere is shifted to the opposite side, just as the sphere itself is shifted in the direction of the particle's new position. The symmetry of the initial configuration would therefore no longer be present and the equilibrium of forces would no longer apply. In the new configuration, a whole layer of matter must be 'removed' on one side of the external surface of the imaginary hollow sphere (in the direction opposite the particle's displacement) and added on the other (this is easier to visualize in a closed universe) which, given the distances involved, means that an enormous amount of matter has changed position from the viewpoint of the particle. It must therefore be recognized that, in the final configuration, the void in the imaginary sphere is no longer centered on the center of mass of the sphere, but is actually located away from the center of the sphere. As a consequence, the spherical symmetry, from which depended the conclusion that there would be no net gravitational force inside the sphere, is no longer to be found in the final configuration experienced by the particle and therefore it must be expected that there would be a net gravitational force on the particle and an acceleration relative to the matter distribution.

It is important to understand that, however large you consider the imaginary sphere encompassing the matter distribution (the size of the universe) to be when dealing with the effects of the gravitational interaction with the whole universe, if the center of the sphere is shifted to one side there would be a non-negligible effect from the displacement of its center of mass. This is true even if the distance to the periphery of the sphere (where the changes occur) is very large and the strength of the gravitational interaction decreases with the square of the distance, because the larger the distances (the larger the sphere) considered, the larger the quantity of matter that is shifted from one side to the other and thus the larger the changes involved in the local gravitational field. We should not be surprised, then, that, even the retarded interaction with matter so distant could have an effect similar in magnitude to the effect that would be exerted by the matter missing from a void located near some particle experiencing those forces. If the center of mass of the universe is always located at the position of the particle experiencing the gravitational effects of all the infinitesimal elements of matter in this universe, then the local effect of the absence of gravitational attraction from that portion of matter which is missing when a nearby void is present in the positive-energy matter distribution would necessarily result in a net force which, for positive-energy matter, would be directed away from the void and which would arise from the gravitational attraction of that portion of positive-energy matter located outside the void in this same direction. But such a force would be completely equivalent to a repulsive gravitational force arising from the void itself.

The fact that, from a practical viewpoint, the formation of a local void in a uniform positive-energy matter distribution would actually have to occur through the ejection of positive-energy matter outside the region that is to become the void and therefore would necessarily produce a compensating overdensity of negative-energy matter in the region surrounding the void would not forbid the existence of a net repulsive force on positive-energy matter inside the void, even though it does, in effect, mean that there would be no resulting force on matter located some distance away from the void. If we consider, for example, the ideal situation of a spherical void produced through the creation of a surrounding spherical shell of positive-energy matter at higher than average density, then, as long as a positive-energy particle is located outside this shell, it would feel no net force, because any reduction of attractive force from the void would be compensated by an increased attractive force arising from the presence of the shell. But as soon as the

particle would enter the shell it would begin to experience the equivalent gravitational repulsion, because the outer layers of the shell would no longer provide any net force on the particle, while the void, for its part, would still exert its net effect, because the equivalent repulsive force it produces is attributable to *all* the surrounding matter (whose distribution is centered on the position of the particle) and not just to the spherical shell. Thus, the case of the particle which experiences no gravitational force *at the center* of a void in a uniform matter distribution is merely a particular case of the more general description according to which there is actually a net force everywhere inside the void, except at the exact location of its center, as would be the case if we were considering the gravitational attraction existing inside an isolated sphere filled with matter (like a planet or a spherical gas cloud) present in an otherwise empty universe. This is an important result which will have decisive consequences for a consistent description of the nature and properties of negative-energy matter.

Concerning the conclusion just reached, it is important to note that even if, under certain circumstances, there may be an equivalence between an imbalance in the sum of gravitational attractions attributable to all the positive-energy matter elements in the universe and what would appear to be a gravitational repulsion exerted on a positive-energy body, we are nevertheless always dealing with gravitational attraction. Indeed, there is no question that it is the gravitational attraction of positive-energy matter that is responsible for the apparent gravitational repulsion which would be exerted on a positive-energy body by a void in the otherwise uniform positive-energy matter distribution. It is clearly as a consequence of the fact that positive-energy matter is missing in the direction where the void is located, while the matter present in the opposite direction still exerts its gravitational pull, that there exists a net force directed away from the void.

Thus, what looks like a gravitational repulsion exerted in a given direction by some matter configuration and which could, from a certain viewpoint, be equivalent to it, would actually be the product of a gravitational attraction arising from an absence of matter exerting a compensating attraction in the opposite direction. This is particularly significant in the context where local inertial reference systems are to be considered as always arising from a perturbation of the equilibrium of large-scale inertial gravitational forces by the gravitational forces attributable to local matter concentrations, as I have emphasized in the preceding section. Yet the fact that we are here

dealing only with gravitational attraction does not rule out the validity of the analogy which may exist, from a classical viewpoint, between the presence of true gravitationally-repulsive, negative-energy matter and an absence of positive energy of some sort. In fact, it rather seems that what allows an interpretation of negative-energy matter as being equivalent to an absence of positive energy to be valid as a general feature of classical gravitation theory is the possibility that always exists (not only in the case of voids in a uniform matter distribution) of attributing an apparent gravitational repulsion to uncompensated gravitational attraction.

To explain what motivates that conclusion, it is necessary to recall the previous discussion concerning the occurrence of negative energy in certain experiments described using ordinary quantum field theory. There, I pointed out that the absence of some positive energy states from the vacuum in certain limited regions of space (between the plates of two parallel mirrors for example) can actually give rise to a vacuum with negative energy density in the volume considered, because removing positive energy from a vacuum state whose energy is already minimum is like decreasing the energy below its zero point, into negative territory. The fact that the vacuum is known to have only a very small average energy density should not be considered an obstacle to the occurrence of large negative energies in such a way, because, as I will explain later in this chapter and in section 4.2, this small energy density appears to be the outcome of very large (actually maximum) but (mostly) compensating opposite-energy contributions, which could be reduced to an arbitrarily large extent by the conditions which are responsible for locally decreasing (under particular circumstances) the energy of the vacuum below the equilibrium point. But if we may, in effect, attribute a negative energy to certain configurations in which positive energy states are missing from the vacuum, then there is no reason why we could not consider that negative energy states in general are equivalent, in some ways, to a local absence of positive energy from the vacuum, if from a phenomenological viewpoint there is no distinction between those two situations.

I must again mention, in this regard, that many authors have expressed doubts concerning the validity of the concept that energy should exist in the vacuum that would be the outcome of the presence of zero-point fluctuations involving virtual particles and have suggested that there may be nothing real with the processes so described outside of the context where they are occurring as part of otherwise real processes involving ‘real’ particles. But I think that what really motivates this mistrust is precisely the fact that the

existence of those processes would imply the reality of negative energy states, because it is no secret that, for most physicists, the theoretical possibility of the existence of negative energy states is not well-viewed. However, I believe that this aversion is merely a consequence of the fact that the traditional concept of negative-energy matter is, in effect, not viable and that it has not yet been realized that a better description of negative-energy matter is possible and even necessary, as I emphasized before.

In any case, the idea that virtual processes would only occur as part of otherwise real processes, thus explaining why we must nevertheless consider the effects of such fluctuations when calculating transition probabilities, is meaningless, because, in a given universe, anything that occurs is related (directly or indirectly) to everything else and even in empty space, far from any ‘real’ matter, the virtual processes of particle creation and annihilation characteristic of the quantum vacuum would occur as an integral part of the surrounding real processes to which they are causally related, as a consequence of their common origin in the Big Bang. In fact, I will explain in section 4.9 why those considerations actually constitute a decisive element of a consistent cosmological theory, even aside from the issue of vacuum energy. Therefore, the argument that the negative energy states predicted to occur in the vacuum under the right conditions are not real, because our description of the vacuum is itself not appropriate in general, cannot be retained. Also, the fact that it has been confirmed that the cosmological constant is not absolutely null is a strong motive to conclude that the rejection of the reality of vacuum fluctuations, as essential aspects of our description of empty space, is not vindicated from the viewpoint of observations and therefore that negative energy states are a real possibility.

I have already explained why we should expect to observe mutual gravitational attraction between two bodies with the same sign of energy and gravitational repulsion between opposite-energy bodies. But on the basis of my conclusion concerning the nature of the gravitational force on a positive-energy body that would be attributable to voids in a uniform positive-energy matter distribution, we now also have the possibility to assert what would be the effects of missing positive energy from the vacuum. Indeed, given that, in the absence of local perturbations at least, the vacuum is to be conceived as involving a constant and uniform density of energy on the largest scale, any local variation in its density must have consequences similar, from a gravitational viewpoint, to those of inhomogeneities in a uniform matter distribution. It therefore appears that if the presence of voids in an oth-

erwise homogeneous positive-energy matter distribution does, in effect, produce an equivalent gravitational repulsion on positive-energy bodies, then the absence of positive vacuum energy in localized regions should itself also exert an equivalent gravitational repulsion on the surrounding positive-energy matter. This would occur as a result of the fact that an absence of positive energy from a region of the vacuum would result in an uncompensated gravitational attraction from the surrounding positive portion of vacuum energy that would pull positive-energy matter *away* from the region where the energy is missing. From that viewpoint, we can deduce that the physical properties (related to the gravitational interaction) that we should expect to be associated with missing positive vacuum energy are the same properties which I explained we should expect to be associated with the presence of negative-action matter, which confirms that, from a phenomenological viewpoint, negative-energy matter is gravitationally equivalent to an absence of positive energy from the vacuum.

Given this equivalence between the presence of negative energy matter and an absence of positive energy from the vacuum, it follows that if states of negative vacuum energy are allowed by current theories, then we must conclude that negative-energy matter is itself allowed to exist and may not always be constrained by the limitations which are observed to apply in the currently considered experiments where it occurs merely as a consequence of a suppression of positive energy from the vacuum that is attributable to singular configurations of matter with otherwise positive energy. It must be recognized, however, that if the presence of negative-energy matter in a region of space is equivalent, for positive-energy matter, to an absence of positive energy from the vacuum, this is simply because in general, for an equilibrium state of any kind, the presence of a negative contribution is equivalent to the absence of a positive contribution of the same magnitude and it just happens that the vacuum is a physical system that appears to arise from precisely such an equilibrium state (as I will explain later). But we must remember that a void in a uniform positive-energy matter distribution (not involving the vacuum) is physically different from a local absence of positive vacuum energy, even if in both of those cases the effects are equivalent (from a gravitational viewpoint) to the presence of an excess of matter of opposite energy sign, because, in the first case, we are dealing with an *absence* of matter of positive energy sign, while, in the latter case, we are actually dealing with the *presence* of matter (of opposite energy sign).

At this point it is important to mention that there would occur a phenomenon of gravitational repulsion similar to that described above, but which would apply from the viewpoint of negative-energy matter in the presence of voids in a negative-energy matter distribution or in the negative-energy portion of the vacuum. Indeed, using the same logic that allowed me to derive the consequences of the presence of a void in a uniform distribution of positive energy, it is possible to deduce that the absence of negative energy from an otherwise homogeneous matter distribution would actually be equivalent, from a gravitational viewpoint, to the presence of a concentration of positive-energy matter. One assumption that will be crucial for my derivation of the modified general-relativistic gravitational field equations is that the equivalence described here is valid both ways and that positive-energy matter can always be considered to actually consist of voids in the *negative-energy* portion of the vacuum, which makes the whole situation symmetrical in a way that does not even depend on the viewpoint of the observer. It must be clear, however, that I'm not suggesting that positive-energy matter is equivalent to voids in a filled distribution of negative-energy matter, even if I do suggest that we must assume that an absence of negative-energy matter from an otherwise uniform distribution of such matter would indeed have effects similar (from a gravitational viewpoint) to those attributable to the presence of positive-energy matter. I must emphasize, once again, that a void in a uniform *matter* distribution remains clearly distinct from a void in the uniform distribution of *vacuum* energy. This means that my proposal is distinct from Dirac's failed, hole theory (proposed as an attempt to solve the negative energy problem), in particular because what I'm suggesting is that *all* positive-energy matter particles (and not just antimatter particles) are actually equivalent to voids in the negative-energy portion of the vacuum, rather than in a filled continuum of negative-energy matter.

What Dirac proposed, in effect, is that all negative energy states are already occupied, so that positive-energy fermions, at least, should not be expected to make transitions to those negative energy states. But even if the existence of such a filled, uniform continuum of negative-energy matter was to have no effect on positive-energy matter (perhaps due to its uniformity), the fact that, from my viewpoint, there would be no reason to assume that positive energy states are not completely filled in the same way means that this hypothesis would not agree with observations. Indeed, it is not possible to assume, in a theory that respects the requirement of a purely relational definition of the sign of energy, that positive-energy antiparticles are merely voids

in a completely filled negative-energy matter continuum, as Dirac proposed, without also assuming that negative-energy antiparticles would be voids in a completely filled positive-energy matter continuum. But, given that positive energy states are obviously not all occupied by matter particles, it appears that this requirement cannot be satisfied. We may then instead assume that *all* positive-energy particles are voids in a filled negative-energy matter continuum, but again in such a case we would have no reason not to assume that all negative-energy particles are also voids in a filled positive-energy matter continuum. The problem, however, is that it seems impossible to assume that we could have a completely filled distribution of negative-energy matter and at the same time a completely filled distribution of positive-energy matter if negative-energy matter is to also consist of voids in a filled distribution of positive-energy matter, because so many voids in the positive-energy matter distribution as would be necessary to describe the filled negative-energy matter distribution would leave no possibility for the positive-energy matter distribution to itself be nearly completely filled.

What cannot be assumed, therefore, is that negative energy states are completely filled and positive-energy particles consist of voids in this negative distribution of energy, while positive energy states would also be completely filled and negative-energy particles would consist of voids in this positive distribution of energy, because those two possibilities are mutually exclusive (cannot occur together). But while it may perhaps appear appropriate from an observational viewpoint to assume that we simply have a filled negative-energy matter continuum combined with a nearly empty distribution of positive-energy matter, there would also be problems with such a proposal. Indeed, what reason would we have not to assume that it is only the positive-energy matter distribution that is filled (even though this assumption would clearly contradict observations)? The problem is that we cannot, in effect, postulate that both positive- and negative-energy matter are voids in their respective opposite-energy matter distributions if we also postulate that there is no absolute (non-relational) difference between positive- and negative-energy matter. In other words, it is not possible to assume symmetry under exchange of positive- and negative-energy particles if matter of a given energy sign is to be conceived as voids in the matter distribution of opposite energy sign and this, simply because matter cannot be at once present and absent. The truth is that any description of matter or antimatter as voids in a matter distribution of opposite energy sign would require giving preferred status to negative-energy matter as being the matter

whose distribution is completely filled (because obviously the positive-energy matter distribution, at least, is not completely filled) and this would break the requirement that only differences in the energy sign of particles are to be conceived as physically significant.

What must be clear, therefore, is that if we were to make use of such a description, we would allow the identification of a preferred sign of energy as being that which would be associated with the filled matter distribution, while, from a theoretical viewpoint, that should be considered impossible. A theory of matter particles as voids in a uniform matter distribution would, in effect, imply that the requirement of symmetry under exchange of positive- and negative-energy matter is violated in a way that cannot be allowed if the sign of energy is to be conceived as a relationally defined physical property. Thus, it must be recognized that while it would appear possible to explain the presence of matter with a given energy sign as being equivalent to missing *vacuum* energy with an opposite energy sign, it is nevertheless forbidden to consider that the presence of matter with a given energy sign could be explained as resulting from the presence of voids in a *matter* distribution of opposite energy sign, even if there does exist a phenomenological equivalence (from the viewpoint of the gravitational interaction) between the effects of missing positive or negative vacuum energy and those attributable to a local absence of matter from a homogeneous distribution with the same sign of energy, because, again, those are two distinct phenomena.

The contradiction which would occur if we were to assume that positive-energy particles are voids in a filled uniform distribution of negative-energy matter, while negative-energy particles are voids in a filled uniform distribution of positive-energy matter is that, in the first instants of the Big Bang at least, a lot of particles of both energy signs would be required to fill the matter distributions and at the same time a limited number of particles of both energy signs could be present due to the presence of all the voids attributable to the presence of the nearly filled matter distributions themselves. By contrast, when matter particles are merely equivalent (from a gravitational viewpoint) to a local absence of vacuum energy of opposite sign, it becomes possible for both positive- and negative-energy particles to actually exist as real observable particles independently from the presence of one another. Thus, if the voids in the negative-energy portion of the vacuum, which I assume to be equivalent to the presence of positive-energy matter, are not equivalent to voids in a hypothetical filled distribution of negative-energy matter it is simply because, in fact, voids in the vacuum cannot be

equivalent to an absence of voids in the vacuum.

I may add that, from the viewpoint of a consistent interpretation of negative-energy matter, there would also be a problem with Dirac's original proposal that a void in the filled negative energy continuum could be created, along with a positive-energy particle (as would a particle-antiparticle pair), when photons provide enough energy to raise a negative-energy particle to a positive energy level. Indeed, as I mentioned before and for reasons I will explain in section 2.8, a consistent theory of negative-energy matter would require that negative-energy matter be dark, which means that there would be no electromagnetic interactions between opposite-energy particles and therefore a positive-energy photon could not even interact with a negative-energy electron to provide it with the required positive energy. Thus, even if we insist on assuming the existence of a filled negative energy continuum, we could not use this hypothesis to explain the existence of antimatter.

It is essential to understand, therefore, that the situation we would have if all negative energy states were filled is different from that we would have when dealing with a vacuum in which there would be a very large negative contribution to the average energy density of zero-point fluctuations. Indeed, in contrast with the vacuum, a negative-energy matter distribution which would be filled at one particular epoch would no longer be filled at a later time, given that space is expanding. This is reflected in the fact that the uniform portion of vacuum energy obeys an equation of state which is different from that of a homogeneous matter distribution. Also, even if there is a large negative contribution to the energy of the fluctuating vacuum, there is no reason to expect that it gives rise to a situation similar to that which would occur if space was filled with negative-energy matter, because in such a case there must also be a large positive contribution to the energy of empty space (the motives behind this conclusion will be clarified in section 4.2). A space filled with positive- or negative-energy matter would be as different from the true vacuum as the primordial soup which existed in the first instants of the Big Bang is different from the space nearly devoid of matter particles that currently exists between galaxies. Thus, if a theory of voids is to have any relevance in a gravitational context, it must involve a description of matter of any energy sign as consisting of voids in the opposite-energy portion of the vacuum, so that the presence of matter with a given energy sign does not imply an absence of *matter* with opposite energy sign.

When the energy distribution in which the voids equivalent to the presence of positive-energy matter occur is the negative-energy portion of the

vacuum it becomes possible to assume the presence of arbitrarily high or arbitrarily low densities of matter of both energy signs, all at once, in the same region of space, because the presence of matter of one energy sign in a given location does not preclude the presence of matter with an opposite energy sign in the same location (at least when the matter distributions are smooth enough). Thus, we do not need to assume the presence, at all times, of a nearly filled negative-energy matter continuum combined with a distribution of positive-energy matter of arbitrarily low density, which would otherwise be the only (perhaps) observationally acceptable configuration, but which would also have allowed to establish an absolute (non-relational) distinction between positive- and negative-energy matter, as I just explained. But what makes the vacuum particularly suitable for accommodating the above proposed description of matter as consisting of voids in some uniform energy distribution is the fact that we are actually allowed to assume that there are both positive and negative contributions to vacuum energy density, even as arise from the presence of otherwise identical virtual particles. We can, therefore, expect a certain level of compensation between the gravitational effects of those two contributions that may give rise to an arbitrarily small residual value for the cosmological constant. Indeed, in sections 4.2 and 4.5 I will explain that one of the consequences of the assumption that there exists a distinct component to the energy of the vacuum, arising from the presence of those virtual particles that directly interact (other than through the gravitational force) only with negative-energy matter, is that the natural value of the cosmological constant which we can expect to observe is zero, even though this value can be altered so as to compensate any imbalance that might develop between the scale factors experienced by positive- and negative-energy observers, which are required by the weak anthropic principle to be indistinguishable, in the first instants of the Big Bang.

When it is understood that all positive- and negative-energy particles are actually equivalent to voids in their respective opposite-energy portions of the *vacuum*, as I propose, then it also follows that the unsatisfactory categorical distinction between matter and vacuum becomes meaningless. This is because, in such a context, all matter can actually be considered to consist in a particular aspect of zero-point vacuum fluctuations. It is by building on this insight that I will be able to provide a unified and totally symmetric description of the gravitational dynamics of positive- and negative-energy matter, according to which the measure of energy of matter is significant merely in relation to an energy scale associated with objective properties of

the vacuum. I was able to obtain those results only at a relatively late stage of my reflection, because I had initially assumed that only the nearly vanishing total energy density of the vacuum could have an influence on matter of any energy sign and that the positive and negative contributions to vacuum energy could not be considered independently from one another. But once I realize the inappropriateness of this hypothesis, the above discussed results emerged as clearly unavoidable and extremely significant. The notion that both positive- and negative-energy particles actually consist of voids in their respective opposite-energy portions of the vacuum, therefore, appears to be the ultimate embodiment of the requirement of a relational definition of all physical properties, understood as a basic consistency condition that must apply to any physical theory.

Concerning the effects which I'm suggesting should be attributed to energy missing either from a homogeneous matter distribution or the vacuum, we may ask to what extent a void may actually be considered as physically significant, in the sense of being merely an anomaly in an otherwise uniform distribution of matter or energy. If we examine the situation carefully it becomes clear, in effect, that, given that, for both matter and vacuum, it must be the surrounding energy that exerts the *outward* directed gravitational pull that would be experienced as a gravitational repulsion, then it follows that, as we consider voids of larger sizes, there may come a point when there would be no matter left outside the void to produce the uncompensated attraction that must exist to produce the equivalent repulsion. Normally, this is not an issue, as any void that forms in a matter distribution that is arbitrarily smooth initially (and this appears to be a necessary feature of our universe at the Big Bang, as I will explain in chapter 4) will necessarily involve the creation of a surplus of matter in its surroundings, which, for a remote observer (away from the void), would have the effect of compensating the equivalent force arising from the presence of the void itself. Such voids, regardless of how large they may become, would, therefore, leave the universe at large in a state equivalent to that of a uniform matter distribution, which would be allowed to continue to exert its influence on the matter that is present inside or around the voids.

But if we are to consider the equivalence between missing positive vacuum energy and the presence of negative-energy matter to be generally valid, then the presence of a uniform negative-energy matter distribution would imply the existence of a void in the positive-energy portion of the vacuum which

would actually extend to the whole universe. This void would have been present in the vacuum from the very beginning of the universe's history and would not have developed through the production of some inhomogeneity. In such a case we would no longer be able to assume the existence of an uncompensated gravitational pull on positive-energy bodies from the surrounding positive vacuum energy, because indeed there would be no surrounding vacuum energy with higher positive density to generate the attraction. Under such conditions, therefore, I'm allowed to conclude that no outward directed gravitational force, which would be equivalent to gravitational repulsion, could exist, that would be experienced by positive-energy bodies.

Now, given that I will later argue that the equivalent gravitational repulsion exerted on positive-energy matter by voids in the positive-energy portion of the vacuum actually constitutes the only form of gravitational interaction between this matter and negative-energy matter, it would appear that the preceding conclusion imposes very strong limitations on such an interaction. Indeed, it transpires that the absence of equivalent gravitational repulsion on positive-energy matter from a completely homogeneous negative-energy matter distribution, is a very general and unavoidable feature of the description of the gravitational interaction between positive- and negative-energy matter. This is because such a limitation would also be verified in the case of a uniform distribution of positive-energy matter from the viewpoint of negative-energy bodies, if the gravitational repulsion exerted on those objects by positive-energy matter can be attributed to an absence of *negative* energy from the vacuum.

Thus, if opposite-energy bodies can be shown to interact only through their respective same-energy-sign vacuums, we would be allowed to conclude that negative-energy matter interacts with positive-energy matter only in the presence of inhomogeneities in any of the two matter distributions. But given that only an inhomogeneity that develops over the initially-smooth negative-energy matter distribution (if we may suppose that negative-energy matter was as homogeneously distributed as positive-energy matter in the primordial universe) can contribute to the gravitational dynamics of positive-energy matter and given that the formation of such an inhomogeneity would involve the formation of a compensating one, involving an opposite variation of density in the surroundings of the first, we must then conclude that the presence of an average density of negative-energy matter has absolutely no effect (at least from a gravitational viewpoint) on the gravitational dynamics of positive-action matter (and vice versa). This would mean, in particular,

that the rate of universal expansion experienced by positive-energy observers cannot be influenced by the presence of negative-energy matter and similarly that the rate of expansion experienced by negative-energy observers is not affected by the presence of positive-energy matter. This is a very significant result, which will have an impact on many aspects of cosmology theory and whose implications will be developed in chapter 4.

I may add that the conclusion discussed here is the one on which is founded the hypothesis, discussed in section 2.5, which allowed a relational description of the phenomenon of inertia. There, I explained that if both the large-scale positive- and negative-energy matter distributions were to exert an influence on positive-energy bodies, then the hypothesis that accelerated motion is relative would be invalidated in the presence of negative-energy matter on a cosmological scale. Indeed, under such circumstances there would be compensating imbalances in the sum of gravitational forces (to which we would try to attribute the resultant inertial force) arising from the acceleration of a positive-mass body relative to the two opposite-energy matter distributions, whose average states of motion should correspond with one another on the largest scale. But if only matter with a positive energy sign has a gravitational effect on positive-energy bodies on the cosmological scale, then the global inertial reference system experienced by a positive-energy body could actually be determined by the average state of motion of positive-energy matter, given that the inertial force exerted on such a body would result only from its gravitational interaction with the large-scale distribution of positive-energy matter. Thus, we can now see why the rejection of the assumption that a uniform, large-scale distribution of negative-energy matter can exert a force on positive-energy matter (and vice versa), which appears to be required in order to arrive at a relational explanation of the phenomenon of inertia based on the principle of relativity, was, in effect, justified. The preceding discussion actually shows (when we recognize that positive- and negative-energy particles can interact only as a result of the fact that their presence is equivalent to missing vacuum energy) that this hypothesis is not only desirable, but actually constitutes an unavoidable consequence of the description of negative-energy matter as being equivalent to missing positive vacuum energy.

But in the context where the description of negative-energy matter as being equivalent to voids in the positive-energy portion of the vacuum is similarly applied to positive-energy matter (in the sense that positive-energy matter would be equivalent to the presence of voids in the negative-energy

portion of the vacuum) a further distinction would arise. Indeed, despite the fact that a uniform distribution of negative matter energy would have no effect on positive-energy matter, it can be expected that positive-energy matter, as voids in the negative-energy portion of the vacuum, would have to interact with the uniform portion of negative vacuum energy, for the same reason that such voids can also be expected to interact among themselves. In fact, even if the missing negative vacuum energy was itself uniformly distributed throughout space, it would still exert an influence on positive-energy matter, despite the fact that a similar distribution of missing positive vacuum energy would have no effect on this positive-energy matter, because negative vacuum energy does interact with itself. In other words, the fact that a void in the uniform negative-energy portion of the vacuum, which is equivalent to the presence of positive-energy matter, could leave no outside, surrounding negative energy to affect the behavior of negative-energy matter (if this void is itself uniformly distributed over the entire volume of the universe) would not affect the ability for such a void to gravitationally attract *positive-energy* matter, that is to say, other voids in this uniform, negative portion of vacuum energy, because in such a case the interaction is actually occurring between the matter particles themselves and not between a particle and the surrounding vacuum with the same energy sign.

Finally, it may be of interest to mention that if we were to consider the effect on a positive-energy body of a void in a uniform negative-energy *matter* distribution, then, based on the above discussed insights, we should deduce that the outcome would be a gravitational attraction directed toward the center of the void. This could be predicted to occur in two different ways. First, given that we can now expect negative-energy matter to exert a gravitational repulsion on positive-energy bodies, then, on the basis of what has been learned concerning the effects of voids in a uniform matter distribution, we could conclude that the absence of gravitational *repulsion* in the direction of the void, consequent to the absence of negative-energy matter in this void, would give rise to an uncompensated *repulsive* force directed toward the center of the void, which would be equivalent to a gravitational attraction directed toward the center of that same void, but which would actually arise from the gravitational repulsion of the surrounding negative-energy matter. But given that we now also know that a uniform distribution of negative-energy matter has no influence on positive-energy bodies, it would seem preferable to derive the consequences of an absence of such negative-energy matter based

on an alternative approach which borrows from the results discussed in the preceding paragraphs.

Indeed, what allows me to conclude that a uniform negative-energy matter distribution has no effect on positive-energy bodies is the fact that such a matter distribution appears to be equivalent to the presence of a void of universal proportion in the positive-energy portion of the vacuum, which therefore leaves no surrounding positive energy to produce uncompensated gravitational forces. But then, if you remove negative-energy matter in a portion of this void the resulting configuration would be that of an imperfect void or an imperfect distribution of *absence* of positive energy from the vacuum. But a local absence of *absence* of energy is really just the same as a local presence of energy and if the energy that was absent (when negative energy was present) was positive, then the energy that is locally present will itself be positive. This local absence of negative-energy matter will thus be totally equivalent to the presence of an equivalent amount of positive-energy matter and should therefore be expected to produce on positive-energy bodies a gravitational attraction directed toward the void. This is an effect which may have interesting consequences on the cosmological scale, in situations where variations in the density of negative-energy matter would have a magnitude comparable with the average density of the matter itself. I will explore the practical consequences of this important result in section 4.3. But for now, let me mention that the effectiveness of the preceding description is a further confirmation of the existence of a close relationship between vacuum energy and matter energy, while the high level of symmetry involved also indicates that the description of negative-energy matter proposed above fully agrees with the requirement of a relational definition of the physical attribute of energy sign.

2.7 Six problems for negative-energy matter

The preceding discussion may already make us feel more comfortable with the possibility that there could exist negative-energy matter, despite the traditional reluctance to accept the reality of negative energy states. But at the current stage of my account, this confidence would not yet be totally appropriate. Even in the context of the new understanding unveiled in the previous sections, there remain many problems associated with the possibility that negative-energy matter may exist in our universe. First of all, we do

not observe in the universe any matter or celestial object which would clearly appear to be involved in repulsive gravitational interaction with other material bodies. This is a very basic, but also very constraining fact. Associated with this problem is the fact that the current predictions of quantum field theory are based on a systematic rejection of the possibility of a transition to negative action states (as states of negative energy propagated forward in time or positive energy propagated backward in time) and yet they appear to produce results which agree very well with observations in all situations where the nature of the interactions involved is well understood and computational methods are sufficiently well developed to allow the derivation of such verifiable predictions. This could provide an additional motive for arguing against the possibility of the existence of negative-energy matter. Such pieces of evidence certainly cannot be dismissed without very good reasons. Any theory involving particles propagating negative energies forward in time must explain why it is that we can safely ignore the existence of those particles in formulating a quantum theory of elementary particles and their interactions, even while we would presumably have to take some of their effects into account in an astronomical context, where gravitational forces are not negligible.

A second category of difficulty has to do with the possibility that seems to be allowed, in the context where negative-energy particles would exist, for the annihilation of particle-antiparticle pairs to occur in which one of the particles would have negative action, therefore permitting matter to vanish, leaving absolutely nothing behind. This would, of course, require the annihilating opposite-energy particles to also have opposite electric and other non-gravitational charges, because charge must still be conserved. We have no reason, however, to assume that negative-action matter does not also come in two varieties, one propagating negative energy and all non-gravitational charges forward in time and the other propagating positive energy and the same charges backward in time (so that we have opposite charges from the forward time viewpoint). Therefore, we cannot *a priori* reject the possibility that such annihilations could take place. But that is a much worse problem than may perhaps appear to be, because if such annihilations were possible there would then be no reason why the time-reverse processes could not also take place. If that was the case, it would actually mean that pairs of opposite-action particles could be spontaneously created out of nothing without immediately returning to the vacuum like ordinary particle-antiparticle pairs, given that the process could occur without requiring a violation of

energy conservation.

The fact that opposite-action particles would gravitationally repel one another cannot be expected to prevent an annihilation process involving such particles from taking place, as the gravitational interaction is very weak and quantum fluctuations in energy would still allow the process to occur. Indeed, if the electrostatic attraction between opposite charges does not prevent ordinary particle-antiparticle pair-creation processes from occurring, then there is no reason why gravitational repulsion would need to be taken into account in the case of pair-annihilation processes involving opposite-action particles. In any case, the fact that the gravitational repulsion between opposite-energy particles would not affect the possibility for the associated *creation* processes to occur means that the problem is real. It may, therefore, seem like positive-energy matter particles could annihilate to nothing at an arbitrarily large rate upon encounter with negative-energy particles, or else be created out of nothing abundantly, even under ordinary circumstances, while both kinds of phenomena would clearly violate observational constraints, which actually provide no evidence at all that such events are taking place. This category of difficulties may then appropriately be called the energy-out-of-nothing problem.

A third potential problem has to do with the possibility that appears to be offered, as a consequence of the existence of negative energy states, for ordinary positive-energy matter particles or even any pre-existing negative-energy matter particles to ‘fall’ into the allowed negative energy states in a continuous, unstoppable process during which they would either release positive radiation or absorb negative-energy radiation and reach ever ‘lower’ energies. This is a difficulty which would also affect negative-energy matter as it is traditionally conceived and which is known as the vacuum decay problem. It would arise from the fact that the zero-energy level would no longer constitute a minimum level of energy (the ground state) from which there can be no transition to lower energies. Here we appear to have a situation where the existence of negative energy states raises the specter of allowing an arbitrarily large amount of work to be generated out of nearly nothing (by letting matter fall into the negative energy states and using the energy difference to produce work), as if energy conservation alone was not enough to restrict the evolution to negative energy states. This is clearly another issue of incompatibility with observation, because such decays are not observed to occur, even under the previously discussed conditions where negative energy densities are allowed to occur in a limited way by ordinary quantum

field theory. In this context, we may, in effect, ask what it is that prevents positive-energy particles from falling into the lower negative energy levels which are predicted to exist, under particular circumstances, by quantum field theory? This is all by itself a legitimate question which has remained unanswered. Even from the viewpoint of the traditional interpretation of negative energy states this situation looks like a deep mystery.

But what is probably the most serious problem which one must face upon recognizing the necessity of introducing a notion of negative-energy matter obeying the requirements of a relational definition of physical quantities (which imply that opposite-energy bodies must gravitationally repel one another) is that the existence of such matter may appear to allow violations of the principle of conservation of energy. This issue arises as a consequence of the fact that it seems possible for energy and momentum to be exchanged between positive- and negative-energy systems in a way that is similar to that by which positive-energy systems exchange energy among themselves. Basically, it appears that the positive energy of a positive-energy body can be turned into an equal amount of negative energy belonging to a negative-energy body (or vice versa) when a ‘collision’ between two such opposite-energy bodies would occur. For example, positive energy could be lost by a positive-energy body colliding with a negative-energy body initially at rest, while negative energy would be gained by the negative-energy body with which the first body has interacted. This would give rise to a net variation in the total energy of the two bodies that would be equal to twice the individual changes of energy (rather than allowing a cancellation of changes, as is observed when two positive-energy bodies collide). The solution to that problem will have to arise from a proper understanding of the fact that it is not possible, in a general relativistic context, for the energy of matter to be conserved independently from that of the gravitational field.

A further difficulty could arise in the context where the inertial force on a negative-mass body has the same direction as that which applies on a similarly accelerating positive-mass body, despite the reversal of inertial mass, which I have argued must occur when gravitational mass itself reverses. Indeed, from the viewpoint of an improved conception of the phenomenon of inertia based on a generalized formulation of Newton’s second law, it is no longer possible to consider that acceleration would take place in the direction opposite the applied force for a negative-mass body, given that the equivalent gravitational field due to acceleration would be reversed for such an object, which means that the inertial force it would experience is identical to

that which is experienced by a similar positive-mass body. It would therefore appear that while the presence of a negative-mass body could contribute to reduce the gravitational mass in a region of space in which positive-mass matter is also present, it would still provide the same resistance to acceleration, despite the fact that it would also provide a negative contribution to the inertial mass contained in this volume. This may not be a problem when we are dealing with independent physical systems with opposite masses, but as I previously mentioned, when a bound system is involved, the energy contained in the field of interaction between its constituent particles would be opposite that of the system as a whole and in such a case it would seem that while the energy of the field should reduce the gravitational mass of the system, it should nevertheless contribute to increase its resistance to acceleration. Given that bound systems with various force field configurations are quite common, it would seem that objects made of different materials should experience distinct accelerations when submitted to a gravitational force, but no such variations are observed. Some much-needed clarification is required here, if the concept of negative mass which I have proposed is to be considered viable from an observation viewpoint.

One last potential category of arguments which one might believe could disprove the validity of the idea of gravitationally-repulsive, negative-energy matter does not actually have to do with the concept of negative-energy matter developed here, but merely with more traditional concepts of antigravity and gravitational repulsion. The problems involved would be difficulties for a theory according to which ordinary antimatter is gravitationally repulsive. They would also constitute a challenge for the traditionally favored interpretation of negative energy states, according to which gravitational repulsion is an absolute property of negative-energy matter itself, while gravitational attraction is an absolute property of positive-energy matter (so that negative-energy matter repels positive-energy matter and is attracted to it). If such conceptions were to be retained as valid, they would allow paradoxical situations such as perpetual motion and time travel to arise. Given that, for most people, those difficulties are associated with the general concept of negative energy, it is important to explain why the issues involved here would not affect a more consistent theory of gravitationally-repulsive, negative-energy matter such as that which will emerge from the developments I introduced in the preceding sections.

We are then faced with six categories of problems, which appear to undermine a conception of physical reality according to which matter would

be allowed to occupy energy levels below zero. I have wrestled with the questions raised by those difficulties for a long time and on many occasions I had nearly given up on the possibility to ever be able to find appropriate answers that would perhaps explain why negative energy is not an inappropriate concept for physical theory. But, gradually, I came to understand that each of those problems really has to do with one or more incorrect, implicit assumptions we make when considering the expected behavior of matter, in the context where those negative energy states are actually allowed to be occupied. In the next six sections I will explain the nature of the insights required to appropriately deal with those severe problems.

2.8 The origin of repulsive gravitational forces

When, as a young man, I first started to contemplate the possibility that there could exist matter in a state of negative energy, I soon realized that if such matter was to attract matter of the same type while it would repel ordinary matter and be repelled by it (as I had intuitively assumed should occur, ignorant of the dominant paradigm), then this matter would have to be dark, because nowhere was it mentioned that we observe gravitational repulsion arising from the presence of any planet, star or galaxy. While I was working on improving my understanding of physics in general and trying to develop a theory incorporating the concept of negative mass, I simply assumed that negative-mass particles were such that they would interact with ordinary matter only through gravitation. I remember that I had read that Feynman once said that we must not question *why* things are the way they are, but simply try to describe in the most accurate way possible *how* they behave. Thus, for a while, I was comfortable with the idea that negative-energy matter simply does not interact, other than through the gravitational force, with ordinary matter (although it could interact with itself through the whole spectrum of forces), even if I had no idea why that should be the case and had to assume that this is just the way things are. The only concern I had regarding this situation is that it appeared odd that negative-energy matter should not interact with ordinary positive-energy matter through the same interactions by which positive-energy particles were interacting among themselves, given that negative-energy matter could be assumed to actually be composed of the exact same particles as positive-energy matter. But then came the shock.

I had for some time tried to figure out what determined the repulsive or attractive nature of an interaction, which clearly depends on the signs of the charges of the interacting particles, and had slowly come to realize that this property seemed to be related to the sign of energy of the field of interaction, not yet fully aware that it is actually rather the attractive or repulsive nature of an interaction (determined by the sign of the charges involved) that determines the sign of energy of the field and not the opposite. In any case, I had understood that the energy of a field associated with a repulsive interaction between positive-energy particles, for example the energy of the electromagnetic field between two electrons, is always positive, while the energy of a field associated with an attractive interaction between positive-energy particles, for example the energy of the electromagnetic field between an electron and a positron, is always negative. But it also had to be the case (as I will explain below) that the energy of a field associated with a repulsive interaction between negative-energy particles is always *negative*, while the energy of a field associated with an attractive interaction between negative-energy particles is always *positive*. What this means is that when two negative-action particles are attracted toward one another or bound together in a single system, the contribution of the attractive field mediating the interaction to the energy of the whole system should be positive, while for positive-action particles it would be negative.

As I was trying to make sense of this observation in the context where the interaction involved would be that between a positive- and a negative-energy body, I suddenly realized that a catastrophe had just happened. The problem is that, if this relation between the sign of energy of the field and the attractive or repulsive nature of the related interaction is right in general, it means that any gravitational interaction between positive- and negative-energy bodies should be either repulsive for positive-energy matter and attractive for negative-energy matter (if the field is attributed positive energy) or repulsive for negative-energy matter and attractive for positive-energy matter (if the field is attributed negative energy), but never repulsive for both the positive- and the negative-energy bodies involved in the interaction. This is because a repulsive field would have to have positive energy for a positive-energy matter particle, while this same positive-energy field would have to exert an attractive force from the viewpoint of a negative-energy matter particle for which the same relation would exist in general between the *difference* between the signs of energy of the matter particle and its field on the one hand and the repulsive or attractive nature of the associated interaction on the other (the

problem is not restricted to gravitation). This is again a consequence of the requirement of relational definition of the physical properties associated with attraction and repulsion which cannot be considered to be determined by the energy sign of the interaction field only, but must be a consequence of the difference between the energy sign of the field and that of the matter particles submitted to the force associated with this field.

But it was just nonsense to conclude that an interaction could be both attractive and repulsive at the same time and it is even more so now, in the context where we must recognize that the hypothesis of the mutual gravitational repulsion between positive- and negative-energy matter is also necessary for a relational description of the gravitational interaction between those two types of objects. The conclusion I had to draw was thus very clear: no definite energy sign could be attributed to the fields of interaction between positive- and negative-energy particles (as must be the case for any interaction involving particles with the same sign of energy) and therefore there simply cannot be any interaction between those two types of particle, not even gravitational. This appeared to be a fatal blow, because if there are no interactions of any kind between positive- and negative-energy matter, then how could negative-energy matter have any relevance to the world we experience?

When I realized the existence of this difficulty for a theory of negative-energy matter, I had already come to appreciate the many advantages that there would be if such matter was allowed to exist (if it could indeed gravitationally interact with ordinary matter). This is because I had been able to solve important problems even while merely using the incomplete description I had by then managed to develop and it seemed improbable to me that the whole idea could simply be wrong. I know that this may look like it was more a hopeful wish, than a rational conclusion, but in fact it was actually both hope and reason. Indeed, we had struggled with the problems I was able to solve for a very long time and there really appeared to be no viable alternative solutions to those problems, while, theoretically, the hypothesis that there could exist matter in negative energy states had a lot of appeal. It is as a consequence of the fact that I had so much confidence in the validity of the basic concept of a symmetry between positive and negative energy states that I did not stopped working on developing the idea when I encountered the difficulties discussed here. And as it turned out, the problems encountered became just another challenge on the way to a satisfactory solution to the problem of negative energy.

So, I went from having to explain why there would be no electromagnetic interactions between positive- and negative-action matter to having to explain why there can be any interaction at all between the same two kinds of matter. Of course, I was glad that, at least, I now had an explanation for why there is indeed no electromagnetic or other non-gravitational interactions between opposite-energy particles, because it was clear that, on the basis of the above discussed observations, it had to be recognized that there cannot be any direct quantized interactions (mediated through the exchange of interaction bosons) between such particles. But gravitation is different, because it is not yet described as a quantized field and I had hope that it might be its singular classical character that would allow the existence of *some kind* of interaction. It must be clear, however, that the problem described above is very real and unavoidable and its significance should not be underestimated, as it actually means that there can be no direct interaction between positive- and negative-action particles. It must also be understood that this is not a hypothesis, as no consistent theory could describe such an interaction and this must simply be understood to imply that those hypothetical interactions are, in effect, nonexistent⁵.

At this stage you may remember that, when I explained that there must be an equivalence (for a positive-energy body) between the effects arising from the presence of a void in a uniform positive-energy matter distribution and those which we may identify with a gravitational repulsion directed away from the void, I insisted that this repulsion was really the consequence of an uncompensated gravitational *attraction* directed away from the void. Therefore, when dealing with matter distributions which are uniform on a cosmic scale, we can observe gravitational repulsion to arise from what are actually purely attractive gravitational interactions. I also insisted that negative-energy matter would be equivalent, from a classical gravitational viewpoint, to the presence of missing positive energy from the vacuum, while the vacuum can itself be considered as being equivalent, to some extent (only in this respect), with a uniform matter distribution. But this means that the gravitational repulsion experienced by a positive-energy body, and which

⁵Those conclusions are the reason why I did not argue in this report that the gravitational interaction between opposite-energy particles must be considered repulsive merely on the basis of the fact that gravitation is mediated by a spin-two interaction boson (the graviton), because, obviously, if there cannot be any direct interaction between opposite-energy matter particles, then it is pointless to argue that it is the spin of the particles they exchange that determines the repulsive nature of such an interaction.

we would expect to arise from the presence of negative-energy matter, actually results from an uncompensated gravitational attraction attributable to the surrounding positive-energy portion of the vacuum. In other words, we can explain the gravitational repulsion apparently exerted by negative-energy matter as really consisting of a gravitational attraction involving only positive energy sources.

Thus, even if we assume an absence of direct interaction between positive- and negative-energy bodies, we can nevertheless expect to obtain an equivalent repulsive gravitational force between these objects. It is in this particular sense that the concept of gravitationally-repulsive matter developed here can be assumed to involve effects which are analogous to the situation we have in the case of voids in a uniform matter distribution. But under such circumstances the above discussed problem of the impossibility of direct interactions of either gravitational or non-gravitational kind between positive- and negative-energy particles is turned into an advantage, because it actually forbids any interactions to occur between opposite-energy particles, except for the equivalent gravitational repulsion just described, and this is precisely what we need. It must be clear, in effect, that the conclusion that there should exist indirect interactions between opposite-action particles only applies to gravitation, because even if a local absence of energy in the vacuum may always be correlated with a local absence of non-gravitational charge, it is not opposite charge particles which cannot interact with one another, but really opposite-energy or opposite-action particles.

What's important to understand is that while the negative-energy particles which exist as a result of a local absence of positive energy in the vacuum must have charges opposite those of the virtual particles which are missing as a result of this local absence of energy (because voids in a uniform positive or negative *charge* distribution must be equivalent to the presence of opposite-sign charges), such negative-energy particles do not interact only with particles from the surrounding negative-energy portion of the vacuum which have the same sign of charge, even though they do interact only with the negative-energy portion of the vacuum. But given that the vacuum is electrically neutral, then, even the negative-energy portion of it carries both positive and negative electric and other non-gravitational charges, which means that the surrounding negative-energy portion of the vacuum cannot exert an electrical force on a negative-energy particle with reversed electric charge, even in the presence of a void in the distribution of negative vacuum energy. As a result, even if a particle interacts only with the negative-energy portion of the vac-

uum, it experiences no external non-gravitational force of either attractive or repulsive kind arising from those interactions. One must, therefore, conclude that only charged particles with the same sign of action can interact through non-gravitational forces and that even non-gravitational forces of the indirect kind cannot exist between opposite-energy particles.

Those results should be encouraging, as the category of problems they allow to solve was the most basic and the most serious of those which I have identified above as facing a theory of negative-energy matter. Thus, it is now possible to explain why it is that we have never observed gravitationally-repulsive matter, because indeed such matter, if it exists, should not be visible, as it would not interact with ordinary positive-energy matter through the long-range electromagnetic interaction. It is also possible to explain why it is that the predictions of quantum field theory, made under the hypothesis that negative energy states are not allowed in the formalism, produce accurate results which correspond with observations to a very high degree of precision. Because if, in effect, only the equivalent repulsive gravitational interaction just described exists as a kind of influence of negative-energy matter on the processes involving positive-energy particles which are described by quantum field theory, then, given the weakness of the gravitational interaction, there should only be a marginal impact from the existence of this negative-energy matter on estimations of physical observables currently made under the assumption that negative-energy particles do not exist. Indeed, if we do not need to take into account the effects of the attractive gravitational interaction between ordinary *positive-energy* matter particles in such calculations, then we should certainly not expect to have to take into account any effects from the equivalent repulsive gravitational interaction with the very sparse amount of negative-energy particles that could perhaps be found to wander around apparatuses located on Earth. Thus, if I'm right, we would have here the solutions to two quite serious problems which were never addressed by any of the authors that previously discussed the possibility of gravitationally-repulsive matter, because it can now be understood, at once, why negative-energy matter is dark and why it would nevertheless appear to interact gravitationally with positive-energy matter.

It must be noted, however, that, even in the context where we have to assume that there is no direct interaction between positive- and negative-energy particles, it would be wrong to consider that positive-energy matter interacts only with the positive-energy portion of the *vacuum* and not with the negative-energy portion of it, because, as I explained in section 2.6, positive-

energy matter must itself be assumed to consist of voids in the negative-energy portion of the vacuum and as such, certainly cannot be considered to behave independently from this negative portion of vacuum energy. Yet it should be clear that we are not really dealing with an interaction between opposite-energy particles here, but merely with the gravitational interaction of this negative-energy portion of the vacuum with itself. Such a phenomenon is somewhat similar to the gravitational dynamics of a uniform negative-energy matter distribution, in which voids may also be present that would exert attractive gravitational forces on each other and repulsive forces on the rest of the negative-energy matter. In such a case it is clear, indeed, that even if the voids were equivalent to the presence of positive-energy matter, their effects would actually be the outcome of the interaction of negative-energy particles among themselves. We may, therefore, still consider that there is no *direct* interaction of any kind between positive- and negative-energy matter or vacuum, but again, this does not mean that positive-energy matter does not experience the gravitational effects of even the uniform portion of negative vacuum energy or that negative-energy matter does not experience the gravitational effects of even the uniform portion of positive vacuum energy, because, if positive-energy matter really constitutes missing negative vacuum energy, it cannot be expected that this portion of the vacuum does not interact with itself and the same can be said of negative-energy matter as missing positive vacuum energy. This conclusion will obviously have enormous consequences for the description of the cosmological effects of vacuum energy that will be discussed in chapter 4.

Finally, I may add that a further justification for the fact that we do not yet have strong evidence for the existence of negative-energy matter is that, as I will explain in chapter 4, it seems that, even though there was as much negative-energy matter as there was positive-energy matter in the initial Big Bang state, it must be recognized that only a negligible amount of baryonic negative-energy matter has survived the early annihilation of matter with antimatter and is still present in the universe today. I believe that this is what explains that no stellar- or galactic-size negative-energy matter overdensities large enough to exert a significant influence on the propagation of positive-energy photons (as could be detected by weak gravitational lensing experiments for example) has ever been observed. Of course, regardless of its abundance, negative-energy matter can be expected to migrate away from concentrations of positive-energy matter and to concentrate itself in regions of the universe where there is a lesser density of positive-energy matter, because

such matter is gravitationally repelled by positive-energy matter and is also gravitationally attracted to itself. It will, therefore, be difficult to observe anomalous gravitational effects arising from the presence of gravitationally-repulsive negative-energy matter overdensities in a region of the universe like ours, where positive-energy matter can be assumed to be the dominant form of matter, given its relatively large density. But once it is recognized that there is not much baryonic negative-energy matter left at the present time, then those considerations are not as significant as they would otherwise be. Anyhow, as we progress, it will transpire that the lack of evidence for negative-energy matter is now so well justified that it appears that, if we are to ever obtain direct confirmation of its existence, it will be necessary to use alternative methods of investigation and to concentrate on the possibility which may be offered to derive observational consequences of the presence of non-baryonic forms of negative-energy matter.

2.9 No energy out of nothing

Before we can conclude that there should indeed be no interference with current predictions, made using quantum field theory, from allowing the existence of negative-energy particles in stable states, we must first explain why it is that there should be no creation or annihilation processes involving pairs of opposite-energy particles with opposite charges, as such a phenomenon could also disrupt current predictions. This is the second category of problems I previously identified as potentially affecting the viability of the negative-energy-matter hypothesis. Given the plausibility of the hypothesis that negative-energy particles should be very rare in our region of the universe, it may seem that the problem of the annihilation of opposite-energy particles does not constitute a decisive issue. But, as I previously mentioned, we cannot avoid having to face the related problem of the creation of pairs of opposite-energy particles, because in such a case it would appear that no favorable initial conditions are required for the discussed processes to occur. Thus, an explanation must be provided for why matter is not being created out of the vacuum in massive amounts, even under normal conditions, despite the fact that the processes involved could occur without violating the principle of conservation of energy, because this prediction clearly disagrees with observations which indicate a complete absence of such processes.

One may perhaps suggest that it is the fact that the opposite-energy par-

ticles emerging from a creation event in opposite directions would have their momenta both pointing in the same direction (because we must assume that a negative-action particle would have momentum opposite the direction of its velocity) that prevents the creation of such pairs, when we require momentum to be conserved. But it does not seem that this would constitute a strong enough constraint under appropriate circumstances, because the pairs could be created without much momentum, or through an input of momentum from the environment, as is the case for ordinary particle-antiparticle creation processes arising from the disintegration of a single boson. It is not possible, therefore, to conclude that it is the requirement of momentum conservation which prevents pair-creation processes involving particles with opposite energies from occurring.

The fact that the kind of creation (or annihilation) processes which would require no energy input (or output) could be described as processes during which a particle reverses its direction of propagation in time while retaining the sign of its energy, may suggest another explanation for why such events would be forbidden. Indeed, we may ask why it is that when a particle changes its direction of propagation in time, in the course of all those particle-antiparticle annihilation processes which do occur under the right conditions, the energy is invariably reversed relative to the new direction of propagation in time (so that it appears to be unchanged from the forward time perspective)? Why must it be imposed that a reversal of the direction of propagation in time be combined with such a reversal of energy which leaves the sign of action invariant, so that the energy of the annihilating pair needs to be compensated by the emission of photons carrying away the energy? Could it be that it is a requirement of continuity of physical properties along the world-lines of elementary particles that prevents a positive-action particle from turning into a negative-action particle? Such a change would, in effect, involve the transformation of a particle experiencing the gravitational interaction in a given way into a particle experiencing it in a different way, but perhaps that a particle cannot change the way it gravitationally interacts with the rest of the universe along a continuous world-line.

I must acknowledge that I once contemplated the possibility that action-sign-changing reversals of the direction of propagation in time may be forbidden by a requirement of continuity of physical parameters along a particle's world-line. But I later came to understand that what such a requirement of continuity imposes is merely an absence of interruption of the flow of the fundamental time-direction parameter, which can be satisfied even when the

energy of a particle does not reverse upon a change of its direction of propagation in time. In section 4.3 I will explain what constraint a condition of continuity of the flow of time along an elementary particle world-line would impose on the transformation of physical parameters and it will be clear that a reversal of the action is not forbidden by such a requirement. In any case, if the *charge* of a particle can vary discontinuously (can reverse) from the forward time viewpoint when the particle reverses its direction of propagation in time in a continuous fashion (during a process perceived as an ordinary particle-antiparticle annihilation process), then there is no *a priori* reason why the *action* of a particle could not reverse in a similar manner when the particle reverses its direction of propagation in time, if the reversal also occurs in a continuous way (without the direction of the flow of time being interrupted along the path of the particle in spacetime), which would simply mean that the particle does not actually experience the usual reversal of its energy sign at the bifurcation point, when it reverses its direction of propagation in time. But such a transformation is precisely what is never observed to occur. Must one, then, conclude that there exists an inexplicable decree, simply banning negative-action particles (carrying positive energy backward in time) from existing? This would again be the easy way out: there is a difficulty, so let's just forget about the whole thing. But if we recognize that the existence of particles carrying positive energies backward in time is theoretically inevitable, then a satisfactory explanation for the absence of spontaneous matter creation out of nothing is required.

In fact, the problem of the creation of pairs of opposite-action particles from nothing and the related problem of the annihilation of pairs of opposite-action particles to nothing turned out to be much simpler to solve than I had originally envisaged. To understand what imposes a limit on the creation and annihilation of pairs of opposite-action particles, we simply need to take into account the results obtained in the preceding section. Indeed, one may ask how it is supposed to occur that a positive-action particle with positive charge, say, could annihilate with a negative-action particle with negative charge if positive- and negative-action particles are to be considered as equivalent to voids in opposite-energy portions of the vacuum? How could the two particles ever annihilate one another, when annihilation is to be considered a kind of interaction and there is absolutely no direct interaction of any kind between opposite-action particles? Had I taken the lesson learned while solving the problem of the nature of repulsive gravitational interactions more seriously, I would have understood much more readily that what prevents the

creation and the annihilation of particles with opposite energy signs is the absence of any direct interaction between such particles. Indeed, in the absence of any direct interactions between them, two opposite-action particles with opposite charges would not be able to annihilate so as to produce a final state of null energy, even if they were to find themselves arbitrarily close to one another, and the same limitation would also make it impossible for two such particles to be created together out of nothing.

It is true, though, that opposite-action particles would, according to the results I derived in the preceding section, be subject to some indirect gravitational interaction, as a consequence of the equivalence between the presence of a particle with a given energy sign and an absence of energy of opposite sign from the vacuum. Thus, even though there are no direct interactions between opposite-action particles, it may perhaps seem that it would be possible for such opposite-action particles to be created from nothing, or to annihilate to nothing, as a result of an exchange of gravitational energy, arising from the indirect gravitational forces they exert on one another. But given that opposite-action particles cannot be produced together as a pair, then they would need to be produced as pairs of particle-antiparticle *pairs* with opposite action signs, because individual particle world-lines never emerge from nothing or vanish to nothing in the vacuum (this is what a condition of continuity of the flow of time along an elementary particle world-line actually forbids, as I will explain in section 4.3).

The problem, however, is that, given that positive- and negative-action particles do not interact directly with each other, then the energy of a pair of positive-action particles produced in such a way would have to come from an exchange of energy with the environment with which it interacts, because when energy is exchanged between opposite-energy systems it can only rise in magnitude for one of them when it diminishes in magnitude for the other one, as I will explain in section 2.11. Thus, if no radiation energy is present beforehand, then no particle-antiparticle pair can be produced in such a way, even if it may appear that energy would be conserved when an opposite amount of energy would be created at the same time as a result of the production of a particle-antiparticle pair with opposite energy, because even radiation particles cannot interact with one another, so as to be produced out of nothing, if they have opposite-energy signs. But this also means that when a particle-antiparticle pair annihilates, positive or negative energy must necessarily be released in the environment with which it does interact, even if the energy change involved could be compensated by a similar,

but opposite change of energy that would occur as a result of the annihilation of a particle-antiparticle pair with opposite action sign, because those radiation energies have opposite signs and positive-energy radiation cannot interact with negative-energy radiation, so as to perhaps annihilate to nothing. It would therefore appear that the absence of direct interaction between opposite action particles constitutes a necessary and sufficient condition to prevent matter and radiation energy from being created out of nothing or from annihilating to nothing.

Given those conclusions, it should be clear that the problem of the creation of opposite-action particles out of nothing, arises only when one fails to recognize that there can be no direct interactions between such particles. But it is important to realize that the solution proposed here to the problem of creation out of nothing does, in effect, also apply to hypothetical processes of annihilation to nothing, because, even though it may seem that the problem of the annihilation of opposite-action particles does not constitute a decisive issue in the context where there cannot be very many negative-action particles and antiparticles in our region of the universe (so that encounters between opposite-action particles should be rare), the situation was different in the very first instants of the Big Bang. When the magnitude of the densities of positive and negative matter energy is arbitrarily large, or actually maximum, as must have been the case in the initial singularity, and the positive and negative energies must be very homogeneously distributed in space (for reasons to be discussed in section 4.9), it follows that if the annihilation of opposite-action particles to nothing was not forbidden, then most of the matter, as well as most of the energy contained in radiation, would vanish within a very short instant, leaving absolutely nothing behind. This is certainly not an outcome that would agree with astronomical observations and therefore one must recognize that the validity of the explanation proposed here for the absence of creation of energy out of nothing is also confirmed by the observation that matter does exist in our universe at the present time. In fact, once this is understood, it transpires that it may not even be necessary for all matter to be created out of nothing at the Big Bang, if it is possible to assume that time extends past the initial singularity following a hypothetical quantum bounce, as the most promising, tentative quantum gravitation theories seem to indicate. The plausibility of the results discussed above would therefore merely provide one more reason to acknowledge that it is necessary to take those predictions seriously.

2.10 The problem of vacuum decay

There is an unavoidable question that arises whenever one proposes that negative energy states may be physically allowed. What is it, in effect, that prevents particles from falling into those ‘lower’ energy states? It has been argued that positive-energy matter particles may not be able to do so because they would first have to surmount the limit imposed by the irreducible value of their positive mass. But that would clearly not prevent particles already in a negative energy state from reaching even ‘lower’ energy states and given that I’m here working under the assumption that negative-energy matter can exist in stable form, this would appear to be a serious problem. Under such conditions it would seem that if even a small amount of matter was to ever find itself in one of the available negative energy states this would give rise to a catastrophic process of creation of negative matter energy and positive radiation energy, because the matter would radiate energy in going from the ‘higher’ energy states (with negative values nearer to zero) to the allowed ‘lower’ energy states (with larger negative values) without ever reaching a minimum energy in which it could settle down. Thus, as I mentioned before, it would seem that if negative-energy matter can exist, we could produce an infinite amount of work by simply harvesting the positive-energy radiation produced when negative-energy particles fall into lower negative energy states. But given that quantum field theory already allows for states of negative energy to occur in limited portions of space, it would seem that we have a very serious problem, even in the current theoretical context, because if negative energy can be made to exist under such conditions (which have already been produced in the laboratory) it should immediately collapse to even lower negative energies and in the process produce an arbitrarily large amount of positive-energy radiation, while, of course, no such a phenomenon has ever been observed.

The insights gained while studying the problem of energy creation out of nothing discussed above, however, provide the elements needed to tackle this additional difficulty from a different angle. Indeed, according to the preceding discussion, an important consequence of the absence of any direct interaction between opposite-action particles is that it is actually impossible for a particle to annihilate with one of its opposite-action antiparticle counterpart, which is another way to say that an already existing particle cannot reverse its direction of propagation in time without also reversing its energy sign (relative to its new direction of propagation in time), therefore describing

an ordinary particle-antiparticle annihilation process. But another, perhaps less obvious, consequence of the absence of any direct interactions between opposite-action particles is that a negative-energy particle cannot emit a real (by opposition to virtual) positive-energy interaction boson regardless of what energy changes the original particle goes through, because the positive-energy boson is not even allowed to interact with the negative-energy particle it is assumed to get its energy from.

Therefore, a negative-energy particle could not ‘lose’ energy (gain negative energy) through the production of a compensating amount of positive radiation energy and the same limitation also implies that a positive-energy particle couldn’t absorb negative-energy radiation and diminish its own positive energy in the process. This constraint must apply even if such processes could occur without violating conservation laws, when the energy change of the matter particle involved would be compensated by the emission or the absorption of an opposite amount of radiation energy. But this means that even the emission of *positive-energy* radiation by a positive-energy matter particle could not occur in such a way that the positive-energy particle could turn into a negative-energy particle, given that this would imply that there would have been a direct interaction between the matter particle that now has negative energy and the positive-energy radiation it would have released (at the particular point in spacetime where the reversal would have taken place), while according to my analysis this must be considered impossible, even for a massless positive-energy particle, because once the particle reaches a null energy, by releasing positive energy, it must continue to release positive energy if it is to reach a negative energy state, but once it crosses the zero-energy limit it would, in effect, be forbidden from emitting positive-energy radiation.

Thus, the same constraint whose existence allowed me to conclude that a particle cannot change its direction of propagation in time without reversing its energy sign, also implies that it is impossible for a particle to reverse its energy without reversing its direction of propagation in time (in which case the particle would not continue to exist with opposite energy in the future). The existence of such a limitation suggests that no interaction vertex involving particles with mixed action signs needs to be taken into account in determining the transition probabilities associated with quantum processes. A certain limitation against the possibility of transitions to negative energy states therefore actually exists, because a positive-energy particle cannot ‘fall’ into a negative energy state by releasing positive-energy radiation. The only

reversal of energy which may occur on a continuous particle world-line would have to involve a reversal of the direction of propagation in time, in which case the energy of the particle in its final state would not be reversed relative to the forward direction of time and we would merely observe a conventional antiparticle in a positive action state annihilating with the ‘original’ positive-action particle.

The limitation imposed on vertexes that they cannot involve particles with mixed action signs would therefore actually prevent a particle that is already in a negative energy state from falling into even ‘lower’ energy states by releasing positive-energy radiation, because such a negative-energy particle could never have interacted with the positive-energy radiation it is assumed to emit. In fact, this explanation works both ways, as it is also true that a particle in a negative energy state could not ‘gain’ energy and turn into a positive-energy particle by releasing a compensating amount of negative-energy radiation, because the interaction bosons so released could not have been emitted by the particle that now has positive energy at the particular point in spacetime where the reversal would have taken place, given that they cannot interact with such a particle. What must be understood, again, is that while the requirement of energy conservation may not alone forbid transitions involving a reversal of the sign of energy, the fact that those transitions would involve the emission or the absorption of radiation with an energy sign opposite that of the final or original particle (respectively) actually prevents them from occurring in the context where a negative-energy particle (be it matter or radiation) is not allowed to directly interact with a positive-energy particle.

Yet it must be remarked that the constraint described here would not prevent a negative-energy particle from absorbing negative energy radiation and ‘falling’ into ever more negative energy states, if such an evolution is favored from a thermodynamic viewpoint. In this particular sense it may therefore appear that a certain aspect of the problem of vacuum decay remains unsolved. I believe that the situation we have here is analogous to that which was faced upon the introduction of the Rutherford atom model, which was initially rejected despite its apparent empirical inevitability, because it was assumed that the electrons in orbit around the nucleus would lose energy in the form of electromagnetic radiation and end up collapsing into the nucleus, while no such catastrophe was observed. But just like the Rutherford model, it appears that negative energy states are unavoidable and thus a solution to the problem of vacuum decay that does not simply amount to reject the phys-

ical nature of those states must be provided. Based on the results achieved in the preceding sections, I would suggest that the difficulties described here arise, again, from the fact that we ignore the requirements imposed by the necessary relational definition of physical quantities.

What is happening is that we are attributing a preferred direction to energy variations without referring to a physical aspect from our universe relative to which that direction could be compared. In other words, we use an absolutely defined direction on the energy scale which we arbitrarily define as ‘lower’ and we attribute distinctive physical properties to energy variations occurring along that absolutely defined direction, despite the fact that it actually has no objective significance. This traditional assumption seems to be justified by the observation that, for positive energy states at least, there does exist a singled-out direction on the energy scale which is related to the natural tendency for matter to disintegrate and to reach thermal equilibrium. This direction can be associated with a well-defined physical aspect of our universe, which is the direction of time in which entropy is growing. In the absence of such a relationship, we would have no motive to assume the existence of a preferred direction on the positive energy scale, that would not necessarily be opposite any such direction on the negative energy scale.

However, when I examined what the motives are, exactly, which allow us to consider the existence of this objectively defined ‘lower’ direction on the positive energy scale, arising in relation to the direction of time in which entropy grows, I realized that there is absolutely no reason to assume that this direction on the energy scale can be extended into negative-energy territory without being subjected to a reversal like energy itself. The only assumption necessary to assert the validity of this conclusion is that the thermodynamic arrow of time points in the same direction from the viewpoint of both positive- and negative-energy observers, which certainly constitutes a plausible hypothesis, especially in the context of the explanation that will be proposed in chapter 4 for the origin of time asymmetry. Therefore, it seems that the objectively defined ‘low’ energy direction on the positive energy scale cannot be extended into negative-energy territory, but would actually be effective toward smaller, *less negative* energy states (toward the zero-energy ground state) for negative-energy matter.

Basically, what allows me to conclude that the low-energy direction, for negative-energy matter, is toward the zero energy, as is the case for positive-energy matter, is that the singled out, objectively defined direction on the energy scale is simply that relative to which the energy tends to dissociate

itself and to become less concentrated, so as to spread into a larger number of independent particles which thus necessarily have smaller (nearer to zero) energy as time goes. What explains this tendency is the fact that such a final configuration is associated with a larger number of microscopic degrees of freedom and a higher entropy (when gravitation can be neglected) and therefore is more likely to be reached in this direction of time in which entropy is actually allowed to grow. But, if the direction in time of entropy growth is the same for positive- and negative-energy systems, then the direction that would emerge as the low direction on the negative energy scale would have to be the opposite of that which constitutes the equivalent, objectively or relationally defined low direction on the positive energy scale, because the spreading of energy into a larger number of particles with smaller *negative* energies, which is necessarily associated with a higher entropy, occurs in the direction on the energy scale which is opposite that in which smaller positive energies are reached. Thus, what we traditionally called ‘low’ energies, far below the zero level of energy, are in fact high energies for negative-energy matter and what we called ‘higher’ energies, nearer to the zero level on the negative energy scale, are actually lower energies for negative-energy matter. This is in perfect agreement with the previously discussed requirement to the effect that there should be a symmetry under exchange of positive- and negative-energy matter, so that the sign of energy can be defined as a relational property.

Such a conclusion is significant, because it allows one to deduce that it is not to be expected that matter should have a tendency (arising from a thermodynamic necessity) to decay into larger negative energy states past the zero-energy level. Negative-energy matter must be expected to have the same tendency as positive-energy matter to decay to energy states which from the perspective of an observer made of such matter would be lower energies and therefore to produce a larger number of particles with smaller negative energies and reach for the vacuum ground state in the future direction of time. If matter was found in a negative energy state, it would not have a natural tendency to decay in a direction on the energy scale which is actually upward for a negative-energy observer. It would be incorrect to assume that negative-energy particles have a tendency to decay by spontaneously gaining negative energy through absorption of negative-energy radiation as time goes, because such configurations are not thermodynamically favored, but are actually less likely to occur, for the same reason that positive-energy matter particles are not likely to reach states where energy would become more concentrated

into fewer particles as a result of the absorption of positive-energy radiation. As a consequence, regardless of the energy level in which a positive-energy particle is to be found at a given time, it can only release radiation until it reaches the energy contained in its rest mass and if it disintegrates and loses its mass it is not to be expected that it would continue to decay past the zero level of energy by gaining more negative energy through absorption of negative-energy radiation.

The unavoidable character of the conclusion that there is no preference for states of larger negative energy means that there should be no continuous decay to more concentrated, negative energy states, especially in the context where there already exists a constraint on the emission of positive radiation energy by matter entering a negative energy state. It would not be possible, therefore, to produce a large amount of work by making use of processes during which particles would gain larger negative energies, either by releasing positive-energy radiation, or by spontaneously absorbing negative-energy radiation, despite the fact that matter is actually allowed to occupy those negative energy states. I should, finally, mention that the fact that we observe no catastrophic collapse to larger negative energies, under the conditions where small negative energy densities are routinely produced in a limited way (as when a negative pressure is observed between two parallel mirrors in a vacuum), is a confirmation of the validity of the conclusions discussed in this section.

Thus, the outcome of the progress achieved in the last two sections is that it is possible to conceive of a fully consistent interpretation of negative energy states that would allow to at least preserve the validity of the current framework of quantum field theory. Indeed, it would appear that what we obtain are two more or less independent frameworks describing two more or less independently evolving categories of systems with opposite energies, which interfere with one another only under those special conditions where it is possible for an observer made of matter with one energy sign to indirectly deduce the existence of opposite energy densities as they occur in the context where constraints are imposed which forbid the presence of certain states which would otherwise be present in that portion of the vacuum which is directly experienced by such an observer. This particularity allows the near perfect agreement between the predictions and the observations related to the small-scale realm of quantum theory to naturally be maintained, despite the fact that it is possible for matter to occupy the available negative energy states; which is also remarkable.

2.11 Energy and momentum conservation

I would now like to discuss the case of that most serious of problems, which could have proved fatal to the alternative concept of negative energy developed here and which I have identified above as being that raised by the apparent possibility of a violation of the law of conservation of energy, under conditions where interactions (even if merely of the indirect kind envisaged here) are allowed to occur between positive- and negative-energy matter. The nature of the issue can be illustrated through the use of a simple thought experiment. I briefly discussed, in a previous section, the problem that would arise in the case where a ‘collision’ would occur between a positive-energy body and a negative-energy body. I explained that such a collision would involve a loss or gain of positive energy by the positive-energy body that would not be compensated, but instead be made worse by the associated gain or loss (respectively) of negative energy by the negative-energy body. This is because, instead of witnessing a loss of energy by one particle that would be gained by another, as when two particles with the same energy sign collide, we would here seem to have equal variations of energy, either both positive or both negative, depending on which particle accelerates and which decelerates as a result of the collision. For example, a negative-action body could lose negative energy, while the positive-action body it repels would gain positive energy, resulting in a net overall increase of energy twice as large as the individual changes. It would then seem that energy conservation is not possible under such circumstances.

The problem discussed here is also apparent when we consider the variations of momentum involved in such a process. Indeed, if action is to be assumed negative for a particle propagating negative energy forward in time, then it means that the sign of its momentum relative to its direction of propagation in *space* must be negative, that is, momentum must be opposite the direction of the motion for a negative-energy particle (because action has the dimension of an energy multiplied by a time or that of a momentum multiplied by a distance). In such a context, it is easy to deduce that the variation of momentum occurring upon a collision between two opposite-energy bodies would be twice as large as the absolute values of the changes in each particle’s momentum rather than be zero, as when two positive-energy bodies collide. This is a problem that does not exist in the context of the traditional conception of negative-energy matter, according to which positive-energy bodies attract negative-energy bodies, which repel them (if we assume that only

gravitational forces exist between opposite-energy bodies) and therefore the existence of such a difficulty could be used as an argument in favor of this traditional viewpoint, despite the fact that it also raises other problems of its own, as I previously explained.

But given that we now understand that there are no direct interactions between opposite-energy particles, we have to recognize that the only way a collision between opposite-energy bodies could occur would be through the indirect gravitational repulsion that would arise as a consequence of what are actually attractive gravitational forces, attributable to the distribution of vacuum energy that surrounds those objects and which are made to exist as a consequence of the equivalence between the presence of matter of one energy sign and an absence of energy of opposite sign in the vacuum. In this context, it would in fact appear unlikely that there could occur violations of energy conservation arising from a collision between positive- and negative-energy bodies, if indeed there are no direct interactions between such objects. Mathematically at least, it certainly seems that a general-relativistic theory of negative-energy matter which would involve only gravitational interactions should not give rise to violations of the law of conservation of energy, given that energy conservation in such a context is actually a constraint concerning the exchange of energy between matter and the gravitational field.

Thus, if opposite-energy bodies do interact only through those indirect gravitational interactions, then it means that from the viewpoint of a general-relativistic description of those interactions any variation in the energy of matter would, in effect, come from a variation in the energy of the gravitational field attributable to the changes generated by those interactions in the energy of the vacuum. The absence of any other interaction between positive- and negative-energy bodies should indeed allow one to expect that it would be variations in the energy of the gravitational field that would balance the variations of energy occurring in the course of the interaction of such opposite-energy bodies. The problem I initially had, however, is that I was not able to figure out how this could come about in the more intuitive context of a Newtonian description of such interactions and I'm always suspicious of conclusions drawn solely on the basis of mathematical deductions, which often conceal totally inappropriate assumptions. So, where exactly does the positive energy go, which is lost by a fast-moving positive-energy body colliding with a negative-energy body initially at rest and where does the negative energy come from, which is gained by the negative-energy body that is accelerated during such a collision?

I was allowed to understand what is going on when a positive-energy body interacts with a negative-energy body only when I became aware of the possibility that the energy of matter and its gravitational field may be null for the universe as a whole. Indeed, as certain authors now recognize, it appears that when positive-energy matter collapses into a spacetime singularity its negative gravitational potential energy becomes equal in magnitude to the energy of the matter itself (even though, in the case of a future singularity, this can only happen after an event horizon has already formed). Thus, if the initial Big Bang state must be considered to consist of a spacetime singularity (which is required even in the presence of negative-energy matter, for reasons I will discuss in chapter 4), then it means that the gravitational potential energy of positive-energy matter was initially the exact opposite of the energy of this matter. As space expanded this potential energy immediately began to decrease (toward the zero value) along with the positive kinetic energy of expansion, but it remains that under such circumstances there naturally occurs a compensation between the energy of matter and its gravitational potential energy (although it is actually the kinetic energy of expansion that must compensate the gravitational potential energy at all times, as I will explain in section 4.5).

In the case at hand, what happens, therefore, is that, given that the depth of the void in positive vacuum energy that is equivalent to the presence of the negative-energy body grows larger, along with the negative energy of the object, as a result of the gravitational force exerted by the void in *negative* vacuum energy that is equivalent to the presence of the positive-energy body, then one must conclude that, following the interaction, more gravitational attraction goes missing between all positive-energy matter in the universe and the *positive* portion of vacuum energy. But the potential energy of the attractive gravitational interaction between positive vacuum energy and positive matter energy is negative, so that if it goes missing the outcome is a positive variation of energy that exactly compensates the negative energy gained by the void, that is to say, by the negative-energy body⁶. What's crucial to understand here, is the fact that the interaction of the positive-energy

⁶In fact, the missing positive vacuum energy would actually interact with all positive energy matter, but also with the rest of positive vacuum energy and all negative vacuum energy. However, given that the gravitational potential energy attributable to interaction with the vacuum cancels out for the most part, then it is appropriate to consider that most of the potential energy that is lost concerns interaction with positive matter energy and must, therefore, be negative.

body is with the positive-energy portion of the vacuum and when less positive vacuum energy is left to interact with all the positive-energy matter in the universe *as a result of the local influence exerted by the positive-energy body on that portion of the vacuum*, then more interaction goes missing globally between all the positive-energy matter in the universe matter and the positive portion of vacuum energy. Yet given that a negative gravitational potential energy would have been associated with those interactions, it follows that the negative gain in energy by the negative-energy body, whose presence is equivalent to an absence of positive vacuum energy, is exactly compensated by the loss of this negative gravitational potential energy: more missing positive vacuum energy, more missing negative gravitational potential energy. And the same conclusion holds for the loss of energy by the positive-energy body: less missing negative vacuum energy, less missing positive gravitational potential energy. That's all there is to it.

It should be clear that I'm not saying that it is the gain of positive gravitational potential energy attributable to the interaction of a negative-energy body with all the matter (with the same energy sign) in the universe that would compensate a negative gain of energy by that very same body, which would mean that no interaction may be required to trigger those changes, which could then occur without any identifiable cause (for both positive- and negative-energy matter). It is merely the changes in negative gravitational potential energy which are attributable to the influence exerted by a positive-energy body on local measures of positive-vacuum energy *in its environment* and which are produced by the indirect gravitational interaction between positive- and negative-energy bodies, which can compensate a variation in the negative energy of matter. From my perspective, what happens during such an interaction is that the loss of positive kinetic energy by the positive-energy body is compensated *not* by a gain in positive gravitational potential energy of the negative-energy body, but by the loss in negative gravitational potential energy of the positive-energy portion of vacuum energy whose density decreases locally in proportion to the energy gain of the negative-energy body, while the gain in negative kinetic energy experienced by the negative-energy body is itself balanced not by the loss in negative gravitational potential energy of the positive-energy body, but by the gain in positive gravitational potential energy of the negative-energy portion of vacuum energy whose density increases locally in proportion to the amount of energy that is lost by the positive-energy body, as a result of its indirect interaction with the negative-energy body.

The fact that, under all circumstances, only as much energy as is present in a field of interaction can actually be exchanged between the particles interacting through that force field means that the energies exchanged during the process of indirect gravitational interaction between a positive- and a negative-energy body are relatively small, but they are not completely negligible and it is possible to understand how it is exactly that they are compensated, when one takes into account the variation of gravitational potential energy attributable to the related local changes in vacuum energy density. It must be clear, however, that we are not dealing here with the gravitational potential energy that could be associated with a repulsive force field mediating the interaction between the positive- and negative-energy bodies themselves, which, in fact, cannot exist, as I explained before, but merely with *independent* measures of gravitational potential energy associated with the interaction of each of the two opposite-energy portions of zero-point vacuum fluctuations with all the matter in the universe that shares the same sign of energy.

What must be understood, therefore, is that, following any interaction between a positive-energy body and a negative-energy body, there actually occurs a variation in the total energy of matter of both positive- and negative-energy sign, but this is only half of the equation, as to any such change there must be a related compensating change in the gravitational potential energies attributable to the variations that take place in the energy of the vacuum, from the viewpoint of observers with the same sign of energy as that of the portion of vacuum energy that must be assumed to vary locally. But it must be clear that this is only a reflection of the exchanges of gravitational energy occurring between positive-energy matter and the positive portion of vacuum energy, on the one hand, and between negative-energy matter and the negative portion of vacuum energy, on the other, because there is no actual exchange of energy between those two kinds of matter.

One must, therefore, conclude that kinetic energy is exchanged between positive- and negative-energy bodies as if it was a positive-definite quantity (from the viewpoint of positive-energy observers, or as a negative-definite quantity from the viewpoint of negative-energy observers), which means that if the magnitude of the energy of a body with a given energy sign varies as a result of its interaction with a body with opposite energy sign, then the magnitude of the energy of that second body should necessarily vary in the opposite way, except while the interaction is under way and changes in the kinetic energy of matter are compensated by *local* changes in the

gravitational potential energy associated with the interaction of each of the two bodies with their same-energy-sign portions of the vacuum. It must also be mentioned that the variation in the momentum of matter which would be observed to take place during such an indirect interaction is also compensated by an opposite variation in the momentum of the gravitational field (or the equivalent component of space curvature), which occurs as a consequence of those local changes in the distribution of vacuum energy. The fact that the gravitational interaction is very weak means that the energy flow between matter and gravitational field that occurs in the course of any indirect gravitational interaction between opposite-energy bodies is relatively small, but it nevertheless exists and it appears to be what allows energy to be conserved during such interaction processes.

2.12 Absolute inertial mass

One last objection which could be raised against the interpretation of negative energy states which I proposed has to do with the fact that, from my viewpoint, negative-energy matter would offer the same resistance to acceleration as would positive-energy matter. Traditionally, this would occur whenever we would assume that inertial mass is positive, even for negative-energy matter otherwise characterized as having a negative gravitational mass. Of course, as I already explained, the inertial mass must be considered to actually be reversed, along with the gravitational mass, from the viewpoint of a consistent description of the gravitational dynamics of negative-energy matter. But in the context of the previously discussed, improved conception of the phenomenon of inertia that emerged from my generalization of Newton's second law, it was shown that acceleration would not occur in the direction opposite the applied force for a negative-mass body. In fact, once it is recognized that the equivalent gravitational field experienced by such an object must be opposite that experienced by a positive-mass body, it is necessary to conclude that negative-mass matter would actually experience the same resistance to acceleration as positive-mass matter, when submitted to the same forces, despite the reversal of its inertial mass. Thus, negative-mass or negative-energy matter would appear to violate the principle of equivalence as it is traditionally conceived.

In fact, there could also be situations where the gravitational mass in a volume of space would be relatively small or even zero, despite the pres-

ence of a potentially large amount of matter in this volume, as when two opposite-mass bodies are present all at once in the same location (which would be allowed in the absence of strong interactions between them). Yet such configurations would not be equivalent, from an inertial viewpoint, to the case of a system with nearly vanishing total mass, because the matter that is present would be more difficult to accelerate than if it actually had such a small mass. To better describe such vanishing energy configurations, which are clearly different from the vacuum, we may define a measure of inertial mass that would be related to the physically significant properties with which it is traditionally associated and that would correspond to the true amount of matter present under such circumstances, independently from the total amount of mass, which may partially or totally cancel out. The *absolute inertial mass*, obtained by adding the absolute values of the masses of all material bodies present in some volume of space (or by adding all masses as negative from the viewpoint of a negative-mass observer), would constitute such a measure of the true amount of matter present.

Now, while it is clear that the acceleration of a negative-energy body in the gravitational field attributable to a local matter inhomogeneity (such as the gravitational field which exists on the surface of the Earth) would not be that which is shared by all objects made of positive-energy matter, experiments provide very strong constraints on the degree of violation of the equivalence principle and to date there is, in fact, no evidence at all that any such violations have ever occurred when systems of various different compositions are utilized. However, I did say, in a previous section, that negative energy is as common as bound systems of particles such as atomic nuclei and molecules, due to the negative energy of their attractive force field. Why, then, do we never observe an altered level of resistance to gravitational acceleration? We may, for example, consider atomic nuclei formed of many protons and neutrons bound together by the strong nuclear interaction, with various measures of negative force-field energy, associated with various configurations, involving a variable number of component particles. It may then appear that the gravitational acceleration of such bound systems should be reduced by the negative value of the energy of the field, while the inertial resistance would be proportionately larger, as the absolute inertial masses attributable to the component particles and the force field would not cancel out like the gravitational masses. If we measured the acceleration of a whole body composed of one such type of nucleus on the surface of the Earth and compared it with the acceleration of another body made of another kind of

nucleus, containing a lesser proportion of such negative energy, we may then expect to discern a difference. But it appears that this is precisely what the experiments discussed above rule out to a very good degree of precision. Shall we then once again abandon everything and conclude that negative energy, even though it is definitely present in bound systems, must be described in a non-relational manner (so that the sum of forces associated with positive and negative inertial masses is allowed to cancel out like those associated with gravitational mass)?

It must be understood that, in fact, this conclusion would constitute a theoretical problem as grave as apparently is the empirical difficulty revealed by the absence of differences in the acceleration of various bound systems. But can we ever hope to solve a problem by creating a 'new' one and assume that, despite all indications to the contrary, the latter difficulty is not real, simply because it only affects consistency on a more general level? This is not the path I chose to follow, because I realized that, despite what is often suggested, there is simply no reason to expect the kind of violations of the principle of equivalence which are described here, even if inertial forces do not cancel out when we consider two masses with opposite signs. What is wrong, I believe, with traditional assumptions is that, when we are considering a bound system and its force field, we assume that we have two masses with opposite signs, while what we really have is one single mass with one overall magnitude and one polarity, both from the viewpoint of inertia and from that of the response to local gravitational fields. What motive do we have, then, for considering that there could be independent contributions to the mass of a bound system (inertial or otherwise) when, in fact, the energy of the subsystems forming it (in particular the particles mediating the attractive force fields) could not be measured independently, given that they arise as virtual processes which do not even have classically well-defined physical properties?

It is a fact that the particles mediating an interaction are virtual and as such, exist merely by virtue of quantum uncertainty, which allows them to carry energy, but only for a time that is short enough that this energy cannot be determined. The virtual particles involved in giving rise to interactions must then be considered unobservable, even if only because, to actually establish their presence in any one particular instance would require a time length greater than the duration of the exchange process. But, under such circumstances, how could we be talking about an *independent* contribution of those particles to the energy or the mass of the bound systems in which

they materialize? I think that this would, in effect, be non-sense and that it must be recognized that any component of a bound system whose physical properties cannot be directly and independently observed does not contribute independently to any of the properties associated with the mass of the system as a whole, when those are actually measured. Failure to understand this decisive requirement would mean that we again allow one more inconsistency to obscure our conception of negative energy in a way that could only be made acceptable by rejecting one or another of the fundamental constraints identified above. In the present context, this could not even be avoided by assuming that negative energy does not exist at all, because the issue is no longer merely about deciding if negative energy exists, but about determining its properties in a context where we must definitely accept that it is occurring.

There is no contradiction here, because there is definitely a negative contribution to the energy of bound systems, only this energy contribution cannot be independently measured in any specific case and this is the crucial distinction we must take into account when estimating the absolute inertial mass of such a system. Thus, the difference between the situation described above of the two superposed opposite-mass objects with large absolute inertial masses and that of a composite system with absolute inertial mass smaller than that of its constituent particles is that, in the former case we are actually dealing with two independent systems, which may be interacting only negligibly with one another, while in the latter case we have one single bound system, which is physically different from the sum of its parts and to which must therefore be associated one single combined measure of mass, gravitational and inertial. In any case, the fact that we do not observe violations of the principle of equivalence for bound systems whose observable total energy is positive confirms that this conclusion is appropriate.

2.13 A few other misconceptions

Before finishing this discussion concerning the potential problems facing a theory of negative-energy matter I would like to provide arguments to the effect that a few other problems which are often associated with the possibility that there could exist gravitationally-repulsive matter are actually of no concern, because they are significant only in the context of a traditional

conception of negative energy and gravitational repulsion⁷. It is nevertheless important for me to discuss those issues, because I have come to realize that the perception of negative energy as being associated with all sorts of strange phenomena that defy common sense is responsible, more than anything else, for making the perfectly acceptable idea of negative-energy matter look like a pseudo-scientific concept without any relevance to physical reality. I will thus try to make clear that what is wrong is not the hypothesis of matter in a negative energy state, but merely the current assumptions regarding what would be the properties of such matter.

One of the problems I would like to discuss arose as an outcome of the first attempts at finding an interpretation for the negative energy states which were predicted to occur by relativistic quantum theories. Indeed, when the existence of antimatter was experimentally confirmed, it was suggested that this kind of matter may perhaps actually give rise to antigravity, in the sense that antimatter would experience repulsive gravitational forces in the presence of ordinary matter. But only theoretical arguments could be given to disprove this possibility when it was first suggested, because no experiment had yet been performed to demonstrate that antimatter would not fall upward in the gravitational field of the Earth. One of those arguments was based on the recognition that if antimatter was to repel or be repelled by ordinary matter, this would allow perpetual motion machines to be built that would extract more energy from a process than was initially available. Indeed, under such circumstances, it would take no energy to slowly raise a particle-antiparticle pair in the gravitational field of our planet (because there would be as much gravitational repulsion as attraction). But when this would be accomplished, the pair could be made to annihilate and the positive energy of the photons so produced could fall back to a detector on the ground where they would be measured as carrying more energy than the pair initially had as a consequence of the frequency increase to which the positive-energy photons would be submitted on their way down (this would

⁷It is not possible to provide a detailed review of all the papers which claim to offer a proof that gravitationally-repulsive, negative-energy matter cannot exist in our universe, but I can assure the reader that, even though I have carefully analyzed many of the so called ‘theorems’ concerning the positivity of energy, I have never found any that does not contain one or another implicit or explicit assumption which would not apply to the kind of approach developed in this report and which invalidates them as theoretical arguments against the possibility of developing a consistent model based on the assumption that matter is allowed to occupy the available negative energy states.

be allowed in the context where the energy of the gravitationally repelled antiparticle would nevertheless be assumed to be positive relative to the forward direction of time, so that the annihilation process is allowed to produce positive energy radiation). It would then seem that energy can be freely produced if antimatter ‘falls’ up.

I think that this argument is perfectly valid, only it cannot be used to justify the rejection of anomalous gravitational interactions in general, but rather simply means that, given that antimatter does not have negative energy (as observed in the forward direction of time), then it should not be expected to be submitted to anomalous gravitational forces. Now, could the same experiment be performed with negative energy (actually negative action) antimatter and then what would it mean for energy conservation? The answer to that question is to be found in the developments introduced in previous sections, while solving other aspects of the problem of negative energy states. First of all, it must be understood that, given that there are no interactions between positive and negative-energy matter, other than the indirect repulsive gravitational interaction which I have already described, it seems that it would be much more difficult to raise a pair of opposite-energy particles together in the gravitational field of a planet without doing work on at least one of them. Yet, this may not constitute an insurmountable difficulty, because it is possible to imagine arrangements which would allow a negative-energy body to achieve the task of raising a positive-energy body in the gravitational field of a positive-energy planet by making use of the indirect, repulsive gravitational forces existing between the two bodies (which could also be composed of matter with opposite charges). But, in fact, the same limitation concerning the absence of any direct interaction between opposite-energy particles would also imply that it is not possible to make a pair of opposite-action particles to annihilate. However, other means would probably exist for harvesting the energy contained in those particles, so that this limitation does not really constitute a decisive constraint that would allow to rule out the kind of processes discussed here.

The real difficulty for any incipient free-energy harvesters would actually arise from the fact that, in the context of a concept of gravitationally-repulsive, negative-energy matter such as the one I have proposed, even if the experiment described above could be performed with a pair of opposite-action particles with opposite charges, upon annihilating one another the particles would release no energy at all. Indeed, if the particles have equal, but opposite energies initially and they do not gain or lose any kinetic energy

as a result of their ascension, then their respective final energies would still be equal in magnitude. As a consequence, even if those particles could annihilate one another (which is not the case, as I have explained in section 2.9), no energy would be released, so that there would be no photons to fall back toward the surface of the planet with a net gain of energy. It must be clear, however, that it is not the limitation imposed on the annihilation of opposite-action particles which alone prevents the production of free-energy, because we could arrange things so that the positive-energy particle annihilates with a positive-energy antiparticle already in place at the destination point, while the negative-energy antiparticle would annihilate with a negative-energy particle already in place. But if the positive-energy photons produced by the annihilation of the positive-energy particles could actually gain positive energy while falling back to a detector on the ground, the negative-energy photons produced by the annihilation of the negative-energy particles, for their part, would lose negative energy while reaching the same detector and would therefore end up with less negative energy than they would have had if the negative-energy particles had been submitted to annihilation before rising to a higher altitude. Thus, while positive radiation energy would be gained during such a process, negative radiation energy would be lost and this means that no work can be performed in such a way.

In order to better understand the significance of the changes involved, we can consider the variations occurring in the potential energy of two opposite-energy bodies as they are raised in the gravitational field of a positive-energy planet. From this more general perspective what would be observed, in effect, is that any potential energy that would be gained by one of the two bodies (the one that was actually lifted by the other) would necessarily be lost by the other body, thereby preventing any useful energy from being produced in the course of such a process. Indeed, while the positive-energy body would gain positive potential energy, the negative-energy body would lose negative potential energy. Now, this may seem to imply that a forbidden net increase of (positive) energy can be obtained despite the fact that no work would have been done to take the system to its final state. Yet, as I have explained in a preceding section, this variation is not significant, because any change in the energy of matter resulting from an interaction between positive- and negative-energy bodies is compensated by an opposite change in the energy of the gravitational fields attributable to the local changes occurring in the positive and negative portions of vacuum energy, as a result of those interactions.

What must be understood here is that, even if a positive change may occur in the potential energy of matter, this would not mean that we have gained the ability to perform more work, as would be required to produce perpetual motion, because what the loss of negative potential energy by the negative-energy body means is precisely that there was a loss of useful energy (energy that could be used to do work) for that object during the process by which it would have performed work to raise the positive-energy body and increase the ability of this positive-energy body to perform work. In other words, despite the net gain in potential energy for the pair as a whole, the ability to do work would not have increased, because the negative-energy body, having been raised by the repulsive gravitational field it experiences, would have exhausted its ability to perform work (even though its kinetic energy would remain unchanged), which is precisely what its loss of *negative* potential energy implies, because indeed the object would have lost energy of the same sign as its own and therefore would actually end up with less energy available to perform work after the lifting process has occurred. The gain in useful energy by the positive-energy body would actually have been provided by the negative-energy body which would have lost its own useful energy and in fact, if the usual friction and other degradation of energy had been taken into consideration, it should be observed that the positive-energy body would have gained less useful energy than the negative-energy body would have lost, thereby precluding any perpetual motion from being achieved.

The fact that positive energy seems to have been created, on the other hand, is a simple consequence of the fact that the process discussed involves an indirect gravitational interaction between the two opposite-energy bodies and between the negative-energy body and the positive-energy planet during which the total energy of matter may indeed vary, as I remarked above, given that it is compensated by an opposite variation in the energy of the gravitational field attributable to the local changes occurring in the energy of the vacuum as a consequence of those indirect gravitational interactions. No additional difficulty is involved here and therefore it seems that the perpetual motion argument against gravitational repulsion cannot be considered significant, other than as an argument against the possibility of an anomalous gravitational interaction between ordinary matter and ordinary antimatter.

A more exotic and hypothetical phenomenon, which according to certain accounts could have interesting practical applications, but which would raise serious problems from a theoretical viewpoint, given that it may provide the means of achieving faster-than-light space travel and therefore, also, time

travel, is that of wormholes. It is often thought that wormholes would naturally occur in the presence of some types of black-hole singularities and may allow remote regions of space to be directly connected in some way, so that traveling through such wormholes would enable to bypass the limitations associated with the passage of time experienced under normal circumstances when traveling over such long distances at slower-than-light velocity. It is not clear exactly what regions of space could be connected in such a way, or if we are really talking about connecting regions of our own universe, but if we leave aside those uncertainties, then it would seem that all that is required for unlocking the potential of faster-than-light space travel is the existence of traversable versions of such hypothetical shortcuts through space and time. What must be provided, therefore, is a means to maintain the ‘throat’ of a wormhole open for a long enough period of time that space travelers can safely traverse it, despite the tendency for the matter configurations involved here to collapse under the effect of the gravitational attraction exerted by the singularity. The idea is that gravitationally-repulsive, negative-energy matter (often called exotic matter) may allow to achieve that goal, given that it could be used to exert a gravitational repulsion that would compensate the attraction exerted by the spacetime singularity at the center of the black hole. But again, when we look at the details of such proposals, it becomes clear that the conditions necessary for achieving the desired results are incompatible with a consistent notion of negative-energy matter. That may not be good news for science fiction lovers, but if I’m right negative-energy matter could never be used to achieve such a goal.

To help identify what’s wrong with current expectations, I would suggest that we ask how it is exactly that negative-energy matter could be brought, not just inside some black hole, but toward the point of maximum density of positive-energy matter (the singularity), despite the enormous gravitational repulsion that this positive-energy matter would exert on the exotic matter? It should be clear that it is merely because we traditionally assume that negative-energy matter would be attracted by a positive-energy black hole and its singularity, even while it would repel it, that this appears to constitute an achievable goal. But the truth is that any negative-energy matter approaching a large concentration of positive-energy matter, such as an ordinary black hole, would be submitted to repulsive forces as large as those maintaining positive-energy matter trapped inside the same black hole. In this context, the only way by which negative-energy matter could find itself inside the event horizon of a positive-energy black hole would be by having

already been present inside the region destined to collapse into that positive-mass black hole, before it formed. But even if that was to happen, there is no way that the negative-energy matter could be made to remain near the black-hole singularity, where repulsive forces would be the strongest. This situation is simply unstable and given that stability is precisely what is required for a traversable wormhole to exist, we must recognize that negative-energy matter could not provide the necessary element for allowing spacetime singularities to be used for faster-than-light space travel and time travel. The possibility that the kind of phenomenon discussed here could actually have been used for achieving theoretically problematic, causality-violating processes may seem far-fetched, but I think that it is nevertheless important to show that, even under such extreme conditions, there is no reason to expect that the existence of negative-energy matter could facilitate such an outcome (in section 5.11 I will explain why it is exactly that closed time-like curves, of the kind that could have been allowed by the existence of traversable wormholes, are to be considered problematic and it will become clear that the difficulty is not that they may allow a time traveler to alter his or her own past).

The same argument I have used to rule out the possibility of engineering traversable wormholes can also be utilized to solve a more down-to-earth problem that is not often discussed, but which would contradict one of the most unavoidable constraints applying to the evolution of irreversibly evolving physical systems, such as black holes. The problem is that negative-energy matter, as it is traditionally conceived, could be used to reduce the mass of a black hole and therefore, also, the area of its event horizon. This could be achieved by simply throwing negative-energy matter into a black hole, which would presumably absorb it, given that negative-energy matter is usually assumed to be gravitationally attracted by a positive-energy black hole. This would be possible, even if negative-energy matter repels a positive-mass black hole, because we could throw negative-energy particles in small amounts and their gravitational fields would be too small to resist the much larger gravitational attraction of the black hole. But the surface area of a black hole has been shown to constitute a measure of the entropy of such an object, so that reducing the area of the black hole is similar to reducing its entropy. Again, however, if we reject the traditional conception of negative-energy matter, the problem does not exist, because a negative-energy particle cannot even get near a positive-energy black hole without experiencing extreme gravitational repulsion, so that it certainly cannot be absorbed by the object, as would be necessary for reducing its mass and the area of its event horizon.

If negative energy states are to be considered a true possibility, then the fact that the traditional concept of negative-energy matter would allow such violations of the second law of thermodynamics, while the alternative approach developed in this report would forbid them, constitutes a strong indication to the effect that this latter proposal is more appropriate.

In fact, we are dealing with a much more general problem in this case, because, from a traditional viewpoint, it is actually assumed that when negative-energy radiation would come into contact with positive-energy matter (not necessarily a black hole), it could be used to withdraw positive thermal energy from this matter (as if it was providing negative heat), therefore raising the possibility of allowing temperature to decrease in both a positive-energy system and the negative-energy system with which it is interacting, without any heat being released in their environment, which again would violate the second law of thermodynamics. But the situation would only be worse if we also assumed that the absorption of positive-energy radiation by negative-energy matter would itself allow negative thermal energy to decrease (toward less negative values) in a negative-energy system. Of course, given that, from my viewpoint, negative-energy radiation cannot even come into contact with positive-energy matter, the possibility raised here appears to be mostly irrelevant from a practical viewpoint. We may nevertheless examine the situation which would arise following an exchange of thermal energy between positive- and negative-energy systems occurring as a consequence of the indirect repulsive gravitational forces they exert on one another.

The conclusion we must draw, in such a case, is that negative energy is not equivalent to negative heat for a positive-energy system. Indeed, according to my conception of negative-energy matter, from the viewpoint of a positive-energy observer, kinetic energy is exchanged between opposite-energy particles as if it was a positive-definite quantity. This is allowed given that the energy of matter particles is not conserved independently from certain opposite contributions to gravitational potential energy which vary as a result of local changes in the energy of the vacuum produced by those interactions, as I explained in section 2.11. But the fact that the kinetic energy of matter appears to be conserved as if the interacting particles all had the same sign of energy means that thermal energy itself can only be gained as a positive-definite quantity by positive-energy systems, or equivalently as a negative-definite quantity by negative-energy systems, even when the exchange involves opposite-energy systems. Thus, when heat is provided by a negative-energy system it can only raise the positive temperature of

a positive-energy system (as if positive thermal energy was provided) and the same is true for the heat provided by a positive-energy system to a negative-energy system, which can only raise the negative temperature of the negative-energy system toward more negative values (as if negative thermal energy was provided by the positive-energy system). It is necessary to assume, in effect, that temperature, as a measure of the local intensity of thermal energy, is negative for negative-energy matter, even under normal circumstances (when the number of possible microscopic states is allowed to rise without limits as the negative energy of a system rises), given that when the energy of matter rises (into positive or negative territory), the entropy of matter itself rises, so that if such a change takes place as a result of the absorption of negative heat (as may be the case for a negative-energy system), then it can only mean that the temperature of the system in which those changes are taking place is negative.

Thus, we have no reason to expect that even the indirect gravitational interactions between opposite-energy systems could be used to transform useless forms of energy into more useful forms and in such a way reduce the entropy of matter. Negative thermal energy cannot reduce the temperature of a positive-energy system any more than positive thermal energy could diminish the magnitude of the temperature of a negative-energy system. The temperature of a positive-energy system can only be reduced through the emission of positive heat, just like the temperature of a negative-energy system can only be reduced (toward less negative values) when it releases negative heat. For a positive-energy system to lose thermal energy at the expense of a negative-energy system, the magnitude of its temperature must be larger than that of the negative-energy system and under such conditions the magnitude of the temperature of the negative-energy system would be raised by an amount proportional to that which is lost by the positive-energy system, as when all temperatures are positive. What must be understood is that transferring heat from a negative energy source to a positive-energy system is not equivalent to removing positive heat from that system. In fact, it rather seems that adding heat from a negative-energy system to a gas of positive-energy matter would actually raise its temperature (unlike most people considering the possibility of the existence negative-energy matter usually assume). This is all a consequence of the fact that negative kinetic energy can be turned into positive kinetic energy and vice versa, even when energy is assumed to be conserved, as I previously explained.

It appears, therefore, that the positive thermal energy of a gas of positive-

energy particles can actually be raised through contact with a gas of negative-energy particles, when the magnitude of the negative temperature of the negative-energy gas is larger than the positive temperature of the positive-energy gas, because thermal energy is a measure of the average kinetic energy of gas molecules and this energy would become more evenly distributed between the two gases (independently from energy signs), if they could be put into contact through the indirect gravitational interaction. In this context, it transpires that all that matters from a thermodynamic viewpoint, for a positive-energy system which interacts with a negative-energy system, is whether negative thermal energy is actually gained or lost by the negative-energy system and not whether the sign of this energy is positive or negative. The rule that emerges is that when heat is lost by a negative-energy system in contact with a positive-energy system, it is gained as positive heat by the positive-energy system, while when heat is lost by a positive-energy system in the same situation, it is gained as negative heat by the negative-energy system.

Once again, the traditional expectation can be seen to arise from a misconception. You should take note, however, that I'm not just trying to debunk myths here. The opposite conclusion, that a low temperature gas made of positive-energy particles would be cooled even further upon contact with heat from a negative-energy gas, regardless of the magnitude of the temperature of this negative-energy gas, and the above discussed assumption that the mass of a positive-energy black hole could be reduced through the absorption of negative-energy matter, would constitute serious problems for a gravitational theory integrating the concept of negative-energy matter. There are very strong motives behind my desire to demonstrate that the possibility of such entropy decreasing processes can be rejected and they are actually related to those which one might raise against the above discussed possibility of causality-violating processes. I will explain what is the profound significance of the results discussed here in the multiple sections of chapter 4 that deal with the problem of time irreversibility.

2.14 An axiomatic formulation

Before I complete the process of integration of negative-energy matter to classical gravitation theory, I would like to provide formal statements of each of the significant rules I have derived in relation to this issue and which were

discussed in the previous sections of the current chapter. Basically, there are ten fundamental rules which clarify the situation regarding the nature and the behavior of negative-energy matter itself, as well as the behavior of positive-energy matter in the presence of negative-energy matter. Those rules actually constitute the axioms on which a generalized classical theory of gravitation can be based. The axioms are legitimized by the fact that they have been shown to be necessary on the basis of both logical consistency and agreement with experimental facts and thus we may appropriately refer to them as principles. The first principle is the most fundamental and a recognition of its validity opens the way for a derivation of all the other results. The formal statement of this principle goes like this:

Principle 1: The distinction between a positive-energy particle and a negative-energy particle (propagating negative energy forward in time) can only be defined by referring to the difference or the identity of the energy sign of one particle in comparison with that of another, so that the sign of energy or mass has no absolute meaning.

From a gravitational viewpoint, this principle is satisfied when positive-energy particles are submitted to mutual gravitational attraction among themselves (as we observe), while negative-energy particles (actually negative-action particles) also attract one another gravitationally and positive- and negative-energy particles repel one another, as a consequence of the indirect gravitational interaction which actually originates from an uncompensated gravitational attraction between matter of one energy sign and that portion of vacuum energy with the same energy sign. Compliance with this rule means that for a positive-energy particle, a negative-energy particle should be physically equivalent to what a positive-energy particle is for a negative-energy particle. This property will be decisive for deriving the observer-dependent generalized gravitational field equations that will be introduced later.

Another rule applies only in the classical Newtonian context where mass is a significant concept, but given that it allows to derive the rules which must also be obeyed in a general-relativistic context it is necessary to mention it as a basic result. It simply amounts to recognize that:

Principle 2: When mass is reversed from its conventional positive value, both gravitational mass and inertial mass are reversed and together become negative.

This is actually equivalent to assume that there is indeed only one physical property to which we may refer to as being that of mass and that there cannot be any arbitrary distinction between gravitational and inertial mass.

While principles 1 and 2 are for the most part theoretically motivated, the next principle is both theoretically and observationally motivated. Indeed, principle 3 arose as the unavoidable consequence of an analysis of the relationship between the attractive or repulsive nature of a field of interaction and the sign of the energy classically contained in this field, but it is also a necessary requirement of the fact that we do not observe any negative-energy matter, despite the fact that the existence of such matter appears to be allowed from a theoretical viewpoint. The third principle therefore is the following requirement:

Principle 3: There are no direct interactions of any kind (either gravitational, electromagnetic, or nuclear), mediated by the exchange of bosons of interaction, between positive- and negative-action particles (propagating positive and negative energies forward in time, respectively).

Compliance with this principle means that negative-energy observers would also be prevented from directly observing positive-energy matter.

Another important result was discussed at length in a previous section of this chapter, where its validity was shown to be unavoidable despite the fact that it appears to contradict some assumptions which are usually considered to be irrefutable. This result simply states that:

Principle 4: A void of limited size that develops in an otherwise uniform matter or energy distribution gives rises to uncompensated gravitational forces which are the opposite of those which would otherwise be produced by the matter or energy that is missing.

The effect it describes is the consequence of an alteration (caused by the presence of some local void) in the equilibrium of gravitational forces applying on any particle and due to its interaction with all the other particles in the universe (with which this particle actually interacts). The importance of this principle becomes clear when we consider its significance in the context where the uniform energy distribution is actually the distribution of vacuum energy and it is recognized that principle 5 below applies.

The following principle is probably the most decisive after principle 1, given that it is the result that allows the whole concept of negative-energy matter to have a significance despite the validity of principle 3 and the absence of direct interactions between positive- and negative-energy particles. It states that:

Principle 5: Locally, the presence of negative-energy matter is equivalent to the absence of an equal amount of positive energy from the vacuum, while the presence of positive-energy matter is equivalent to the absence of an equal amount of negative energy from the vacuum.

As I explained in section 2.8, those equivalences constitute the particularity that allows opposite-energy bodies to exert gravitational forces on one another despite the absence of direct interactions between them, simply because, according to principle 4, voids in a uniform, positive energy distribution do have an indirect influence on positive-energy matter, despite the fact that those voids are actually equivalent to the presence of negative-energy matter with which positive-energy matter does not directly interact. In fact, it would be appropriate to assume that the presence of matter is the consequence of a local absence of both energy and non-gravitational charges from zero-point vacuum fluctuations, as I mentioned in section 2.8. But, again, even though such an absence of charges is equivalent to the presence of opposite-sign charges, this is without any consequences, given that while negative-energy matter cannot interact directly with positive-energy matter, it does interact with both the positive and the negative charges present in the electrically neutral vacuum, which means that the effects of all those interactions cancel out, even in the presence of voids in the negative-energy portion of the vacuum and the same argument applies for positive-energy matter and the positive-energy portion of the vacuum.

Now, even in the context where we assume the existence of a symmetry between positive- and negative-energy matter, principle 5 would require that it is, in fact, only the inhomogeneities (either overdensities or underdensities) present in the negative-energy matter distribution which can affect the gravitational dynamics of positive-energy matter, while it is only the inhomogeneities present in the positive-energy matter distribution which can affect negative-energy matter. This is because, as previously discussed, the void in the positive portion of vacuum energy that is equivalent to a totally homogeneous distribution of negative-energy matter would leave no

surrounding positive vacuum energy to produce an uncompensated gravitational attraction that would be equivalent (according to principle 4) to the gravitational repulsion otherwise attributable to the negative-energy matter and the same is true concerning a homogeneous distribution of positive-energy matter from the viewpoint of negative-energy matter. An additional principle thus emerges that expresses this limitation applying on principle 5. It amounts to assume that:

Principle 6: Only (positive and negative) density variations in an overall homogeneous, cosmic-scale distribution of negative-energy matter can be assumed to exert gravitational forces on positive-energy matter.

Of course, a similar limitation would also apply, which would actually express the absence of gravitational forces on negative-energy matter from a totally smooth and uniform cosmic-scale distribution of positive-energy matter.

A further particularity could be derived from the already stated principles, but I will provide it as an additional specific rule, because it may not be obvious that it applies in the context where principles 3 and 6 are assumed to constrain the interaction between positive- and negative-energy matter. This ordinance states that:

Principle 7: Despite its energy sign and its assumed uniformity, the negative-energy portion of the vacuum does exert the gravitational influence it should have on positive-energy matter.

As I previously explained, this deduction (which would also apply to the positive-energy portion of the vacuum from the viewpoint of negative-energy matter) follows from the fact that the restriction that applies on the interaction of positive- and negative-energy matter does not prevent positive-energy matter, when it is conceived as voids in the negative-energy portion of the vacuum, from having an influence on that very portion of the vacuum in which the voids are present, just as voids in a matter distribution do exert an influence on this matter. Also, the fact that the energy of the vacuum can be expected to be uniformly distributed, does not restrict the influence of the negative portion of it from influencing positive-energy matter, simply because we are not dealing, in this case, with negative-energy matter and the negative energy of the vacuum itself cannot be considered as being equivalent to a void in this very vacuum, so that whatever the extent of the distribution

of negative energy involved it would still exert its influence on both positive- and negative-energy matter, unlike a uniform distribution of negative-energy matter.

In a previous section I have explained that a consequence of principle 1, in the context where principle 2 (regarding the negativity of the inertial mass of a negative gravitational mass) is considered to apply, is that the usual assumption that reversing all mass (gravitational and inertial) would allow to maintain agreement with the equivalence principle (as it is traditionally conceived) is wrong. Therefore, only an altered principle of equivalence between acceleration and a Newtonian gravitational field can remain valid. The additional condition applying on the equivalence principle would be the following:

Principle 8: The equivalence of the effects of gravitation and acceleration does not apply merely locally, but merely for one single elementary particle (in a given location and with a given sign of mass or energy) at once.

What remains true, in this context, is that the motion of bodies in a gravitational field does not depend on any physical properties of those bodies other than the sign of their mass or energy and this is what will allow the essence of the current theory of gravitation to be retained, while accommodating a consistent concept of negative-energy matter.

Another rule must be obeyed in the context where negative-energy matter is governed by principle 1 above and where the appropriate inertial behavior of this type of matter, which can be derived from principles 2 and 6, is assumed to apply (which actually means that the inertial response of negative-mass or negative-energy bodies to a given force is the same as that of positive-energy bodies, as I explained before). This rule would not apply if the traditional assumptions regarding the inertial response of negative-energy or negative-mass bodies were valid. But given that I have argued that those assumptions are problematic and cannot be justified, then it seems that, even traditionally, we would have a problem if we were not taking the following experimentally motivated principle into account.

Principle 9: When the negative contribution of a field of interaction to the energy of a bound physical system with overall positive energy cannot be independently and directly observed, only the diminished total energy of the bound system contributes to its (previously defined) absolute inertial mass.

Again, this is also valid for bound physical systems with overall negative energy, for which we may say that, when the positive contribution of a field of interaction to the energy of the bound system cannot be independently and directly observed, only the diminished (less negative) total energy of the bound system contributes to its absolute inertial mass. It must be remarked that the validity of this rule does not mean that the opposite contribution to the total energy of a bound system by the field of interaction responsible for the mutual attraction of its component particles cannot be well-defined, only that if it cannot be isolated and independently measured then it also does not independently contribute to the inertial properties of the whole system.

One last constraint is observed to apply when negative energy states are allowed to be occupied (can be propagated forward in time). While this rule is theoretically motivated, I originally derived it based on purely phenomenological arguments. It is the following:

Principle 10: A particle cannot reverse its direction of propagation in time on a continuous particle world-line without also reversing its energy and equivalently, a particle cannot reverse its energy on a continuous particle world-line without also reversing its direction of propagation in time.

Here by ‘negative energy’ I mean negative energy relative to the true (even though relationally defined) direction of propagation in time, as in the case of the positron as a negative-energy electron propagating its negative electric charge backward in time. This rule is equivalent to assume that it is impossible for pairs of opposite-action particles to be created out of nothing, or to annihilate to nothing, which is an indirect consequence of principle 3 (regarding the necessary absence of direct interactions between opposite-action particles).

The ten principles enunciated above embody the essence of the insights I have gained through an analysis of the problem of negative energy in light of the requirement of relational definition of the physical properties of mass and energy signs. They will now be used to help derive a generalized formulation of the gravitational field equations that will allow to describe the motion of particles with a given sign of energy in the gravitational field of an object with opposite mass or energy.

2.15 Generalized gravitational field equations

I previously indicated that equations would be scarce in this report. But the point has now been reached where it is absolutely necessary to provide some level of quantitative detail regarding the manner by which the concept of negative energy that was developed in the preceding sections of the current chapter is to be integrated into a classical theory of gravitation. The objective I'm seeking here, though, is not to provide a complete treatise on the subject, but merely to introduce the modified gravitational field equations which constitute the core mathematical structure of the generalized theory that emerges from the alternative set of axioms introduced in the preceding section. The essential requirement that must be imposed on a formulation of the gravitational field equations, in the context where the principles enunciated in the preceding section are to govern the behavior of negative-energy matter, is that the gravitational field attributable to a given local source is not to be considered attractive or repulsive depending only on the sign of energy of the source. This can be satisfied by assuming that the gravitational field experienced by a negative-energy particle and attributable to a given matter distribution is actually different from the one experienced by a positive-energy particle. In such a context only the difference or the identity between the energy signs of two masses would be physically significant to determine the character of their gravitational interaction, so that any one mass could be considered to have positive energy, while masses with an opposite energy sign would then have to be the ones to which a negative energy is to be attributed. But the choice of which of two opposite-energy bodies has positive energy is itself completely arbitrary.

Thus, an observer formed of matter with a given energy sign is free to attribute positive energy to particles with the same sign of energy, even though an observer formed of matter of opposite energy sign may attribute a negative energy to the exact same matter. The only requirement is that the value of the gravitational field (which, in a general-relativistic theory, is associated with the metric properties of space and time) always be adjusted as a consequence of the arbitrary choice which is made regarding the attribution of energy signs to various objects. There is, however, a natural choice for the attribution of energy signs by a given observer, which consists in assuming that matter with the same sign of energy as that of the observer itself is always to be considered positive by this type of observer. The viewpoint under which what we traditionally call positive-energy mat-

ter actually has positive energy is therefore the natural viewpoint of what we traditionally consider to be a positive-energy observer, while the viewpoint under which what we traditionally call positive-energy matter actually has negative energy is the natural viewpoint of what we would traditionally consider to be a negative-energy observer. When this convention is adopted, we can write observer-dependent gravitational field equations which replace the traditional equations. According to this alternative formulation, the motion of matter with a given energy sign is determined by the gravitational field associated with observers having the same energy sign. The gravitational field, therefore, varies as a function of both the energy sign of the sources and the energy sign of the particles submitted to it, so that only the difference or the identity between the energy sign of the source and that of the matter submitted to the observer-dependent gravitational field determines the repulsive or attractive nature of the interaction.

In a relativistic context, the observer dependence of the gravitational field would imply that observers of opposite energy signs actually experience space and time in a different way. But despite the awkwardness of this possibility from the perspective of our conventional perception of spatial relationships, from a mathematical viewpoint this requirement does not constitute an insurmountable difficulty. We merely have to assume two spaces, related to one another by the fact that the same unique set of events is taking place in both of them, but which may nevertheless have distinct metric properties, in the sense that the events which are taking place in the universe are separated by space and time intervals which are dependent on the energy sign of the observer. Indeed, as I mentioned before, the equations which will be proposed here merely constitute a generalization of the existing mathematical framework of relativity theory and we will therefore be in familiar territory. I'm, in effect, assuming that the reader already has a proper understanding of the current general-relativistic theory of gravitation and of the physical significance of the various mathematical objects which are relevant to the conventional formulation of this theory. Also, given that attempts at formulating a relativistic theory of gravitation that would allow for the existence of observer-dependent gravitational fields were the subject of earlier publications by various authors and since it would be pointless to simply reproduce what has already been discussed elsewhere, I will leave to experts the task of introducing the general framework in which the developments I will propose are to be formulated and concentrate instead on describing the essential, distinctive mathematical features unique to the theory I'm proposing.

This choice is appropriate, despite the fact that the approach I favor involves several distinctive aspects, because the most general features of the kind of framework involved are not dependent on the specific assumptions of the model considered. The reader may refer in particular to a paper published sometime ago by Sabine Hossenfelder [25] in which were introduced meaningful developments essential to any theory according to which the gravitational field is assumed to be dependent on the nature of the matter experiencing it. But keep in mind that even the most suitable of the currently available mathematical frameworks still involves theoretical constructs and assumptions which I would consider inappropriate for the formulation of a fully consistent, generalized, classical theory of gravitation integrating the concept of negative-energy matter and therefore only the general structure provided by those developments must be retained. I will here provide an interpretation of such bi-metric theories that is different from those which were tentatively proposed by the few authors that preceded me and this will have significant consequences which will be reflected in the fact that the final equations at which I have arrived are actually distinct from those which had been proposed until now.

In any case, it must be mentioned that the gravitational field equations which appear in the above cited paper were not the first equations of that kind to have been developed. Gravitational field equations involving conjugate metrics had already been proposed that simply amounted to allow for negative contributions to the stress-energy tensor of matter⁸, while implicitly (but unsatisfactorily) trying to conform to the requirement of symmetry under an exchange of positive and negative energy signs. But even in the more recent publications, no justification has ever been provided for the assumptions on which are based the emerging theories and the only experimental consequences that were derived from those developments actually appeared to disagree with observations or were again unjustified on the basis of the hypotheses which were assumed to characterize the behavior of the gravitationally-repulsive matter. To my knowledge, no author was ever able to recognize the exact nature of the anomalously gravitating matter they sought to describe, or to explain how the various problems related to the existence of such matter could be solved. In fact, none of them even suc-

⁸I became aware of those developments mainly through the early writings of a Frenchman named Jean-Pierre Petit, but given that I have never read any official research publication from him that contains the set of equations to be discussed here, then I will not attempt to provide specific references to his work on the subject.

ceeded in justifying the validity or the superiority of an approach to classical gravitation based on the requirement of exchange symmetry, in comparison with the traditional viewpoint according to which gravitational attraction and repulsion are absolutely defined properties of matter.

Meaningful equations were, nevertheless, derived, which happened to be compatible with the simplest of the conditions I have identified above as characterizing a consistent theory of negative-energy matter. Those equations, therefore, constituted a step forward in deriving a quantitative model for the gravitational dynamics of negative-energy matter, even if they failed to provide a totally appropriate framework and had to be assumed to apply only under particular circumstances, as they were clearly inappropriate to describe the early phases of cosmic evolution. In any case, the equations which were initially proposed were of the following form:

$$\begin{aligned} R_{\mu\nu} - \frac{1}{2}g_{\mu\nu}R &= -\frac{8\pi G}{c^4}(T_{\mu\nu} - T_{\mu\nu}^-) \\ R_{\mu\nu}^- - \frac{1}{2}g_{\mu\nu}R^- &= -\frac{8\pi G}{c^4}(T_{\mu\nu}^- - T_{\mu\nu}) \end{aligned} \quad (2.1)$$

Here and in what follows G is Newton's constant, c is the speed of light in a vacuum, and the Greek indexes μ and ν run over the four general coordinate system labels (assuming a metric with diagonal elements $+1, +1, +1, -1$ in an inertial coordinate system). The usual notation is used for the curvature tensors $R_{\mu\nu}$ and R experienced by positive-energy observers and for the stress-energy tensor $T_{\mu\nu}$ of what we conventionally consider to be positive-energy matter, as measured by a positive-energy observer. The curvature tensors experienced by negative-energy observers are for their part denoted as $R_{\mu\nu}^-$ and R^- , while the stress-energy tensor of what we would conventionally consider to be negative-energy matter, as measured by a negative-energy observer, is here denoted as $T_{\mu\nu}^-$. The first of those two equations can thus be used to determine the geodesics followed by positive-energy particles, while the second determines the geodesics followed by negative-energy particles. Here, all stress-energy tensors would have to be assumed to correspond with positive-definite energy densities if it was not for the negative sign in front of the second stress-energy tensor on the right-hand side of each equation, which allows for a negative contribution to the total stress-energy tensor of matter that is dependent on the particular measure of the sign of energy associated with one or the other type of observer. The negative sign of stress-energies can thus be attributed alternatively to what we would usually consider to be

negative-energy matter and to what we usually consider to be positive-energy matter.

This actually means that what appears to be negative-energy matter to a conventional positive-energy observer would really be positive-energy matter for an observer we would normally consider to be a negative-energy observer, while what appears to be positive-energy matter to a positive-energy observer would really be negative-energy matter for an observer usually considered to be made of negative-energy matter. Therefore, all energy signs must now be assumed to depend on the energy sign of the observer, which is itself assumed positive as a matter of convention. The viewpoint I previously identified as equivalent to a reversal of the sign of mass and according to which it is the gravitational field itself (represented here by the curvature tensors) which actually varies, while the sign of mass (replaced here by the sign of energy) of the observer which experiences that gravitational field is to be considered positive definite, is thus applied and this is certainly appropriate given that it gives rise to equations of the simplest form. It is because there are two different measures for the gravitational field, associated with the two different ways by which the positive and negative contributions to the total energy of matter can be attributed, that there are two equations for the gravitational field, instead of the single one that is usually considered. Otherwise, however, those equations are fairly conventional and were certainly the most straightforward that one could derive for a bi-metric theory, as they were the closest to Einstein's own equation that one could propose.

The fact that, in the context of those equations, the sign of energy contributed by a given mass must now be assumed to depend on the sign of energy which we would normally attribute to the observer determining the associated gravitational field has important consequences. Indeed, if variations in the gravitational field (which is represented by the curvature tensors) are to compensate variations in the stress-energy of matter (as the general covariance of the equations require) then it means that the gravitational field attributed to some matter can actually be either attractive or repulsive depending on the sign of energy of the observer that measures the energy of this matter.

Four situations may therefore arise when we limit ourselves to merely permute the energy signs of a pair of interacting bodies. First, the source of the field could have what we traditionally consider to be positive energy and the field be attractive, because the particle submitted to it also has positive energy. Next, the source of the field could have what we traditionally

consider to be negative energy and the field be repulsive, because again the particle submitted to it has positive energy. Another possibility is that the source of the field could have what we would traditionally consider to be positive energy and the field nevertheless be repulsive, because we consider its effects on what we would traditionally consider to be a negative-energy particle and from which viewpoint the source actually has negative energy. Finally, the source of the field could have what we traditionally consider to be negative energy and the field nevertheless be attractive, again because we consider its effects on what we would traditionally consider to be a negative-energy particle and from which viewpoint the source actually has positive energy. This is certainly appropriate from the viewpoint of the principles identified in the preceding section. But given the insights I had already obtained when I first learned about the mathematical developments which can be used to articulate those requirements, it appeared to me that what the available framework provided was, at best, an incomplete formulation of the gravitational field equations to associate with a theory of negative-energy matter.

To try to address those shortcomings, I thus proposed (in a preprint [26] published in early 2006) the following equations which allowed to express the particularities of the indirect gravitational interaction of positive- and negative-energy matter that I had come to consider as unavoidable:

$$\begin{aligned} R_{\mu\nu}^+ - \frac{1}{2}g_{\mu\nu}R^+ &= -\frac{8\pi G}{c^4}T_{\mu\nu}^+ \\ R_{\mu\nu}^- - \frac{1}{2}g_{\mu\nu}R^- &= -\frac{8\pi G}{c^4}T_{\mu\nu}^- \end{aligned} \quad (2.2)$$

Here $R_{\mu\nu}^+$ and R^+ are simply the curvature tensors experienced by positive-energy observers, while $R_{\mu\nu}^-$ and R^- are the curvature tensors experienced by negative-energy observers. But the stress-energy tensors figuring in the equations I proposed are actually different from those entering the previously mentioned set of equations, despite the similar notation I adopted here, because the $T_{\mu\nu}^+$ tensor encompasses all contributions to the energy and momentum experienced by positive-energy observers, while the $T_{\mu\nu}^-$ tensor encompasses all contributions to the energy and momentum experienced by negative-energy observers and I did assume contributions to those stress-energy tensors which were different from those which had previously been considered in the literature. Thus, when written in a more explicit form, with all the components actually entering the stress-energy tensors on the

right-hand side, the equations I proposed are the following:

$$\begin{aligned} R_{\mu\nu}^+ - \frac{1}{2}g_{\mu\nu}R^+ &= -\frac{8\pi G}{c^4}(T_{\mu\nu}^+ + \check{T}_{\mu\nu}^- - \hat{T}_{\mu\nu}^-) \\ R_{\mu\nu}^- - \frac{1}{2}g_{\mu\nu}R^- &= -\frac{8\pi G}{c^4}(T_{\mu\nu}^- + \check{T}_{\mu\nu}^+ - \hat{T}_{\mu\nu}^+) \end{aligned} \quad (2.3)$$

In this notation $T_{\mu\nu}^+$ is the stress-energy tensor of what is usually considered to be positive-energy matter, as measured by a positive-energy observer, while $\check{T}_{\mu\nu}^-$ is the stress-energy tensor associated with the measure of energy of negative-energy matter (effected by a negative-energy observer) below its average cosmic density (toward the zero-energy level) and $\hat{T}_{\mu\nu}^-$ is the stress-energy tensor associated with the measure of energy of negative-energy matter (effected by a negative-energy observer) above its average cosmic density (away from the zero-energy level). Similarly, $T_{\mu\nu}^-$ is the stress-energy tensor of what we would usually consider to be negative-energy matter, as measured by a negative-energy observer, while $\check{T}_{\mu\nu}^+$ is the stress-energy tensor associated with the measure of energy of positive-energy matter (effected by a positive-energy observer) below its average cosmic density (the difference between this average density and the smaller density of positive-energy matter) and $\hat{T}_{\mu\nu}^+$ is the stress-energy tensor associated with the measure of energy of positive-energy matter (effected by a positive-energy observer) above its average cosmic density.

This formulation of the generalized gravitational field equations allows me to take into account the fact that there are two distinct categories of contributions to the total energy density experienced by positive-energy observers, one positive definite for all densities of positive-energy matter and one that can be either positive or negative depending on the value of energy density of negative-energy matter relative, not to the zero-energy ground state, but to the density of this negative-energy matter averaged over the entire volume of the (observable) universe. Basically, what that means is that the energy measures of the second category of contributions experienced by a positive-energy observer are shifted from the traditional zero point of energy to a lower (more negative) energy level below which energies are negative and above which energies are positive, up to a maximum value which is reached when no negative-energy matter is present at all in the considered location. This redefinition of the measures of energy associated with what we conventionally assume to be negative-energy matter simply amounts to subtract the (time dependent) true, negative, average density of energy of this matter

(add the absolute value of this density) from every measure of its energy density that contributes to determine the gravitational field experienced by what we conventionally assume to be positive-energy matter, that is, the gravitational field observed by positive-energy observers. I may add, however, that the required shift in the origin of the measures of energy, for matter with an energy sign opposite that of the observer, becomes significant only on the cosmological scale, because in the case of stars and planets it doesn't make much difference if we instead simply consider the true density of positive- or negative-action matter, given that the typical densities which are then involved are much larger than the mean cosmic energy density, which can thus be neglected.

The refinement discussed here is justified (theoretically) by the fact that, from the viewpoint of positive-energy observers, the description of negative-energy matter as voids in the positive-energy portion of the vacuum requires considering the contribution of negative-energy matter as being merely relative to the average density of this matter distribution (and therefore to actually be positive in the presence of underdensities in the average distribution of negative-energy matter), given that a uniform distribution of negative-energy matter has no effect on positive-energy matter, for reasons I have explained in section 2.6. The equations I initially proposed also allowed to express the fact that a similar requirement exists for the contributions of positive-energy matter to the total stress-energy tensor experienced by negative-energy observers. But, still, I did not find the set of equations I had proposed completely satisfactory. I thought that the right solution should bring about a simplification of the gravitational field equations, while, visibly, the equations I had derived were even less simple than the equations originally proposed by Einstein, despite the fact that, in their compact form, they were similar.

As I now understand, however, the equations I had proposed also fell short of meeting a certain mathematical requirement which I have come to appreciate as being essential to a consistent bi-metric theory of gravitation of the kind I sought to develop. This became clear when the paper [25] I cited above was published and new equations were proposed, apparently based in part on those I had developed, and which introduced a further refinement to bi-metric theories, by not assuming that there is a unique predefined relationship between the metric properties associated with the measurements of positive-energy observers and those associated with the measurements of negative-energy observers (even though for some reason the author of this

paper preferred not to consider that the matter contributing a negative measure to the total stress-energy tensor experienced by positive-energy matter actually constitutes negative-energy matter). As a consequence of this revised assumption, additional variables had to be considered that affected the contribution of negative-energy matter to the total stress-energy tensor experienced by positive-energy observers, or the contribution of what we usually consider to be positive-energy matter to the total stress-energy tensor experienced by negative-energy observers. The equations proposed were the following, in which the additional factors are written in their explicit form, using my notation⁹, and the quantities are now expressed in units where $c = 1$ and $G = 1/8\pi$:

$$\begin{aligned} R_{\underline{\mu\nu}}^+ - \frac{1}{2}g_{\underline{\mu\nu}}R^+ &= -(T_{\underline{\mu\nu}}^+ - \sqrt{\frac{g^{-+}}{g^{++}}}a_{\underline{\nu}}^{\underline{\nu}}a_{\underline{\mu}}^{\underline{\mu}}T_{\underline{\nu\mu}}^-) \\ R_{\underline{\nu\mu}}^- - \frac{1}{2}g_{\underline{\nu\mu}}R^- &= -(T_{\underline{\nu\mu}}^- - \sqrt{\frac{g^{+-}}{g^{--}}}a_{\underline{\mu}}^{\underline{\mu}}a_{\underline{\nu}}^{\underline{\nu}}T_{\underline{\mu\nu}}^+) \end{aligned} \quad (2.4)$$

The decisive additional factors are the determinants of what the author calls the pull-overs, which are the maps $g_{\underline{\nu\mu}}^-$ and $g_{\underline{\mu\nu}}^+$ (originally denoted $h_{\nu\mu}$ and $g_{\underline{\mu\nu}}$), which we may also write as \mathbf{g}^{-+} and \mathbf{g}^{+-} in tensor form. Those determinants are written here as $g^{-+} = \det(g_{\underline{\nu\mu}}^-)$ and $g^{+-} = \det(g_{\underline{\mu\nu}}^+)$, while $g^{++} = \det(g_{\underline{\nu\mu}}^+)$ is the determinant of the usual metric tensor related to properties of positive-energy matter as observed by positive-energy observers and

⁹From now on, I will use a notation that allows to better represent the relative nature of the physical properties associated with spacetime and the gravitational field. In this notation tensors which refer to positive or negative stress-energies, as determined from the viewpoint of positive-energy observers, will be given a plus or minus upper right index, respectively. Tensors which refer to measures of spacetime curvature or metric properties as observed by positive-energy observers will also be given an upper right plus index, while tensors which refer to the same kind of measures as observed by negative-energy observers will be given an upper right minus index. Also, when the distinct, ordinary or underlined Greek letter indexes used in Ref. [25] are not explicitly present to show the nature of the tensor considered, I will simply add another plus or minus index to the right of that which already characterizes this tensor to define it as an object associated with physical properties as they are experienced by positive- or negative-energy observers, respectively, and associated with their own specific metric. For all such tensors, therefore, the first plus or minus index refers to the matter or gravitational field that is observed while the second plus or minus index (to the right) refers to the matter that is observing. The underline which otherwise appears under some letter indexes can thus be considered as a shorthand for what should be additional plus or minus indexes over the letter indexes themselves.

$g^{--} = \det(g_{\mu\nu}^-)$ is the determinant of the metric tensor related to properties of negative-energy matter as observed by negative-energy observers (the map \mathbf{a} is simply used as a means to transform the metric \mathbf{g}^{++} into the \mathbf{g}^{-+} pull-over or the metric \mathbf{g}^{--} into the \mathbf{g}^{+-} pull-over). It is clear, therefore, that the pull-over \mathbf{g}^{-+} is the map which allows to describe the metric properties obeyed by negative-energy matter as they are observed by positive-energy observers, while the pull-over \mathbf{g}^{+-} is the map which allows to describe the metric properties obeyed by positive-energy matter as they are observed by negative-energy observers (which justifies my notation). To better illustrate the relationships involved we may rewrite those equations as:

$$\begin{aligned} R_{\mu\nu}^+ - \frac{1}{2}g_{\mu\nu}R^+ &= -(T_{\mu\nu}^+ - \gamma^{-+}\sqrt{\frac{g^{--}}{g^{++}}}a_{\nu}^{\underline{\nu}}a_{\mu}^{\underline{\mu}}T_{\underline{\nu}\underline{\mu}}^-) \\ R_{\nu\mu}^- - \frac{1}{2}g_{\nu\mu}R^- &= -(T_{\nu\mu}^- - \gamma^{+-}\sqrt{\frac{g^{++}}{g^{--}}}a_{\underline{\mu}}^{\underline{\mu}}a_{\underline{\nu}}^{\underline{\nu}}T_{\mu\nu}^+) \end{aligned} \quad (2.5)$$

where γ^{-+} is the absolute value of the determinant of the previously considered map of the metric properties of space experienced by negative-energy matter as negative-energy observers measure them, to the metric properties of space experienced by negative-energy matter as positive-energy observers measure them and vice versa for γ^{+-} . We can then rewrite those equations in compact tensor form by making use of those *metric conversion factors* as:

$$\begin{aligned} \mathbf{G}^+ &= -(\mathbf{T}^{++} - \gamma^{-+}\mathbf{T}^{-+}) \\ \mathbf{G}^- &= -(\mathbf{T}^{--} - \gamma^{+-}\mathbf{T}^{+-}) \end{aligned} \quad (2.6)$$

where \mathbf{G}^+ is the Einstein tensor $G_{\mu\nu}^+ = R_{\mu\nu}^+ - \frac{1}{2}g_{\mu\nu}R^+$ related to positive-energy observers, \mathbf{G}^- is the similar Einstein tensor related to negative-energy observers, \mathbf{T}^{++} is the stress-energy tensor of positive-energy matter as measured by positive-energy observers, $-\gamma^{-+}\mathbf{T}^{-+}$ is the stress-energy tensor of negative-energy matter as measured by positive-energy observers, \mathbf{T}^{--} is the stress-energy tensor of negative-energy matter as measured by negative-energy observers and finally $-\gamma^{+-}\mathbf{T}^{+-}$ is the stress-energy tensor of positive-energy matter as measured by negative-energy observers.

As is apparent, however, the proposed equations were still of the traditional kind, in the sense that they did not allow to take into account the fact that negative-energy matter is experienced as voids in the positive-energy portion of the vacuum (and vice versa for positive-energy matter from the

viewpoint of negative-energy observers). The complexity of those equations and their lack of symmetry under exchange of positive and negative energy states can be made more apparent by explicitly adding a term for the observed positive value of vacuum energy density:

$$\begin{aligned}\mathbf{G}^+ &= -(\mathbf{T}^{++} + \mathbf{T}_\Lambda^{++} - \gamma^{-+}\mathbf{T}^{-+}) \\ \mathbf{G}^- &= -(\mathbf{T}^{--} - \mathbf{T}_\Lambda^{+-} - \gamma^{+-}\mathbf{T}^{+-})\end{aligned}\tag{2.7}$$

In those equations $\mathbf{T}_\Lambda^{++} = -\Lambda\mathbf{g}^{++}$ would be the stress-energy tensor associated with the positive value of energy density of vacuum fluctuations $\rho_\Lambda^{++} = \Lambda$ measured by a positive-energy observer (with Λ as the positive cosmological constant experienced by such an observer), while $-\mathbf{T}_\Lambda^{+-} = \Lambda\mathbf{g}^{--}$ would be the stress-energy tensor associated with the negative value of energy density of vacuum fluctuations measured by what we would usually consider to be a negative-energy observer (which would consider energy of her own kind to be positive). The density of vacuum energy measured by a negative-energy observer must be the opposite of that measured by a positive-energy observer if the sign of energy is to remain an observer-dependent physical property (which justifies the presence of a minus sign in front of the \mathbf{T}_Λ^{+-} tensor that enters the gravitational field equations for negative-energy observers). But given that we are indeed dealing with vacuum energy, it would seem inappropriate to assign to this tensor the same metric conversion factor γ^{+-} as apply to measures of positive-energy matter density performed by negative-energy observers, even if the outcome of all positive and negative contributions to the energy of the vacuum is a positive energy, because, in principle, all such contributions exert a gravitational influence on both positive- and negative-energy observers on the cosmological scale. Anyhow, it is apparent that once all relevant contributions to the stress-energy tensors are considered, the symmetry of the original equations is lost, as their form becomes dependent on the actual sign of the average energy density of vacuum fluctuations. To me at least, it is obvious that those equations cannot be considered to embody a simplification of Einstein's theory that could be considered a substantial improvement over the original equations.

In order that such a formulation of bi-metric theory be allowed to at least meet the requirements I had already identified and which were not taken into account by the author of this later proposal, I would first suggest that we consider the limitations imposed on the interaction of positive- and negative-energy matter by the fact that the void of infinite extent in the

positive-energy portion of the vacuum that is equivalent to the presence of a homogeneous distribution of negative-energy matter has no gravitational effect on positive-energy matter (and vice versa when we consider the similar void in the negative-energy portion of the vacuum). In such a case, we would simply have to replace the usual stress-energy tensors associated with the measures of energy of negative- and positive-energy matter made by observers of opposite energies with the following *irregular stress-energy tensors*, which provide measures for the observed variations of energy density of negative- and positive-energy matter above and below their average cosmic densities:

$$\begin{aligned} -\gamma^{-+}\tilde{\mathbf{T}}^{-+} &= -\gamma^{-+}(\mathbf{T}^{-+} - \bar{\mathbf{T}}^{-+}) \\ -\gamma^{+-}\tilde{\mathbf{T}}^{+-} &= -\gamma^{+-}(\mathbf{T}^{+-} - \bar{\mathbf{T}}^{+-}) \end{aligned} \quad (2.8)$$

where $-\gamma^{-+}\mathbf{T}^{-+}$ and $-\gamma^{+-}\mathbf{T}^{+-}$ would be the usual measures of stress-energy of negative- and positive-energy matter, respectively (as experienced by observers of opposite energy signs), relative to the conventional zero level of energy and $-\gamma^{-+}\bar{\mathbf{T}}^{-+}$ and $-\gamma^{+-}\bar{\mathbf{T}}^{+-}$ are the measures of average stress-energy of negative- and positive-energy matter which would be determined by observers with an opposite energy sign if they could directly measure those parameters (it is precisely by measuring the irregular stress-energy tensor of negative-energy matter that a positive-energy observer can determine the average value of the stress-energy of that same matter).

In such a context, it appears that negative-energy matter would contribute negatively to the total measure of stress-energy experienced by a positive-energy observer only when the magnitude of its local energy density is larger than the magnitude of its average energy density. Otherwise negative-energy matter would actually contribute positively to the total measure of stress-energy experienced by a positive-energy observer, up to a maximum level fixed by the average density of negative-energy matter (the measure of average negative-energy matter density which would be determined by a positive energy observer, if such an observer could directly measure this density). The same remark would apply for the contribution of what is usually considered to be positive-energy matter to the total measure of stress-energy experienced by a negative-energy observer, which would be opposite the energy contribution of negative-energy matter only when the magnitude of the local density of positive-energy matter is larger than the magnitude of its average cosmic density.

It must be noted, however, that even though positive contributions to

the energy density measured by positive-energy observers may occur which would be attributable to the presence of underdensities in the negative-energy matter distribution, we must nevertheless apply the metric conversion factor γ^{-+} to such energy measures, because they still relate to measurements regarding the density of negative-energy matter, which are subject to the same mapping relationships as apply to other (truly negative) measures of energy related to negative-energy matter and made by a positive-energy observer. Of course, this is also true concerning below average measures of the energy density of what we would usually consider to be positive-energy matter made by negative-energy observers. Indeed, even when the second category of contributions to the energy density of matter is of the same sign as the energy of the matter experiencing the gravitational field, it is still undetermined to the same extent as negative contributions, because what is unknown (due to the impossibility to directly compare the measures of distances experienced by positive- and negative-energy observers) is the exact true density of negative-energy matter (in comparison with that of positive-energy matter) and this indefiniteness also affects the positive value of such contributions. Therefore, positive energy contributions arising from underdensities of negative-energy matter are contained in the same irregular stress-energy tensor as negative energy contributions.

A more appropriate set of gravitational field equations would, therefore, take into account the shifted origin of the measures of stress-energy related to positive- and negative-energy matter as they are experienced by observers of opposite energy signs:

$$\begin{aligned} \mathbf{G}^+ &= -(\mathbf{T}^{++} + \mathbf{T}_\Lambda^{++} - \gamma^{-+}\tilde{\mathbf{T}}^{-+}) \\ \mathbf{G}^- &= -(\mathbf{T}^{--} - \mathbf{T}_\Lambda^{+-} - \gamma^{+-}\tilde{\mathbf{T}}^{+-}) \end{aligned} \quad (2.9)$$

But clearly, for what regards simplicity, we appear to be no better off than with the previous set of equations. Something is still missing from those equations. At this point I suggest that we take a bold step forward and instead of trying to derive the gravitational field equations from a variational principle, as is usually done, we rather follow Einstein's way and simply guess what the final form of the equations should be that would generalize the set of equations (2.9) I have just proposed, which would otherwise constitute the most accurate description of the gravitational dynamics of positive- and negative-energy matter. As I have been able to understand, the crucial step in this process consists in reconsidering the meaning of the vacuum-

energy terms whose contributions I had long suspected were inappropriately attributed, in the context of bi-metric theories. Indeed, I always thought that the cosmological term should arise from an asymmetry between some positive and some negative contributions to the energy budget, while in the current set of equations it occurs only as an additional term, which must merely be attributed the appropriate energy sign depending on whether it is observed by a positive-energy observer or a negative-energy observer, which I do not find satisfactory.

It is only when I recognized the profound significance of my description of positive- and negative-energy matter as voids in their respective opposite-energy portions of the vacuum, that I was able to achieve the breakthrough that allowed me to guess what the appropriate generalized gravitational field equations are that allow the concept of negative-energy matter to be integrated into a general-relativistic framework in a way that actually simplifies Einstein's theory rather than further complicate things. What I realized, basically, is that if the results of the above described analysis is right, then all energy is vacuum energy, either present or missing. An additional insight was then necessary, which consists in recognizing that the magnitude of the natural positive and negative values of vacuum energy density relative to which are measured the missing energies which are equivalent to the presence of negative- and positive-energy matter (respectively) is actually provided by the Planck energy. What must be understood is that when we remove energy from the vacuum, we decrease its energy density from a maximum (positive or negative) value which is fluctuating quantum mechanically (upon measurement) in just the same measure as does the energy of matter itself. Therefore, if the presence of negative-energy matter is to be considered as equivalent to the presence of a void in the positive-energy portion of the vacuum, then locally we should observe a value of fluctuating vacuum energy density that would be decreased from its natural maximum value in just the same measure as that of the energy of the matter that is present.

But given that the level of fluctuation of vacuum energy involved would be as large as the void considered is small, it is possible to assume that there is an exact correspondence between the missing vacuum energy and the energy of the matter ordinarily expected to be present, which is known to be fluctuating (even if it is actually the measure of momentum that is involved) in proportion with the level of spatial confinement to which the matter is submitted. The natural level of energy involved would thus correspond to that which is known to be associated with the highest possible magnitude

of energy fluctuation, which is the Planck energy¹⁰. Therefore, any missing vacuum energy attributable to the presence of matter with an energy sign opposite that of the portion of vacuum in which it arises may be considered to actually be a local decrease over the maximum energy density determined by the Planck scale.

Let me thus introduce the generalized gravitational field equations which allow to fulfill all the requirements I have identified as being essential aspects of a classical theory of gravitation that solves the problem of negative energies. The formula, in all its beauty and simplicity, is the following:

$$\mathbf{G}^{\pm} = -\mathbf{V}^{\pm} \quad (2.10)$$

where \mathbf{G}^{\pm} is the Einstein tensor associated with the metric properties experienced by what we would usually consider to be positive- and negative-energy observers and \mathbf{V}^{\pm} is the *vacuum stress-energy tensor* associated with the measures of vacuum energy effected by those same positive- and negative-energy observers. The similarity with the compact form of Einstein's own equation is very clear, but it is also somewhat misleading, as the right-hand side of the equation proposed here is a much more general object than the stress-energy tensor of matter which appeared in the original theory. I will now define it with various levels of precision and generality. If we first consider the significance of the equation for a positive-energy observer, we would obtain the following equation:

$$\mathbf{G}^{+} = -(\gamma^{-+}\mathbf{V}^{++} - \mathbf{V}^{-+}) \quad (2.11)$$

in which \mathbf{G}^{+} is, again, the Einstein tensor associated with the gravitational field experienced by positive-energy observers, but now the vacuum stress-energy tensor is decomposed into its positive- and negative-energy portions

¹⁰The validity of this assumption could be the subject of controversy, but given that the most advanced and least speculative theoretical developments toward a theory of quantum gravitation indicate that this is an appropriate and unavoidable constraint, I will nevertheless consider it to be universally valid. However, even if the existence of such a limit to the energy associated with quantum fluctuations was to be found irrelevant, there is no *a priori* reason why the following results would have to be considered invalid. I believe that the situation we have here is once again similar to that which existed at the turn of the twentieth century concerning the hypothesis of the existence of atoms, which was often rejected on the basis of an absence of direct observational evidence, despite the fact that this assumption had actually become unavoidable theoretically.

$\gamma^{-+}\mathbf{V}^{++}$ and $-\mathbf{V}^{-+}$ as they can be measured by such positive-energy observers, based on the curvature of space they produce. This is the most basic form of the proposed generalized gravitational field equations for a positive-energy observer.

In accordance with what was explained above we would then obtain the next level of decomposition of the equations, in which the two opposite contributions to the energy of vacuum fluctuations determined by positive-energy observers are given their explicit form:

$$-\mathbf{G}^+ = (\gamma^{-+}\mathbf{V}_P^+ - \gamma^{-+}\mathbf{T}^{-+}) - (\mathbf{V}_P^- - \mathbf{T}^{++}) \quad (2.12)$$

where $\gamma^{-+}\mathbf{V}_P^+$ and $-\mathbf{V}_P^-$ are the *natural vacuum-stress-energy tensors* associated with the maximum, positive and negative contributions to the energy density of zero-point vacuum fluctuations set by the Planck scale and from which are subtracted the missing vacuum energies $\gamma^{-+}\mathbf{T}^{-+}$ and \mathbf{T}^{++} which are equivalent to the presence of negative- and positive-energy matter, respectively. What justifies the attribution of the previously introduced metric conversion factor γ^{-+} to the positive measure of vacuum stress-energy in equation (2.11) and therefore, also, to the maximum *positive* contribution to the energy of zero-point vacuum fluctuations in equation (2.12) is precisely the fact that this is the portion of vacuum energy relative to which the negative measure of matter energy $-\gamma^{-+}\mathbf{T}^{-+}$ is determined and which we can therefore expect to be directly experienced (other than through the gravitational interaction) only by this negative-energy matter, even though it does exert an observer-dependent gravitational force on positive-energy matter. Given that the previously introduced metric conversion factors are made necessary as a result of the absence of fixed relationships between the metric properties of space experienced by negative-energy matter and those experienced by positive-energy matter, it is natural to assume, in effect, that if the density of negative-energy matter itself cannot be directly observed by a positive-energy observer, then the positive measure of vacuum energy density relative to which this matter energy is defined cannot be directly determined either, because if this was not true, then by directly measuring the density of energy contained in the positive measure of vacuum energy, a positive-energy observer could determine the density of negative-energy matter which is experienced by negative-energy observers. What must be understood is that the fact that this portion of vacuum energy density is positive should not be assumed to invalidate the conclusion that it cannot

be directly experienced by positive-energy observers *other than through the gravitational interaction*¹¹.

The preceding equation can then be rewritten in the following form, when we take into account the previously introduced definition of the measure of stress-energy associated with negative-energy matter as it would actually be experienced by positive-energy observers, which are only affected by *variations* in the density of negative-energy matter:

$$-\mathbf{G}^+ = \mathbf{T}^{++} - \gamma^{-+} \tilde{\mathbf{T}}^{-+} + (\gamma^{-+} \mathbf{V}_P^+ - \mathbf{V}_P^-) \quad (2.13)$$

This allows one to isolate a term, in the generalized gravitational field equations, that can be associated with pure vacuum energy and that would be provided by the following stress-energy tensor:

$$\mathbf{T}_V^+ = \gamma^{-+} \mathbf{V}_P^+ - \mathbf{V}_P^- \quad (2.14)$$

where the positive index attributed to this *vacuum-energy term* (associated with the energy that is present in the vacuum independently from the contribution of ordinary matter) now merely denotes the purely conventional energy sign of the observer experiencing it, without referring to an actual energy sign of the vacuum fluctuations themselves, which could in principle be either positive or negative (without affecting the form of the equations) and which is determined solely by the metric conversion factor provided by the previously discussed map of the metric properties of space experienced by negative-energy observers onto those experienced by positive-energy observers. Indeed, given the invariant nature of the maximum positive and negative contributions to the density of vacuum energy associated with the Planck scale, for an observer having an energy sign opposite that of the contribution considered, the above equation means that a non-zero value of

¹¹Yet this is not the conclusion I originally drew when I wrote the first versions of this document, because at that time it seemed to me (for reasons that will be explained in section 4.2) that this hypothesis would be ruled-out from an observational (astronomical) perspective. But I have since realized that there are very good reasons to believe that this is not the case, after all, and that consistency requires that it is the portion of zero-point vacuum fluctuations that gives rise to a maximum positive contribution to the density of vacuum energy that cannot be directly observed by a positive-energy observer, even though it does interact with positive-energy matter gravitationally (while the portion of vacuum fluctuations that gives rise to a maximum negative contribution is the one that cannot be directly observed by a negative-energy observer, even though it does interact with negative-energy matter gravitationally).

vacuum energy density can only be measured by positive-energy observers when there exists a difference between the metric properties of space they experience and those which are experienced by negative-energy observers.

It is now possible to write the generalized gravitational field equations associated with positive-energy observers in their most explicit form as:

$$\mathbf{G}^+ = -(\mathbf{T}^{++} - \gamma^{-+}\tilde{\mathbf{T}}^{-+} + \mathbf{T}_V^+) \quad (2.15)$$

The formal equivalence of this equation with the first member of the equation (2.9), at which I had arrived on the basis of considerations of a physical nature, is quite clear. But while one may be tempted to deduce from this that the vacuum-energy term \mathbf{T}_V^+ is equivalent to the cosmological term \mathbf{T}_Λ^{++} which is present in the original version of the gravitational field equations, this would not be entirely appropriate, because contrarily to the cosmological term (associated with the cosmological constant Λ), which must by necessity provide a uniform and invariant contribution, the vacuum-energy term can vary in space and incidentally also with time, given that it is determined by the locally variable, metric conversion factor γ^{-+} . Thus, only the contribution associated with the average value of the vacuum-energy term at one particular time can be expected to be equivalent with the original cosmological term associated with the cosmological constant. In sections 4.2 and 4.3 I will explain how one must interpret the variable nature of the vacuum-energy term and why it is still appropriate to consider that, in general, the density of vacuum energy does not vary with position, in the absence of local inhomogeneities in the positive- and negative-energy matter distributions.

Anyhow, given that we know that on the cosmic scale, at least, the vacuum-energy term $\mathbf{T}_V^+ = \gamma^{-+}\mathbf{V}_P^+ - \mathbf{V}_P^-$ is very small, compared with the natural energy scale provided by the Planck energy, then it is possible to conclude that the correction provided by the γ^{-+} conversion factor is itself actually very small on such a scale. This observation, therefore, indicates that there is a near perfect level of symmetry between the metric properties of space experienced by positive-energy observers and those experienced by negative-energy observers at the present epoch, on a global scale. It may be added that if we are considering the above equation in a cosmological context, then the irregular stress-energy tensor $-\gamma^{-+}\tilde{\mathbf{T}}^{-+}$ would presumably reduce to zero on average (as the overdensities of negative-energy matter would cancel out the underdensities present in the same matter distribution), so that the relevant equations, for positive-energy observers, would be of the

following form:

$$\mathbf{G}^+ = -(\mathbf{T}^{++} + \mathbf{T}_V^+) \quad (2.16)$$

which is similar to their traditional form, except for the fact that the cosmological term \mathbf{T}_Λ^{++} is here replaced by the vacuum-energy term \mathbf{T}_V^+ that may vary with position. But given that local variations would presumably cancel out for vacuum energy as well, on a very large scale, and given the (relative) success of current cosmological models for predicting the relevant features of our universe's history, then this outcome would appear appropriate from an observational viewpoint.

We may then also write the following set of equations, which would provide the various levels of decomposition of the general equation (2.10) that apply from the viewpoint of negative-energy observers:

$$\begin{aligned} \mathbf{G}^- &= -(\gamma^{+-}\mathbf{V}^{--} - \mathbf{V}^{+-}) \\ -\mathbf{G}^- &= (\gamma^{+-}\mathbf{V}_P^- - \gamma^{+-}\mathbf{T}^{+-}) - (\mathbf{V}_P^+ - \mathbf{T}^{--}) \\ \mathbf{G}^- &= -(\mathbf{T}^{--} - \gamma^{+-}\tilde{\mathbf{T}}^{+-} + \mathbf{T}_V^-) \end{aligned} \quad (2.17)$$

where $\mathbf{T}_V^- = \gamma^{+-}\mathbf{V}_P^- - \mathbf{V}_P^+$ would provide the locally variable (positive or negative) value of vacuum energy density observed by such a negative-energy observer. The last equation, as well the other two, are now manifestly symmetric with the corresponding equations associated with positive-energy observers under a reversal of the sign of energy, as I have argued should be required. But the most remarkable feature of those equations (and the related equations for the gravitational field experienced by a positive-energy observer) is that they are actually obtained from a very simple expression (the first of the three equations) according to which the gravitational field experienced by an observer with a given energy sign is determined merely by the appropriate measures of (positive and negative) vacuum energy densities. This equation alone allows to embody the essence of the emerging framework. Indeed, it turns out that for the general equation (2.10) to give rise to the decomposition of energy contributions exhibited in the first and second of the observer-specific gravitational field equations (in which the metric conversion factors are present), all that is required is that the portion of zero-point vacuum fluctuations which directly interacts (other than through the gravitational interaction) with positive-energy matter produces a maximum value of energy density that is measured to be negative by a positive-energy observer, while the portion of zero-point fluctuations which

directly interacts (other than gravitationally) with what we would normally consider to be negative-energy matter produces a maximum value of energy density that is also measured to be negative by what we would usually consider to be a negative-energy observer (in the sense that the sign of this energy must be opposite that of the observer, which from a conventional Newtonian viewpoint would mean that it is positive).

The quantitative aspects of the proposed integration of negative energy states to classical gravitation theory having been properly introduced, it is now possible to look back and examine whether the equations obtained can actually provide the structure of an alternative model, which would conform to all of the principles enunciated in the preceding section. As I previously remarked, the basic structure of the proposed bi-metric theory was adopted precisely because it allows the kind of arbitrariness of the attribution of the sign of energy that is required for this physical property to be defined in a relational manner. But the ultimate confirmation that the proposed framework is compatible with the fundamental requirement expressed by principle 1 is the fact that, even in the presence of a non-vanishing value for the cosmological constant, the set of equations (2.17) describing the motion of negative-energy matter is now symmetric with the corresponding set of equations describing the motion of positive-energy matter under a reversal of the sign of energy. Furthermore, the requirement set by principle 2, that inertial mass be reversed along with gravitational mass, is also fulfilled by the proposed gravitational field equations, given that my analysis of the physical property of inertia has shown that imposing such a condition should give rise to gravitational attraction between masses of the same sign (whatever this sign is assumed to be) and to gravitational repulsion between masses of opposite signs and this is precisely what we obtain with the proposed equations, even if the sign of energy that replaces the sign of mass is here arbitrary and the gravitational field is a variable property, dependent on the nature of the matter submitted to it.

On the other hand, the validity of principle 3 and the absence of direct interaction between positive- and negative-energy matter particles may seem to be threatened by the fact that the stress-energy tensor associated with negative-energy matter contributes to determine the gravitational field experienced by positive-energy matter. But again, in the context of the more refined set of equations I have proposed, it is explicit that the negative contribution that enters the total measure of the stress-energy of matter that

determines a gravitational field and which we associate with the presence of negative-energy matter is actually a measure of the amount of stress-energy missing from the positive portion of vacuum energy. The effect on positive-energy matter, which must be taken into account in the presence of negative-energy matter, cannot therefore be attributable to an interaction with negative-energy matter (whose presence is not directly felt by a positive-energy observer), but must necessarily come from a *gravitational* interaction between positive-energy matter and the surrounding positive-energy vacuum. The equations, thus, naturally require that there be no direct interactions between particles with opposite energy signs.

The new equations are also the perfect embodiment of the requirements set by principles 4 and 5, because they allow the voids in the positive-energy portion of the vacuum to actually provide a negative contribution to the total stress-energy tensor of matter and in a general-relativistic context a negative contribution to the stress-energy of matter must be matched by a contribution to the gravitational field that is opposite that which is produced by positive stress-energy, so that if positive energy produces an attractive gravitational field from the viewpoint of positive-energy matter, negative energy must produce a repulsive gravitational field from the same viewpoint. The presence of voids in an otherwise uniform distribution of positive vacuum energy should therefore give rise to uncompensated gravitational forces opposite those attributable to the presence of an equivalent amount of positive-energy matter and by analogy the same should also be true for voids in a uniform positive-energy matter distribution.

We can now understand why it would be inappropriate to assume, as some authors do, that the energy of the gravitationally-repulsive matter whose behavior is described by conventional bi-metric theories is positive, even for an observer that measures a negative contribution from it to the total stress-energy of matter (so that the difficulties usually associated with the presence of negative energy matter could perhaps be avoided). Indeed, according to the above proposed equations, such matter would produce a gravitational field that would itself have an energy content (to the extent that a definite energy could actually be associated with the gravitational field) opposite that of the gravitational field which is produced by particles contributing positively to the total stress-energy of matter. But this means that if matter was assumed to always have positive energy, then, when energy is exchanged between the two types of matter, the variation of total gravitational energy (which would occur because opposite variations of *opposite* gravitational en-

ergies are involved) would not be compensated by a variation of the energy of matter (which would involve opposite variations of *positive* energies). Therefore, in the case of our two colliding bodies exerting a gravitational repulsion on one another, it would be impossible for the variation of the kinetic energy of the decelerating positive-energy body to be compensated by a variation of negative gravitational potential energy attributable to the changes occurring in the positive portion of vacuum energy as a result of the acceleration of the negative-energy body, despite the fact that this must be considered necessary if energy is to be conserved, as I previously explained.

Those problems can be avoided, however, when real negative energy states are allowed for matter, because, in a general-relativistic context, changes in the gravitational field can actually balance the changes occurring in the stress-energy of the two interacting matter components and given that indirect gravitational interactions are responsible for all energy exchanges between opposite-energy bodies, then no energy variations remain uncompensated. I think that this is a clear indication that the tentative solution to the problem of vacuum decay (the collapse of matter to ever more negative energy states) through the contradictory proposal of a gravitationally-repulsive matter that would have positive energy (from all viewpoints) is misguided and ineffective. Thus, if an observer is allowed to attribute a positive energy to matter of his own kind, regardless of which matter he is made of, it should be clear that, once this choice is made, the energy sign of the matter which from the viewpoint of this same observer provides a negative contribution to the stress-energy tensor of matter must be assumed negative. In any case, I must mention again that, from a cosmological viewpoint, the growth of negative-energy matter overdensities occurring in an initially homogeneous distribution of such matter will always be compensated by an opposite growth of underdensities in the surrounding environment. But given that from my viewpoint those two kinds of inhomogeneities provide opposite contributions to the total stress-energy tensor of matter experienced by a positive-energy observer, then it follows that there is an additional constraint regarding the conservation of energy contributed by negative-energy matter and this is a further confirmation of the viability of the proposed equations.

Returning to the criteria imposed by the principles enunciated in the preceding section, we can readily assess that the condition set by principle 6 (according to which only density variations over and below the average cosmic density of negative-energy matter have an effect on positive-energy matter) is also reflected in the equations proposed above. Indeed, the modified mea-

sure of negative stress-energy provided by the irregular stress-energy tensor $-\gamma^{-+}\tilde{\mathbf{T}}^{-+}$ which naturally enters the gravitational field equations associated with a positive-energy observer (given that the presence of negative-energy matter is here explicitly equivalent to an absence of positive energy from the vacuum) actually allows to fulfill the requirement set by principle 6, given that it provides a measure of stress-energy from which is subtracted the average stress-energy of negative-energy matter. This compliance of the proposed gravitational field equations may perhaps appear to be of secondary concern, given how negligible the average density of positive-energy matter (and even more so, that of negative-energy matter) really is in comparison with the density variations encountered under most circumstances when we are dealing with astronomical objects of interest, like stars or even galaxies. But, if it was not for the modified measure of negative stress-energy provided by the second term of equation (2.15), or the corresponding term from equation (2.17), serious problems would occur.

In section 2.5 (in which was elaborated the alternative concept of negative mass on which is based the mathematical framework developed here) I mentioned, in effect, that if a body with a given mass sign was to interact with all matter of both positive and negative mass that is present on the cosmological scale, then the classical phenomenon of inertia itself could not even exist (because the inertial forces resulting from acceleration relative to positive- and negative-mass matter would cancel out, either partially or completely). However, a Newtonian model is all about inertia, so that if inertial forces were made impossible by the presence of negative-energy matter, then reduction of the relativistic equations to a Newtonian gravitation theory with gravitationally-repulsive, negative mass densities would actually be impossible, even as an approximation. I believe that ignorance of the requirement to impose a suitable, modified measure of negative stress-energy for the generalized gravitational field equations is, in fact, the ultimate source of the difficulties which, according to certain authors, are encountered in trying to obtain an appropriate Newtonian limit from traditional bi-metric theories. This is in addition to the fact that, without the appropriate measure of negative stress-energy, complex hypotheses (of the kind which are often found in the literature) would have to be introduced concerning the variation in time of the ratio of the average cosmic densities of positive- and negative-energy matter in order to try to maintain the agreement of the proposed models with astronomical observations regarding the rate of expansion of ordinary,

positive-energy matter, which is already predicted with (relatively) good accuracy by traditional cosmological models, when no negative-energy matter is assumed to be present initially.

Finally, the fact that two maximum contributions of opposite signs to the energy density of the vacuum are now explicitly present in the most general form of each of the gravitational field equations means that both positive and negative contributions to the energy of the vacuum itself (ignoring voids) are allowed to contribute to the gravitational field experienced by positive- or negative-energy matter on the cosmological scale, as required by principle 7. From this alternative viewpoint, what allows one to appropriately ignore most of the effects that the vacuum would have on the gravitational field experienced by positive- or negative-energy matter is merely the fact that those opposite energy contributions nearly cancel each other out at the present epoch. I may also mention that the condition set by principle 8 (that the equivalence principle be valid, not merely locally, but really for one unique particle with a given energy sign) is implicitly contained in the structure of the proposed equations at the most basic level, because they describe gravitational fields which are dependent, not merely on the location, but also on the sign of energy of the particles submitted to them. On the other hand, principles 9 and 10, which identify requirements that have to do with the properties of matter particles (namely the absence of independent energy contributions for bound systems and the impossibility of a reversal of action on a continuous particle world-line), are not explicitly contained in the gravitational field equations proposed here, but if we assume the validity of those equations, then experimental facts make those constraints unavoidable.

Chapter 3

Time Reversal and Information

3.1 The problem of discrete symmetries

In this chapter I would like to explain how a more consistent and adequate formulation of the discrete P , T , and C symmetry operations, involving a revised concept of time reversal, can be obtained that integrates the insights gained while studying the problem of negative energy and that offers a better understanding of why and how such symmetries can, under certain circumstances, appear to be violated. Discrete symmetry operations are usually assumed to be relevant only in the context of quantum field theory, but in fact they can also be examined from a semi-classical standpoint. Their level of application is actually right at the interface between the classical world of gravitation theory and that of quantum theory and it should not come as a surprise, therefore, that some of the results which I have obtained will allow progress to be achieved concerning the problem of identifying the origin of the degrees of freedom associated with black-hole entropy, which arises merely in a semi-classical context. In order to do so it will be necessary to introduce an additional category of discrete symmetry operations that relates positive- and negative-action matter particles in a way that is similar in many respects with that by which the charge-conjugation symmetry operation relates ordinary matter and antimatter.

I had long ago realized that it would be necessary to revise our conception of space and time reversals, because the current formulation of those symmetry operations is based on unreasonable assumptions regarding the significance of time reversal and its relationship with the sign of energy and

that of non-gravitational charges. It is indeed presently believed that the charge-conjugation or C symmetry operation is not a discrete space or time symmetry operation, but simply an additional symmetry having to do with charge as an independent concept. But I came to suspect that the relationships which are known to exist between this charge reversal operation and the discrete P and T symmetry operations associated with space and time reversals are an indication that C should be conceived and explicitly defined as a particular instance of discrete spacetime symmetry operation. What constitutes the underlying basis of those considerations is the acknowledgment that the sign of certain physical quantities (including charge) are dependent on their direction of propagation in time. From that viewpoint it would seem, indeed, that both the T and the C symmetry operations should be assumed to involve some form of time reversal and this is reason enough to suspect that they may also both give rise to a reversal of charge.

The problem, however, does not really have to do with our current concept of charge reversal operation as such. What is truly inappropriate is the simple, kinematic representation of time reversal as involving a backward motion of all particles and their angular momenta, which I believe is too rudimentary to characterize a reversal of the fundamental time-direction degree of freedom. I also think that if T is to be assumed to actually reverse time, then it should leave momentum unchanged (despite common expectations) as this is a quantity that should rather be reversed independently, along with the direction of space intervals. In this context, if some reversal of momentum may still be of relevance to T it would clearly have to be due to the fact that it is actually equivalent to the effects we should expect to obtain from an appropriate reversal of time, when we insist on measuring physical quantities against the perceived, rather than the actual direction of the flow of time. In any case, it must be understood that what we observe from our classical historical perspective is not representative of the true evolution that takes place when we are dealing with the propagation of elementary particles. The subtleties of what is going on at the microscopic level are not directly apparent from the superficial viewpoint associated with a global representation of events ‘after the fact’ that provides a historical picture of the spacetime paths followed by elementary particles. Therefore, it is not appropriate to define a reversal of the fundamental (non-thermodynamic) time-direction degree of freedom based merely on narrative aspects of phenomena which are all directly discernible at this superficial level of description. Better formulations of the discrete spacetime symmetry operations are required which

would reflect the actual and sometimes unrecognized variations, or absence of variation of physical parameters associated with each of those reversals of the fundamental space- and time-direction degrees of freedom.

3.2 The constraint of relational definition

To begin this discussion, I must first of all mention that, once again, the most significant constraint which we need to consider and against which our understanding of the discrete symmetry operations must be developed is that of the necessary relational definition of physical quantities and their changes. Those quantities are here the directions of space and time intervals, the directions of momentum and angular momentum and the signs of energy and non-gravitational charges. The main point I want to emphasize is that there can be no meaning in considering a change of any one of those quantities (to its opposite value) that does not occur relatively to some remaining, unchanged parameter of the same kind. Breaking that rule is to be considered logically impossible, simply because if it was allowed it would mean that we can define an absolute (metaphysical) direction or polarity (in the general sense), which, in effect, would not be related to any reference point of a physical nature in our universe. What I'm suggesting is that the profound reason why a certain level of lopsidedness, such as the observed breaking of P symmetry by the weak interaction, can exist is that such asymmetries merely occur when one or two physical parameters are reversed *relative* to a fixed background of unchanged directional parameters of a similar kind. In other words, what makes these violations of discrete symmetry possible is simply the fact that application of a reversal operation to a single parameter leaves some other properties unchanged, which allows the asymmetry to occur as a real feature characterized by a measurable change relative to a distinct physical quantity. In the case of P symmetry, the reversal of space intervals involved occurs relative to the direction of time intervals, which remain unchanged by such an operation and therefore it should be expected that violations of P can be observed, given that the reversal of physical parameters associated with this operation can be measured against the unchanged properties.

But those asymmetries cannot imply the existence of an absolute lopsidedness or directionality at the most fundamental level, for the universe as a whole, because they can be compensated by an appropriate reversal

of the unchanged parameters relative to which the original transformation took place. This is what explains that despite the violation of P symmetry by the weak interaction, it remains impossible to provide an absolute definition of left and right, because indeed reversing the sign of charges allows to regain invariance. Thus, contrarily to what is sometimes assumed, the preferred handedness unveiled by the weak interaction is not more profound than that we observe in certain complex structures. As long as invariance under a more general discrete symmetry operation like CP is observed to hold, it is impossible to communicate the significance of right and left without knowing which of two C -related particles is to be considered as having positive electric charge. But if it is impossible to distinguish an absolute (non-relational) difference between positive and negative charges themselves, as I previously suggested, then only observers which are actually sharing the same universe and which are allowed to directly compare physical quantities, could differentiate between left and right.

This is a very general feature which I think would always be observed to apply, given that it is actually required by the condition of relational definition of physical quantities, which is relevant to any change of direction or polarity (such as a reversal of the sign of charges). The directions of space and time which are singled out by any process which appears to violate a discrete symmetry are significant only in relation to other aspects of reality which must be identifiable from within the universe in which those processes take place. If, in one particular instance, it was to be found that no combination of discrete symmetry operations allowed invariance to be regained, then it would mean that there exist physical properties which can refer to elements of reality not shared only by observers within our universe. In other words, if directional asymmetries not occurring merely in relation to unchanged quantities (not defined as mere relative properties) were allowed, it would, in effect, be impossible to describe the polarities so revealed by referring only to measurable properties of physical reality.

The problem which there would be if such violations of discrete symmetry were possible is that completeness and self-determination are the defining characteristics of the universe concept, in the sense that the universe is precisely that ensemble of physical elements which are all causally related to one another and to nothing else. Thus, if we were to find that the description of our universe can refer to absolute and immaterial notions of direction, not defined merely as relationships between elements of reality which must be part of that universe, then the only logically valid conclusion would have to

be that there exists a *causally related* reality outside what we consider to be the universe (this has nothing to do with the concept of the multiverse, whose elements are not to be assumed as causally related to one another) relative to which the otherwise metaphysical polarities could be properly defined. As a consequence, there is definitely no way our universe could be considered lopsided if it is actually the whole universe and I believe that the fact that it can be shown that the existence of such an irreducible asymmetry would imply that some physical quantities may not be conserved for the universe as a whole, is a confirmation of the validity of this conclusion. It must be understood, however, that the identified requisite does not mean that symmetry could never be preserved following a reversal of one single parameter, like space direction alone, which can be defined in a relational way, but simply that such invariance is not absolutely required to apply under all circumstances.

Given those considerations, we can be totally confident that there is no such thing as an absolute direction of space or time intervals, because, indeed, this would imply a violation of the principle of relativity (as understood in its most general form, which predates relativity theory) and the validity of this criterion is necessary for the consistency of any model concerning physical reality. Even without going into elaborate mathematical arguments, such as those entering the CPT theorem, it is therefore possible to appreciate that the only problem there could be in relation to the observation of an asymmetry under a properly defined discrete symmetry operation, would have to involve a violation of invariance under a combined operation that reverses all parameters and leaves absolutely none unchanged. I will later explain why an appropriately defined *PTC* transformation must be considered as one instance of such a symmetry operation that reverses all parameters and leaves nothing unchanged (by actually reversing all space- and time-related parameters twice) and which we are thus justified to categorize as inviolable. But I believe that the fact that it would be impossible to provide a mathematical framework for quantum field theory that would satisfy the requirements set by special relativity if the equations of the theory are not invariant under *PTC* (which constitute the substance of the argument behind the traditional CPT theorem), confirms that relativistic imperatives (all measures of space and time intervals are relative) are the true constraints which impose invariance under the most general, combined, discrete symmetry operation.

The fact that this simple, but most unavoidable requirement has never been considered as a means to restrict allowed violations of discrete sym-

metry illustrates the fact that our treatment of space and time reversals is incomplete and inadequate, due to multiple misconceptions which do not concern only the aspect discussed here. The often met remarks to the effect that there is no *a priori* reason why the universe could not be asymmetric in a fundamental way and that it is only the above mentioned mathematical requirements, arising from the CPT theorem, that motivate the conclusion that some overall symmetry must nevertheless be obeyed under all circumstances, are therefore inappropriate and misleading. But it should also not come as a surprise that the discrete symmetry operations, when performed independently from one another, may not produce invariance. What justified the unexpectedness of the violations of P and CP symmetries, when they were first observed, is actually the intuitive belief that absolute directionality should not be allowed, while, as I just explained, this is rather the argument that would apply to a more general symmetry operation like PTC whose required conservation, ironically, is usually not believed to be intuitively explainable. The truth is that, for an imbalance under reflection to exist, all that is required is that the world be unbalanced with respect to something. This conclusion is the outcome of the most unequivocal interpretation of the requirement of relational definition of physical quantities, which itself constitutes the one rule we can be most confident need to apply to the physical world we experience. In fact, the argument against the possibility of a violation of symmetry under a combined reversal of all space- and time-related parameters is probably the strongest kind of argument which can be proposed from a theoretical viewpoint.

Regarding time reversal, in particular, and the question of what it would mean to assume that the whole universe is running backward in time and whether there can be any objective meaning to such a reversal operation, I think that, given the preceding discussion, we would have to recognize that such a reversal could, in effect, be physically significant, if it is defined as a reversal that leaves other parameters, such as the direction of space intervals unchanged. But this means that such a time-reversal operation cannot consist in a mere reversal of the motions and rotations of objects taking place in a reverse chronological order. A reversal of time that would be relationally defined would have to be meaningful both globally and locally, as it would allow a distinction between a physical system with unchanged time direction and one with reversed time. This difference could be determined by directly comparing the physical properties of one of the systems with those of the other, if the two systems are part of the same universe. But a difference

could also be identified as occurring for the whole universe in relation to the unchanged direction of space intervals. In any case, the above discussed constraint would require that such a relative backward-in-time evolution be clearly identifiable from the physical properties of the particles involved, precisely because it is only under such conditions that the change of direction in time could be objectively determined by comparing it with that of the unchanged parameters. But given that those differences would then actually be determined in relation to the value of parameters which are themselves reversible, it follows that no absolutely characterized notion of asymmetry would be involved.

In the context where absolute lopsidedness is to be considered impossible, it follows that it is of primordial importance to identify all the physical properties which can be related to one another and which could be affected by transformations of the kind that involve a reversal of space and time directions at the fundamental level. Indeed, if we are to be able to determine whether there remain quantities not reversed when a certain discrete symmetry operation is performed, we certainly have to be able to determine which quantities are actually affected by the operation involved. It is my belief that some of the violations of discrete symmetries which are usually assumed to have been observationally confirmed are actually a consequence of the fact that the effect of the considered reversals on certain quantities are not taken into account, while invariance would actually be inferred if all quantities dependent on the parameters which are assumed to be reversed were appropriately transformed. I already mentioned the fact that there are indications to the effect that we may, in particular, expect the sign of charges to be dependent on the sign of time intervals experienced by the particles carrying them. Yet the traditional definition of the time-reversal operation T does not involve any reversal of charges (from whatever viewpoint) and thus we could observe violations of such a T symmetry that would occur simply because we do not appropriately reverse the sign of charges when we try to verify invariance under a reversal of time (from a certain viewpoint). We must therefore first take care of identifying all unaccounted dependencies which may confuse our assessment of symmetry violations, before we can truly appreciate under which conditions they are actually allowed to occur.

3.3 The concept of bidirectional time

Concerning the problem of discrete symmetries, another essential aspect must be recognized, in addition to that regarding the necessity of a relational definition of all such symmetry operations. Awareness of what it involves is of the highest importance for a proper resolution of all matters associated with time directionality and given that this is the central problem with which this report is concerned, it is crucial to grasp the significance and the implications of the notions involved. Basically, what must be understood is that a distinction is to be made between the traditional concept of time direction associated with changes occurring at a statistically significant level, where the notion of entropy is meaningful, and a concept of time directionality associated with the existence of a fundamental time-direction degree of freedom, independent from the constraints related to entropy variation. The traditional concept of time direction related to statistically significant changes and the growth of entropy gives rise to what I call the unidirectional- or thermodynamic-time viewpoint, while the alternative concept of time directionality, related to the existence of a fundamental time-direction degree of freedom independent from statistical constraints, gives rise to what I call the time-symmetric, or bidirectional-time viewpoint. In chapter 5 the refinement of the concept of time direction associated with the bidirectional viewpoint will be shown to allow the formulation of a principle of causality that is different from the traditional one and which no longer requires the existence of an absolute distinction between causes and effects.

But, associated with this alternative concept of time direction, is also a different notion of time reversal, not limited by the constraints imposed on our description of physical processes by the second law of thermodynamics. Indeed, the traditional notion of time reversal, associated with the thermodynamic-time viewpoint, merely consists in assuming a reversal of the motion of all particles involved in a process, so as to give rise to the same events as observed in the original process, but in the reverse order. However, those events would still be described from the same unique and immutable forward direction of time associated with entropy growth. This is a consequence of the fact that the unidirectional-time viewpoint involves considering that there can only be one direction in time at once for the propagation of all particles, indiscriminately, which actually amounts to ignore the existence of a fundamental time-direction degree of freedom. From that viewpoint, if time was reversed, all particles would have to propagate backward, not relative

to some fundamental time-direction parameter, but in comparison with the direction of motion which they were all *observed* to have originally. Thus, the time-reverse of a process would simply be the equivalent process for which the same observations are made, but in the reverse order. The bidirectional, or time-symmetric viewpoint, on the other hand, is at once less restrictive and more distinctive, in that it actually recognizes the existence of a fundamental time-direction degree of freedom, distinct from the observed direction of motion of particles apparent to an observer constrained by the law of entropy increase. This time-direction parameter must be allowed to vary from one particle to another, even between those of an otherwise identical nature which are involved in the same process at the same time.

Now, of course, I have already discussed the significance of the existence of a fundamental time-direction degree of freedom as being that property which allows to explain the distinction that exists between a particle and its antiparticle, despite the fact that from an observational viewpoint both objects appear to be ordinary particles traveling forward in time, but which merely happen to carry opposite non-gravitational charges. However, I previously made clear that, in fact, the sign of charge is *not* affected by a reversal of the direction of propagation in time which may relate a particle with its antiparticle and therefore, if it is nevertheless observed as being reversed, it can only mean that the direction of time relative to which we measure the charge is not the true direction in which the particle is propagating in time, because an observer measuring the same physical property while following the true direction of propagation in time of the particle would not observe any change¹. It is merely the fact that a backward-in-time observation is impossible that justifies assuming an apparent reversal of charges for a particle propagating toward the past. Indeed, measuring apparatuses always record changes as they occur in the future direction of time due to the fact that the processes involved in the amplification of the signal which gives rise to a measurement can only take place in this direction of time in a universe where a thermodynamic arrow of time governs the evolution of processes involving a large number of independently evolving particles. This constraint is there-

¹I will henceforth use the term ‘propagation’ in place of ‘motion’ to designate the true direction in which a particle is traversing space and time intervals, as occurs from a bidirectional-time viewpoint. This allows to explicitly refer to those aspects associated with the fundamental time-direction degree of freedom which are ignored from the viewpoint of unidirectional time, relative to which all changes refer to a particle’s observed (semi-classical) trajectory.

fore what justifies the use of a unidirectional viewpoint, relative to which all physical properties are given as they would appear relative to the conventional future direction of time, even when the true direction of time in which the processes involved occur is the past direction. Non-gravitational charges, therefore, actually remain unchanged from the bidirectional viewpoint when the fundamental time-direction degree of freedom is reversed, but this is the very reason why they appear to be reversed from the unidirectional-time viewpoint.

A rule, thus, emerges which is that, for any particle propagating in the past direction of time, a time-direction-dependent physical property of that particle which would be positive when considered from the bidirectional-time viewpoint (relative to the true direction of propagation of that quantity in time), would appear as negative from the unidirectional-time viewpoint. But this reversal of observed quantities from their true value is not restricted to charge or energy, which I had already identified as properties dependent on the direction of propagation in time, but would actually have to apply to the direction of space intervals associated with the motion of particles (which are always given in relation to time intervals) and thus, also, to momentum (even if the time intervals entering the traditional definition of momentum were assumed positive definite as a consequence of adopting a unidirectional-time viewpoint). Thus, if momentum was assumed to be left unchanged by a properly defined reversal of time, it would nevertheless appear to be reversed in comparison with its actual value, from the unidirectional-time viewpoint. But given that the direction of momentum is not fixed for a given type of particle, propagating in a given direction of time (it also changes when the direction of propagation of the particle in *space* is reversed), it cannot be taken as a clear indicator of the direction of propagation in time of a particle. That, however, is not the case with charge, which from the bidirectional-time viewpoint remains unchanged, even as a particle reverses its direction of propagation in time (while also reversing its energy sign), and this is why it is possible, from the unidirectional-time viewpoint, to identify the true (even if merely conventionally-defined) direction of propagation in time of a particle, based on the observed value of its non-gravitational charges (in relation to those of an otherwise identical particle)².

What is important to understand is this interdependence of space and

²In fact, even if this relationship between time direction and observable charge was valid only for ordinary particles and antiparticles, while it would be possible to conceive of a distinct operation of charge reversal that would reverse charge independently and

time intervals, even as they would be separately and independently transformed by their respective, discrete symmetry operations. Thus, when we reverse the direction of the motion of a particle in space, we reverse the sign of the space intervals associated with this motion, not merely relative to the position axes, but also relative to time intervals (same time interval, opposite space interval). The sign of space intervals associated with the propagation of a particle submitted to a reversal of space directions would be reversed not merely from what it previously was (or relative to the space intervals associated with the motion of a particle not subject to the reversal), but also relative to the direction of time intervals in which the particle is still propagating. A particle which was propagating to the right, relative to the future direction of time, will now be propagating to the left, relative to the same future direction of time, which was not affected by the reversal of space directions (this is illustrated in Figure 3.1 where I consider the effects of the various discrete symmetry operations as they will be defined below). In other words, the particle is not just propagating left, it is propagating left, forward in time, because indeed we are always concerned with the properties of processes involving particles propagating in space and time and not just with the properties of space or time themselves. What matters, therefore, is not just the direction of space intervals associated with some arbitrarily-fixed spatial coordinate system, but the direction of space intervals for a particle propagating in a given direction of time, as asserted from a fundamental bidirectional viewpoint. Similarly, when time is assumed to be reversed, it must be considered that the time intervals are reversed relative to the unchanged direction of space intervals in which a particle submitted to the reversal is propagating, so that the same positive space intervals are now traveled in the opposite direction of time. This does not mean that a reversal of both space and time cannot have clear meaning, however, because, as I will explain later, even in such a case there would still remain unchanged physical properties relative to which the transformation could be characterized.

This relationship between space and time intervals is what gives a true physical meaning to the notion of time reversal, when it is to be considered as a symmetry operation clearly distinct from space reversal and which should,

not merely as result of a reversal of the direction of propagation in time of particles, this conclusion would still be valid, because, as I will explain in section 4.3, particles carrying such a reversed bidirectional charge would remain clearly distinguishable from ordinary particles and antiparticles, regardless of the direction of time in which they are propagating.

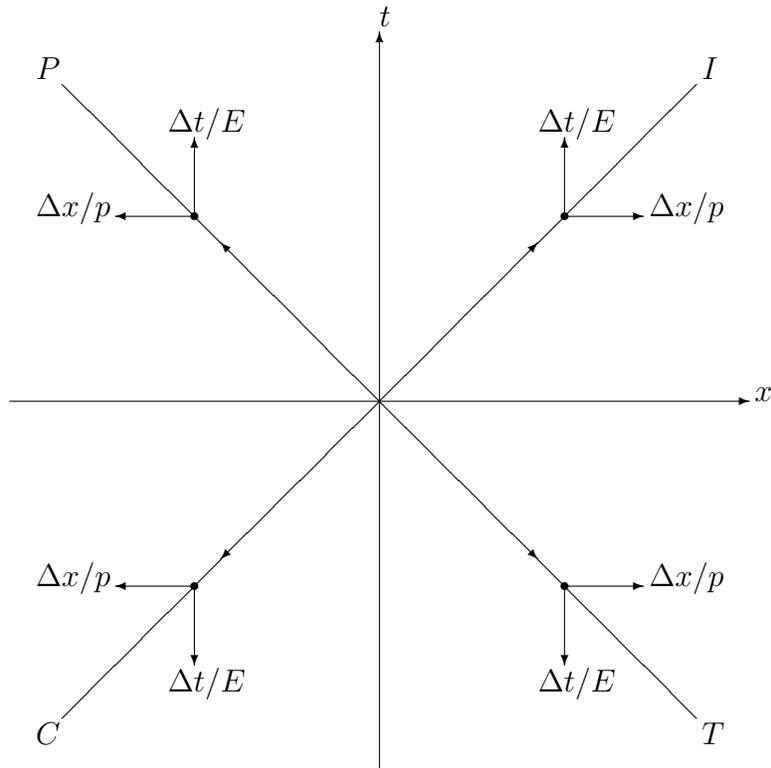


Figure 3.1: Variation of physical parameters under the proposed alternative definition of P , T , and C , as described from the bidirectional-time viewpoint. In this figure and the other related figures, I represents the original state and the diagonal lines correspond to particle trajectories. The space and time intervals Δx and Δt are indicated by vectors whose lengths correspond to the magnitude of the intervals and whose directions indicate the sign of the intervals relative to the space and time coordinates. The direction of the vectors associated with the momentum p and energy E of particles corresponds with the sign of momentum and energy relative to the direction of the space and time coordinates.

therefore, leave momentum unaffected (from the bidirectional viewpoint, at least). In fact, it is what allows the very notion of a fundamental degree of freedom associated with direction in time to have a definite meaning, because it allows to distinguish (as a theoretical possibility) the process by which a particle is going through a given spacetime trajectory forward in time from the similar process by which an identical particle would be going through the exact same spacetime trajectory, only now backward in time. Such a distinction is crucial, given that if we were to ignore it, then from a unidirectional viewpoint in time there would be no meaning to assume that it may be possible for a trajectory to be traversed backward in time, given that from such a viewpoint we always observe particles as if they were necessarily going forward in time. But given that charge can be assumed to be left unchanged by a reversal of time (from the bidirectional viewpoint), we are actually allowed to differentiate between those two situations, from an observational viewpoint, even in the context where all particle trajectories are necessarily followed as if they were occurring in the ‘normal’ chronological order (forward in time) associated with the growth of entropy, regardless of the true direction of propagation in time of the particles. It is, therefore, the relation between space intervals and time intervals that allows to distinguish backward-in-time propagation from forward-in-time propagation and the fact that the observed value of the sign of charge is dependent on that distinction simply confirms that it is appropriate to consider the existence of such a directionality parameter for the time dimension at the fundamental, elementary-particle level.

It must be clear, however, that the coordinate systems for space and time still have a physical significance, because you may reverse the direction of the space intervals traveled by particles in the forward direction of time, as well as the associated momenta, while keeping the positions of the particles in space unchanged (not reversed as they would under a conventional space reversal operation). Indeed, as a comparison of figures 3.1 and 3.2 allows to reveal, it is only from the bidirectional-time viewpoint that the sign of space and time *intervals*, corresponding to the directions of propagation of particles, always change in association with the sign of *positions* on the space and time coordinate axes, while from the unidirectional-time viewpoint that need not be the case. Under such conditions, quantities like angular momentum, which depend on both the position in space and the direction of space intervals, may not always be left invariant, as they would when a complete space reversal operation is performed. This would occur, in effect, for processes

submitted to a reversal of time, when they are described from the unidirectional viewpoint in which time is maintained positive, even for backward-in-time-propagating particles, and all time-direction-dependent quantities like the direction of space intervals and the momentum of a particle consequently appear to be reversed, while the positions are left unchanged (which implies that spin would appear to be reversed). In this context, it seems that space intervals, as properties defined in relation to the direction of propagation in time, can actually be reversed in two different ways. They may be reversed because space directions are reversed (which also reverses positions), or they may be reversed because the direction in which they are assumed to be traversed in time is reversed (which leaves positions unchanged). This distinction is what allows the traditional concept of time reversal, as affecting the directions of momentum and angular momentum, to still be relevant, even in the context of the existence of a fundamental time-direction degree of freedom, when those directions should in fact be left invariant (from a bidirectional viewpoint) by a properly defined time-reversal operation.

Another point must be emphasized regarding the kind of time-reversal operation which can be developed in the above described context. Indeed, if we no longer consider appropriate the picture of time reversal as consisting in a simple reversal of the observed motion of each and every particle, then it must also be recognized that a properly defined time-reversal operation could never give rise to a reversal of the thermodynamic arrow of time for the physical systems involved. In fact, I think that we should already suspect that there is something wrong with the often-met suggestion that a reversal of the motion of every particle in a region of space would give rise to entropy decreasing evolution (in the absence of any external perturbation). For such a proposal to be valid it would have to be shown that the origin of the observed time asymmetry of thermodynamic processes in our universe is to be found in a very precise adjustment of the motion of every single particle in the universe at the present time, which would occur in just such a way as to allow a state of minimum entropy to be reached as time unfolds in the past, right back to the Big Bang state.

However, given the inherently random nature of quantum processes and the extreme sensitivity to initial conditions (here the ‘final’ conditions giving rise to a given past evolution) which are known to exist, even in a classical context, this hypothesis appears highly implausible (I will address this question more thoroughly in section 4.6). But if, in addition, we admit the existence of a fundamental time-direction degree of freedom, distinct from

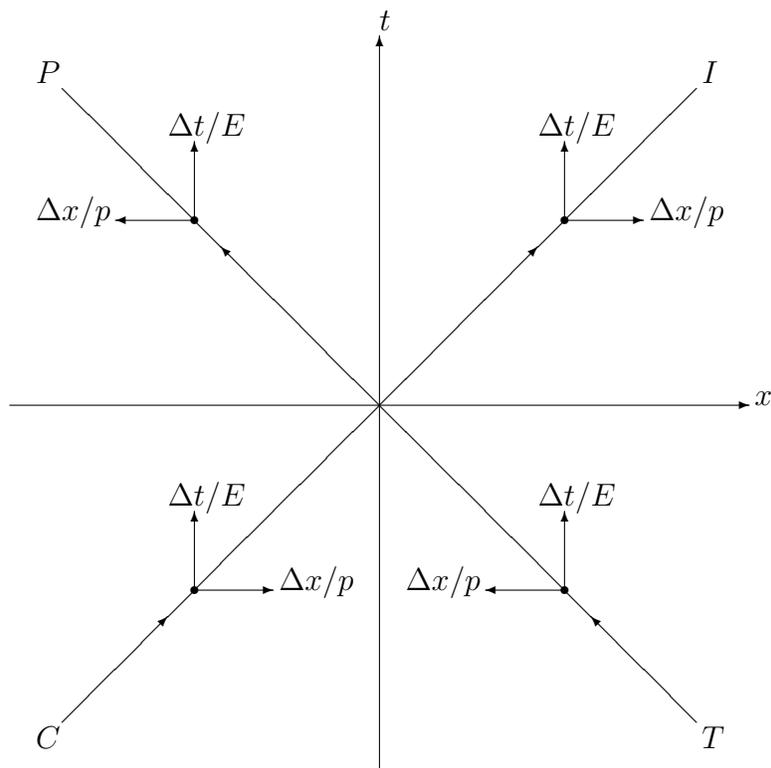


Figure 3.2: Variation of physical parameters under the proposed alternative definition of P , T , and C , as apparent from the unidirectional-time viewpoint. We can see that, from this viewpoint, the only difference between the original process and the T -reversed process is that the space intervals are traversed in the opposite direction, just as would be expected according to the traditional definition of backward-in-time motion. The case of the C -reversed process is also quite in line with traditional expectations, given that such a process should not be different from the original process except for a reversal of the sign of charges (which is not illustrated here) which would in fact also occur for the T -reversed process, despite traditional expectations.

the observed motion of particles, then we clearly have to reject the possibility that a reversal of time may produce anti-thermodynamic behavior, because time-reversed propagation is in fact already taking place in processes for which there is no apparent change to the direction of the thermodynamic arrow of time. This means that the direction of propagation in time of particles (the sign of time intervals associated with a bidirectional viewpoint) is not necessarily that relative to which entropy increases, despite the fact that it may appear unnatural that evolution could proceed in a direction of time other than that in which we do observe time to be ‘flowing’ (as a thermodynamic necessity). The thermodynamic arrow of time and the notion of time directionality occurring from a bidirectional viewpoint are two completely independent concepts.

3.4 Alternative definition of C , P , and T

One last remark is necessary before I can provide a full description of exactly how the fundamental physical properties of matter should be considered to vary under an alternative set of discrete symmetry operations formulated so as to allow the above discussed requirements to be satisfied. I previously hinted at the fact that the direction of momentum should be considered as independent from the direction of time, at least from the most consistent viewpoint, which is provided by a bidirectional perspective on time. I believe, in effect, that momentum, as the attribute conjugate to physical space, should only be considered to reverse along with space and not along with time, just as energy, being the physical attribute conjugate to time, should necessarily reverse when time reverses and only then. There is, however, an additional motivation for requiring this kind of joint variation of all space-related attributes, or time-related attributes (independently), besides the fact that consistency may require that it be imposed when what we seek to assert is precisely the dependence of various parameters under reversal operations which are defined after the quantities they are assumed to reverse. This, perhaps more unavoidable, justification for the joint variation of conjugate attributes is to be found in the requirement that the considered symmetry operations should not change the sign of action of the physical systems on which they operate.

It is my understanding of the true physical significance of a reversal of the sign of action that allows me to recognize the necessity to define the discrete

symmetry operations in such a way that momentum would necessarily reverse as a consequence of a reversal of space coordinates, while energy would necessarily reverse as a consequence of a reversal of the time coordinate. Indeed, in the context where a reversal of space coordinates would necessarily give rise to a reversal of space intervals, while a reversal of the time coordinate would necessarily give rise to a reversal of time intervals, if the sign of action itself is to remain invariant, then it means that a reversal of space must also involve a reversal of momentum and a reversal of time must also involve a reversal of energy. In fact, we always implicitly assume that the P , T , and C reversal operations do not relate physical processes in which the particles involved would have opposite action signs or energies (as measured from the forward direction of time). But the implications this should have for the dependence (under conventional discrete symmetry operations) of the signs of momentum and energy on those of space and time intervals is not always recognized. I believe that this lack of clarity is responsible for a good part of the misunderstanding regarding what parameters should really be affected by any symmetry operation involving a reversal of time. In Tables 3.1, 3.2, and 3.3 I will therefore provide an explicit account of the dependence of the signs of momentum and energy, along with those of space and time intervals, under all relevant discrete symmetry operations. It will be apparent from this account that clear distinctions exist between the traditional and the redefined time-reversal and charge-conjugation symmetry operations. Yet, given that the original definitions actually need to be replaced and cannot even be considered meaningful anymore, I think that it will not be necessary to relabel those operations and associate them with new symbols or letters, so that I will continue to use the T and C notation when referring to those redefined discrete symmetry operations.

In the following tables and in the corresponding diagrams (Figure 3.1 corresponds to Table 3.2 and the bidirectional viewpoint, while Figure 3.2 corresponds to Table 3.3 and the unidirectional viewpoint) the position along the space and time axes are denoted x and t (I'm assuming a one-dimensional space for simplicity), while the space and time intervals corresponding to the motion, or the propagation of the particles involved in the processes which are transformed by the symmetry operations are denoted Δx and Δt respectively. The energy of the particles involved in the same processes is denoted E and can actually vary in sign, while the momentum of those particles along the x axis is simply denoted p . The sign of non-gravitational charges (which allows to distinguish between the state of a particle and that

Tradit.	t	Δt	E	x	Δx	p	q	s	h
I	t	Δt	E	x	Δx	p	q	s	h
P	t	Δt	E	$-x$	$-\Delta x$	$-p$	q	s	$-h$
T	$-t$	Δt	E	x	$-\Delta x$	$-p$	q	$-s$	h
C	t	Δt	E	x	Δx	p	$-q$	s	h

Table 3.1: Variation of the physical parameters associated with a process transformed by the discrete P , T , and C symmetry operations, as they are traditionally defined. The variations of the Δt and Δx parameters indicated here are only implicitly assumed from a conventional viewpoint. The absence of reversal of Δt when time is assumed to be reversed can be noted. The variation of the direction of angular momentum s , as well as that of the handedness h , can be derived from those of the other fundamental parameters, but they are nevertheless indicated here and in the other tables, because in certain cases they differ from what is traditionally expected. The identity operation I which corresponds to an absence of reversal is shown for reference purpose.

of its antimatter counterpart), even though it should be understood not to be reversed by any of the conventional discrete symmetry operations (including C) from the bidirectional-time viewpoint (which provides the most accurate description of the transformations involved), is nevertheless included in the tables and denoted q , as it may actually appear to be reversed from the unidirectional viewpoint by some of those symmetry operations. The sign of angular momentum, related to the motion of the particles involved in the processes transformed by the P , T , and C operations, as well as the spin direction of elementary particles, which again should be understood not to be affected by those operations from a bidirectional-time viewpoint, are together denoted by the letter s , while the associated parameter of handedness (the direction of spin along the axis associated with the momentum of a particle) is here denoted h and should be expected to vary, even from a bidirectional-time viewpoint.

From a semi-classical viewpoint, the displayed tables, giving the variations of the time-related and space-related physical parameters under the traditional or redefined discrete symmetry operations, along with the assumptions which are made concerning the variation of the sign of charge, provide

Bidir.	t	Δt	E	x	Δx	p	q	s	h
I	t	Δt	E	x	Δx	p	q	s	h
P	t	Δt	E	$-x$	$-\Delta x$	$-p$	q	s	$-h$
T	$-t$	$-\Delta t$	$-E$	x	Δx	p	q	s	h
C	$-t$	$-\Delta t$	$-E$	$-x$	$-\Delta x$	$-p$	q	s	$-h$

Table 3.2: Variation of physical parameters under the redefined discrete P , T , and C symmetry operations, as described from the bidirectional-time viewpoint. The necessary reversal of Δt with E , as well as that of Δx with p , can be noted, as also the necessary reversal of t with Δt and that of x with Δx . This is the variation of physical parameters which would be produced by the most appropriately defined discrete symmetry operations that can be formulated in a semi-classical context. Here, all reversals of physical quantities are seen to occur twice or not at all, as required for explicit invariance under a joint PTC operation.

Unidir.	t	Δt	E	x	Δx	p	q	s	h
I	t	Δt	E	x	Δx	p	q	s	h
P	t	Δt	E	$-x$	$-\Delta x$	$-p$	q	s	$-h$
T	$-t$	Δt	E	x	$-\Delta x$	$-p$	$-q$	$-s$	h
C	$-t$	Δt	E	$-x$	Δx	p	$-q$	$-s$	$-h$

Table 3.3: Variation of physical parameters under the redefined discrete P , T , and C symmetry operations, as described from the unidirectional-time viewpoint. Again, all quantities are reversed either twice or never by a combination of all operations, which guarantees explicit invariance under PTC . The equivalent reversal of charge q by both T and C , as well as the apparent absence of any variation of Δt and E , and the absence of joint variation of x and Δx when t is reversed, can be noted.

the most precise definitions that can be achieved of the operations involved. Using those definitions, one can rebuild the quantum operators which are needed to transform the state vectors or the propagators corresponding to specific quantum states or processes. It must be clear that quantum field theory itself does not dictate how the discrete symmetry operations should be defined and it is merely the assumptions used while formulating the related operators (to achieve transformations that match our expectations regarding which parameters should be affected by a given operation) that provide the necessary constraints on which depend their precise mathematical formulation. What I bring to the table, therefore, is an improved knowledge of the constraints that must apply to those transformations, based on a re-examination of the meaning of space and time reversals, as they would occur in a semi-classical context. It is important to recognize, indeed, that despite the apparent freedom, the discrete symmetry operations cannot be arbitrarily defined, but must be the outcome of the most unavoidable consistency requirements (formulated in an empirically motivated context), which I believe are those I have identified in the above discussion. The fact that greater simplicity has been achieved while redefining those symmetry operations is only a further confirmation of the appropriateness of the alternative viewpoint that emerged from the preceding analysis. Indeed, the pattern of variations of physical parameters which is illustrated in Figure 3.1 is strikingly simple in comparison with that we would have according to the traditional definition of the discrete symmetry operations and this simplification was actually one of the objectives I sought to achieve while redefining them. Let me then describe what the elegance of this proposal really embodies.

Looking at the tables in which the outcomes of the various discrete symmetry operations are displayed, one thing we may first remark is that the parity operation P remains as it was originally defined, even in the context of the proposed alternative formulation of those transformations and this regardless of whether we use the bidirectional- or the unidirectional-time viewpoint. Of course, the reversal of space intervals associated with the propagation of particles (which from my viewpoint must occur as a result of the reversal of space coordinates) is now explicitly stated, but, otherwise, the traditional definition of space reversal remains unchanged. There is one good reason for that, which is that the revision I'm operating regards the concept of time direction, essentially, and the P operation is unique for being the only one that does not involve any time reversal, regardless of the approach favored. This is what explains that this operation was properly

defined already, in the form it originally was, despite the failure of the traditional viewpoint in general. What P expresses, indeed, is a reversal of space coordinates that produces a reversal of positions, space intervals, and naturally, also, momentum (as a requirement of action-sign invariance), while it leaves unchanged (now as a matter of definition) the position in time, the time intervals and the sign of energy. No reversal of charge is to be observed in this case (particles are not replaced by antiparticles), from any perspective, because there is no time reversal involved from a bidirectional viewpoint and thus no change to be associated with the adoption of a unidirectional-time viewpoint. There is no reversal of angular momentum either (because both momentum and position are together reversed), which is appropriate given that if angular momentum or spin were reversed, a forbidden reversal of action would occur from the bidirectional viewpoint (because spin has the dimension of an action) that would not be associated merely with the shift to a unidirectional-time viewpoint. But again, this is in perfect agreement with traditional expectations regarding the effects of P . Handedness is to be assumed reversed by such a reversal of space, however, because momentum is reversed while spin is left invariant from all viewpoints.

It should be noted that the explicit mention of a reversal of space intervals Δx under a symmetry operation like P does not mean that a reversal of space intervals must be assumed to occur in addition to that produced by the reversal of space coordinates. In other words, if the space intervals are indeed reversed, it is merely as a consequence of the reversal of space coordinates, as otherwise there would be no real change in the direction of space intervals, that is, no change relative to the new coordinates. We may, in fact, consider it more appropriate to assume that it is the intervals themselves which are reversed along with the position of particles while the coordinates remain unchanged, which would still be equivalent to reversing the coordinates themselves. If I choose to explicitly mention a reversal of space intervals, along with the assumed reversal of positions, it is because there may be situations where the intervals would be reversed independently from the positions on the coordinate axes and we must be able to distinguish between the two situations. What the explicit statement of a reversal of Δx should be understood to imply, therefore, is that there must occur a reversal of the sign of space intervals traversed by the particles involved in the reversed processes, in comparison with the sign of space intervals experienced by particles involved in processes which would not be submitted to the reversal. We must, therefore, assume that those reversed intervals are the space

intervals which are traversed during unchanged time intervals and which we may ordinarily associate with the directions of the momenta of the particles involved. Indeed, the reversal of space intervals associated with the motion of particles is usually assumed to be implied by the reversal of momentum itself, but given that I will later suggest that momentum can be reversed without space intervals being equally reversed (when action is to be considered reversed), then it becomes necessary to explicitly define the variation of space intervals under P and to recognize that momentum direction is an independent quantity, whose specification is not sufficient to determine the sign of space intervals spanned during a given time interval (except if the action sign is, in effect, required to be invariant).

It must be recognized, therefore, that the reversal of Δx is not *merely* a reflection of the reversal of space coordinates, but that it also allows to denote the physical changes that occur when a particle reverses its direction of propagation in space, while retaining its direction of propagation in time and those changes would be significant even if the position in space was to itself remain unchanged. Likewise, what the specific statement about the reversal of momentum p under space reversal P is intended to mean is that the direction of momentum is now the opposite of what it was, not merely relative to the new coordinates, but also relative to the directions of the momenta of particles which would not be subject to the symmetry operation. I may add that the same remarks would apply to time intervals Δt and the sign of energy, because if the reversal of those physical parameters under the T and C operations (from a bidirectional viewpoint) can be understood to occur as a consequence of the reversal of the time coordinate, it is clear that it also arises in relation to the time intervals experienced by particles which would be left unaffected by the reversal.

3.5 The time-reversal operation

Despite a concordance of the rules from which are derived the variation of physical parameters under any one of the redefined discrete symmetry operations, there are important differences between the case of time reversal T or charge conjugation C and that of space reversal P and this is reflected in the fact that those two symmetry operations would produce results which are unexpected from a traditional viewpoint. In the case of T , it must be required, in effect, that the physical time intervals Δt associated with the

propagation of elementary particles and the energy E be together reversed when the time coordinate is reversed (if action is to remain positive when it already is), while it is traditionally assumed (even if only implicitly) that both energy signs and *bidirectional* time intervals are in fact unchanged by T despite the reversal of the time coordinate. Also, it must now be assumed that there is no *a priori* reversal of the space intervals Δx and momentum p when time is reversed (which is allowed when those parameters are recognized to be independent from the time-related parameters Δt and E). This is required, despite the fact that, traditionally, momentum is assumed to be dependent on time intervals (I will explain below how this apparent contradiction is to be resolved). In fact, the traditional assumption that p would be reversed by T , while the position x on the space axis would remain unchanged, would be problematic if, in this context, we did not again implicitly assume an independent reversal of physical space intervals Δx , by presuming an invariance of the sign of action.

What must be recognized, therefore, is that from a consistent bidirectional viewpoint, when the time coordinate is reversed, it must be assumed that the time intervals of propagating particles (associated with the fundamental time-direction degree of freedom) are reversed along with their energies (as defined relative to the true direction of propagation in time), while momentum and space intervals are left unchanged, just like a reversal of space coordinates is assumed to imply a reversal of the space intervals and momenta, but no change to energy sign and no reversal of time intervals. This independence of space- and time-related physical parameters (from one another) is a requirement of the constraint of relational definition of those quantities, which imposes that something remains unchanged when T or P is applied, and those invariant properties are in fact the spatial directions themselves (when the direction of time is reversed) or the direction of time itself (when space directions are reversed).

Now, if we appropriately assume that the spatial positions, the space intervals, and the momenta remain unchanged under a properly defined time-reversal operation, it follows that the spin and the handedness must also remain invariant. Those relationships may appear unnatural (spin is usually considered to be reversed under a reversal of time), but from a bidirectional-time viewpoint they are perfectly acceptable and in the context where we want to define time reversal as really affecting time-related parameters in a specific way, they actually constitute unavoidable requirements. What's more, the discussed invariance is observed from the bidirectional-time view-

point, according to which the values of physical properties are such as they would appear to an observer following the direction of propagation in time of the particles involved in the processes submitted to this reversal. But from a unidirectional-time viewpoint (of the kind that is required from a practical perspective), the only quantities which would appear to be left unchanged when time is reversed would actually be the time intervals Δt and the energies E , because they would be submitted to twice the same reversal, once as time-related quantities, and once as a consequence of the additional reversal occurring when we are forcing a forward-in-time perspective. This is what justifies the validity of the assumption that energy would not appear to be reversed from the conventional forward-in-time viewpoint and it means that if energy was not, in effect, reversed from the time-symmetric viewpoint, then from the unidirectional viewpoint it would actually appear to be reversed by T , which is certainly not desirable.

On the other hand, the physical space intervals and the momenta associated with the propagation of particles do need to be reversed (once) when time is reversed, if we insist on describing the motion of particles as it appears to take place from the conventional forward-in-time viewpoint and this despite the fact that only the physical time intervals experienced by the particles should actually be reversed by T . Indeed, given that the direction of space intervals is defined in relation to the direction of time intervals, if time intervals are followed in the wrong direction, then space intervals are also traversed in the wrong direction, so that the observed directions of the motion of particles are opposite the true directions of their motion, which means that those directions are actually reversed under a properly defined T operation, when the outcome of this operation is considered from a unidirectional-time viewpoint (this is made apparent when we reverse the direction of the arrows associated with the time-reversed states in Figure 3.1 to produce those in Figure 3.2). Thus, when the direction of time is reversed, but the time intervals in which the particles propagate are kept unchanged, as a consequence of practical limitations imposed by the thermodynamic nature of the observation process, the associated space intervals actually appear to be reversed (they are the negative of those really experienced by the particles), even though the spatial positions remain unchanged. This is true, again, despite the fact that at the most fundamental level of description, which is that of bidirectional time, the direction of space intervals is to be considered unchanged by a reversal of time. As a consequence, we obtain results which comply with the traditional definition of time reversal, according to which

momentum (and implicitly also space intervals) should, in effect, be reversed by T , along with angular momentum or spin, because given that momentum is here reversed independently from the position parameter x it follows that angular momentum would also appear to be reversed.

From the unidirectional viewpoint it may, in effect, seem like the traditional conception of time reversal, as involving a reversal of motion which simply allows the particles to follow a trajectory backward, could be valid. We must recognize, however, that just as there is no reason to assume that momentum is affected by a reversal of time from a bidirectional viewpoint (which explains that it is reversed from a unidirectional viewpoint), there is also no reason to assume that the sign of charge, as distinct from that of energy (the gravitational charge), would be affected from this same viewpoint when T is applied, because charge is not constrained to reverse by the requirement of action-sign invariance when the direction of propagation in time reverses. This may also appear to comply with traditional expectations, but in fact (as I previously remarked) it rather constitutes the one aspect which introduces a radical departure from what is normally assumed concerning time reversal. Indeed, it means that the same reversal that does apply to momentum from the unidirectional-time viewpoint, would have to apply to non-gravitational charges as well, because if the direction of propagation in time of the charges is actually reversed as required, then the fact that time is followed in the same forward direction relative to which the charges were originally propagating means that the charges would now appear to be reversed. We must, therefore, consider a reversal of charges to be associated with a reversal of time, as a result of the fact that this physical property is not experienced along the true direction of time in which it is propagated. This is a very important result which is definitely not expected from a traditional viewpoint, given that it asserts that a quantity which was previously assumed to be unaffected by a reversal of time (namely the sign of charge), would actually appear to be reversed under such a transformation, and if the preceding argument is valid then this conclusion would have to be considered unavoidable.

Thus, it seems that considering a reversal of time without assuming a consequent reversal of charge is incorrect and may give rise to violations of symmetry which are a simple artifact of the inappropriateness of traditional assumptions concerning which quantities are reversed along with the time coordinates, from the unidirectional viewpoint. To be meaningful, the experiments which seek to verify invariance under T would actually have to

assume a reversal of momentum and spin retracing a process backward, but combined with a reversal of charge (a permutation of particle and antiparticle). In other words, to test the invariance of physical laws under time reversal, we would have to use antimatter, which may explain why a violation of T symmetry is so difficult to observe despite the fact that violations of the combined CP symmetry were actually observed (which implies that T should also be violated, given that CPT is inviolable). It appears that we are simply not using the right kind of matter to probe for T violation. It is not the invariance of a process relative to the thermodynamic arrow of time which must be probed, but invariance under a reversal of the true directions of propagation in time of elementary particles. I believe that the improved consistency of the interpretation suggested here, from both an observational and a theoretical viewpoint, confirms that the traditional definition of time reversal as involving nothing more than a reversal of the directions of motion and rotation of particles can no longer be considered appropriate.

It may also be noted that, from a unidirectional viewpoint, the reversal of charge and the reversal of spin under a properly defined time-reversal operation are now the only aspects that differentiate this T operation from the P operation, apart from the respective reversals of the time and space coordinates themselves. But given that spin can also vary independently from the direction of propagation in time of a particle, this means that the only unmistakable distinction between the time-reverse of a given state and the space-reverse of the same state is, in effect, the sign of charge, which again emphasizes the importance of recognizing the dependence of this parameter on the direction of time. In such a context, it seems possible that the violations of T which may have been observed despite all the previously mentioned experimental difficulties could actually be violations of P symmetry, or violations of combined symmetries under which charge is left invariant by being reversed twice, because indeed those experiments do not compare matter and antimatter processes. Yet it might be considered that, despite what is commonly believed, violations of time-reversal symmetry had already been observed, even before the violations of traditional T symmetry were reported, because, as I will explain below, the C operation also involves some time reversal and violations of charge-conjugation symmetry do occur. In any case, it is clear that a violation of the time-reversal symmetry operation T , as it was here redefined, would not provide us with an absolute direction of time at a fundamental level, but merely with a preferred direction of time relative to some arbitrarily-chosen direction in space.

Another particularity of the alternative definition of time reversal proposed here is that it implies that it would now be electric fields which would reverse under application of the T operation, instead of magnetic fields, because electric fields depend only on the sign of charge of the source particles and charge must be assumed to reverse under time reversal. Magnetic fields, on the other hand, would now remain unchanged under time reversal, because from the unidirectional viewpoint the direction of motion of the source particles would reverse, as is currently understood, but charge would also reverse, despite what is currently assumed, so that currents (which are the source of magnetic fields) would remain unchanged as a consequence of being submitted to this additional reversal. We must, therefore, assume that a relative change between the direction of an electric field and that of a magnetic field does, in effect, take place under a properly defined time-reversal operation, only it is not attributable to a variation of the magnetic field, but rather to a variation of the electric field. The failure to recognize the dependence of the sign of charge on the direction of propagation in time of elementary particles, therefore, gives rise to an incorrect appraisal of the response of electromagnetic fields to a reversal of time.

A more consistent definition of the operation of time reversal, on the other hand, allows to avoid the troubling conclusion that certain phenomena involving electromagnetic fields would actually constitute a challenge to the necessary relational definition of discrete symmetry operations. Indeed, violations of time symmetry could arise, for example, in the case where neutrons would be observed to have an electric dipole moment and as such could effect a movement of precession around the direction of an external electric field, because this movement would appear to vary depending on the direction of time, but independently from the direction of the field and the sign of the electric dipole. However, while the direction of the dipole is not affected by the reversal of a neutron's spin angular momentum occurring as a consequence of the reversal of time, according to my proposal it would nevertheless be reversed together with it, because it depends on the sign of the constituent particles' electrical charges, which we must now also assume to be reversed as a consequence of applying the T operation. It is not possible, in this context, to assume that a reversal of time would allow a change in the precession motion of the neutron (associated with the direction of the neutron's spin) to occur *independently* from the direction of its electrical dipole in the presence of an invariant external electric field, because in fact both the spin and the dipole must be assumed to be reversed by T , along with the external electric

field. In other words, it is no longer possible to assume that while we should observe the precession motion to occur in reverse upon reversing time, the same dipole would nevertheless be interacting with the same electric field, as would happen if applying T actually reversed spin, but left the direction of the dipole and the external electric field unchanged. When the appropriate time-reversal symmetry operation is considered, only *relative* differences can occur between the direction associated with the precession motion and the direction of the dipole.

Still concerning the T operation, it must be clear that it is not possible to assume that what the traditional definition of this transformation involves is a reversal of the time coordinate that reverses physical time intervals and leaves energy unchanged, combined with a reversal of momentum that leaves both space coordinates and physical space intervals unchanged, even if that may not appear to disagree with the explicit definition of T as it is usually conceived. Such a definition of time reversal would be inapplicable, simply because it would reverse the sign of action of the physical systems involved and this is certainly not desirable knowing that negative-action matter (propagating positive energies backward in time) would be an entirely different kind of matter from a gravitational viewpoint and therefore certainly cannot be involved in those processes which we currently assume to be the time-reverse of processes involving positive-action matter. This has nothing to do with the fact that a unidirectional viewpoint is used traditionally. It is a different problem that would be unique to the T operation, despite the fact that I'm here assuming that C also involves some time reversal, because charge conjugation is simply not assumed to involve any space or time reversal traditionally and as such cannot be mistaken to involve action-sign reversal. From the viewpoint of unidirectional time, we can therefore only assume that the space intervals are reversed by T , along with the momenta, and that the time intervals, along with the energies, are left unchanged by the same operation, despite the reversal of the time coordinate. In other words, an appropriate (action-sign-preserving) time-reversal operation needs to reverse both momentum and space intervals together (from a unidirectional viewpoint) or leave them unchanged together (from the time-symmetric viewpoint) and those constraints must be explicitly stated in the definition of the symmetry operation. This, again, illustrates how important it is to identify the variability of all physical parameters under any discrete symmetry operation, in particular for what regards the sign of charge and that of energy, in relation to the direction of propagation in time, as otherwise we may mis-

interpret ordinary phenomena for potentially forbidden, symmetry violating occurrences.

3.6 The charge-conjugation operation

I think that, in the context of the preceding analysis, it becomes clear that the common assumption that time reversal amounts to simple motion (including rotation) reversal is what prevents a proper understanding of the nature of the charge-conjugation symmetry operation. The problem is that if we ignore the dependence of the observed sign of charges on the true direction of propagation in time of the particles carrying them, then this direction of propagation becomes impossible to assert, which explains that the existence of such a degree of freedom has traditionally been ignored altogether. Thus, I believe that the mistake we do when we consider time reversal as it is traditionally defined (even if we can now recognize that this error is not *only* a consequence of using a unidirectional viewpoint), is that we do not consider an evolution according to which the direction of propagation in time of particles is really reversed, but instead consider processes for which a series of events occur forward in time, merely in the reverse order to that in which they would otherwise be observed to occur. But given that non-gravitational charges are not affected by a reversal of the direction of propagation in time of the particles carrying them (which is distinct from the observed direction of their motion), it follows that we have a means to determine the direction of propagation in time of particles, which therefore becomes a meaningful, well-defined concept³. It would be incorrect to argue that only thermodynamic phenomena allow to distinguish a direction of time (even in the absence of violations of T symmetry), because the sign of charge is always dependent on the direction of time relative to which it is measured. It is simply the fact that the sign of charge itself cannot be characterized in an absolute manner that prevents a direction of time from being singled out as objectively distinct, in the way thermodynamic processes may appear to allow.

Now, what makes the acknowledgment of the existence of a relationship

³This conclusion is also justified by the fact that if an observer was ‘following’ the actual direction of propagation in time of an antiparticle, then this antiparticle would appear to have the same charge as its particle counterpart, but then it would be all the other particles in the universe which would appear to have a reversed charge, which is certainly a significant change.

between direction of time and sign of charge unavoidable is the recognized validity of the interpretation of antiparticles as particles propagating backward in time, which allows to identify reversal of time as the very cause of the apparent reversal of charge occurring from the unidirectional-time viewpoint. I believe, indeed, that despite what is often suggested, the interpretation of antiparticles as particles propagating in the opposite direction of time is not merely a helpful analogy with no real significance. Given the absence of a rational motive for rejecting the existence of a fundamental time-direction degree of freedom equivalent to the space direction degree of freedom and given the simplification made possible by the discussed interpretation of antimatter in a relativistic context, I think that we must recognize that there definitely exists a relationship between the direction of time and the sign of charge. But it must also be clear that, despite what is sometimes proposed, there is no equivalence between a reversal of *space* directions and a reversal of the sign of charge (which could imply that antiparticles are merely the enantiomorphic equivalent of their corresponding particles), even if there does occur situations when reversing the space coordinates may appear to counteract asymmetries associated with the sign of charge, because the relationship between space direction and sign of charge is, in fact, always a consequence of the existence of a relationship between the direction of space intervals and that of time intervals. In any case, if the relationship between time reversal and charge reversal which is suggested by the above-mentioned interpretation, is considered valid, then it would mean that the charge-conjugation symmetry operation must actually be understood as itself involving some time reversal.

What I'm proposing, therefore, is that we should recognize that the charge-conjugation symmetry operation C must actually be conceived as a combined space- and time-reversal operation that leaves the sign of non-gravitational charges invariant relative to the direction of time in which particles would be propagating following such a reversal. Thus, C must be understood to reverse the time parameter t (associated with the 'position' in time), along with the physical time intervals Δt associated with the propagation of particles, and the sign of the energy E of those particles (which is reversed as a requirement of action-sign invariance). But it must also reverse the space position parameter x , the physical space intervals Δx associated with the propagation of particles, and the momentum p of those particles (which is also reversed as a requirement of action-sign invariance). Here, again, we must recognize that the charge q is actually left unchanged, along

with the spin of elementary particles, from a fundamental viewpoint, even by this reversal operation we call charge conjugation. Yet it still makes sense to consider C as a reversal of charge, given that, from the viewpoint of unidirectional time, non-gravitational charges would appear to be one of the few physical properties of elementary particles which would actually be reversed by this symmetry operation, while the space and time intervals, along with the energies and the momenta would appear to remain unchanged.

This must happen for the same reasons that justified assuming that momentum and space intervals are reversed by T from a unidirectional-time perspective, even though they are left invariant by this symmetry operation from the bidirectional viewpoint. Indeed, upon applying C , we are in a situation where all intervals and their conjugate attributes are reversed from a fundamental time-symmetric viewpoint, which means that to satisfy the needs of a unidirectional perspective we must reverse the time-related parameters Δt and E again, but given the relationships that exist between the physical time intervals and the space intervals, this means that the space-related parameters Δx and p must also be reversed a second time, just as they were shown to be reversed (once) by T from this unidirectional viewpoint. If the physical time intervals and the energies must be reversed from what they really are (what they have become as a result of applying the operation in the first place) it is, therefore, due to the fact that from the unidirectional viewpoint we use the wrong direction of time, but given that following time in the wrong direction also implies that the space intervals are followed in the wrong direction (the relational aspect), then this actually means that the space intervals must also be reversed from what they really are (what they have become), along with the momenta. As a result, there appears to be no change to space and time intervals from applying C , even though it is here defined as a space- and time-reversal operation. Yet, as charge is not a spacetime related physical property, because it is associated with interactions distinct from gravitation (unlike energy or momentum, which can be conceived as the charges determining the metric properties of spacetime), it should actually remain unchanged, from the fundamental bidirectional viewpoint, under a space- and time-reversal operation such as the properly defined C , which means that it would appear to be reversed, as we would normally expect, from the unidirectional-time viewpoint (because time is then followed in the wrong direction).

There is a slight difference, however, between the outcome of a properly defined C operation and the expected outcome of a traditionally defined

charge-conjugation operation, because the reversal of the space and time position parameters x and t themselves (which now occurs from both the bidirectional- and the unidirectional-time viewpoint), even if it is without any effect on the sign of the space and time intervals associated with the propagation of particles from a unidirectional viewpoint (given that those intervals must then be reversed a second time), actually implies that angular momentum would appear to be reversed by C (because momentum is indeed unchanged, while the position in space is reversed). Thus, despite common expectations, a C -reversed process would also appear to involve reversed angular momentum or spin, which means that contrarily to what is sometimes suggested, the behavior of spin under charge conjugation is not a mere matter of convention and its reversal (apparent from a unidirectional-time perspective) must be considered an unavoidable outcome of applying this symmetry operation.

The reversal of spin under C is certainly unexpected according to the traditional approach, but from my perspective it appears natural, given that C involves a reversal of time. It must be clear, though, that this reversal of spin is only apparent and does not occur at the most fundamental level of description, in accordance with the requirement that an action-sign-preserving symmetry operation like C should not reverse the sign of action associated with angular momentum. This is to be required, even if, in general, the sign of spin is not uniquely tied to the sign of action associated with energy and momentum, because the only way spin can reverse is when either the position in space or the momentum are independently reversed and an action-sign-preserving reversal operation that reverses momentum would necessarily also reverse spatial position given that it must reverse space intervals (which is not required from the unidirectional-time viewpoint, relative to which momentum can be made to vary independently from the sign of space position, even when action is to remain positive).

We are now therefore in the situation where we must recognize that, from a certain viewpoint, charges are reversed by a properly defined time-reversal operation T , while spin angular momenta are reversed by a properly defined charge reversal operation C , despite what had traditionally appeared to be required from such discrete symmetry operations. Another distinction of the proposed approach is that handedness is now also reversed by C from whatever viewpoint, because either momentum is reversed and spin is invariant (as from the bidirectional viewpoint), or momentum is invariant and spin is reversed (as from the unidirectional viewpoint), so that there is always a

relative change between the direction of spin and that of momentum. The outcome of the proposed charge reversal operation C , as it was here redefined, would therefore differ from that of a properly defined T operation mainly through the fact that, unlike C , T would reverse the momentum and space intervals (from a unidirectional viewpoint), but would not reverse the handedness of particles, just as we would also expect traditionally. Thus, both the P operation and the redefined C operation would alone and from any viewpoint reverse the handedness. In this context, the fact that under certain circumstances, such as when the weak interaction is involved, particles of a given handedness seem to be naturally related to antiparticles with opposite handedness, could be understood to follow from the fact that the handedness is reversed by a properly defined charge-conjugation operation (which still relates particles to antiparticles), so that if there can be invariance under such a symmetry operation, then reversing both charge and handedness should not be expected to produce any change. This is an important result which confirms that the suggestion, usually made on the basis of purely phenomenological considerations, that charge conjugation should perhaps involve a reversal of handedness, was in fact justified from a theoretical viewpoint.

3.7 Invariance under combined reversals

I think that I have appropriately justified the inevitability of the above discussed conclusions regarding which parameters should be expected to reverse under the various discrete symmetry operations (in particular when I discussed the requirement of action-sign invariance and the constraint of relational definition of the reversal operations), but I must nevertheless mention how remarkable it is that the described variations of physical parameters under the redefined P , T , and C operations happen to be just such that they *explicitly* require invariance to occur under a combined PTC operation. This happens because all the parameters which are independently reversed by any of the symmetry operations are actually reversed twice when the operations are combined and this regardless of whether we are considering a unidirectional- or a bidirectional-time viewpoint (a look at Tables 3.2 and 3.3 allows to quickly confirm this fact). Either a parameter such as Δt is reversed twice, or else it is not reversed a single time by a properly defined PTC , and this actually guarantees that there is invariance under a com-

combination of the three discrete symmetry operations, because anything that may be reversed is reversed again and only once. In fact, as I will explain below, what we really need is twice a reversal of *all* fundamental space- and time-related parameters (that is both the time-related parameters t , Δt , and E , and the space-related parameters x , Δx , and p) under a properly defined *PTC* and this actually occurs when the appropriate bidirectional-time viewpoint is considered. Charge and spin, on the other hand, need not reverse at all from such a viewpoint, under a *PTC* operation, as they necessarily transform independently from the action-sign-preserving discrete symmetry operations and only reverse as a consequence of adopting a unidirectional viewpoint and in such a case they do reverse twice, as required. This is in contrast with the traditional definition of the discrete symmetry operations (described in Table 3.1) according to which some parameters, like the space and time coordinates, the charge, and the spin, can be reversed a single time only by the combined *PTC* operation.

We can understand, however, why it is that this combined symmetry operation should be expected to produce invariance, even as it is traditionally defined (as required by the CPT theorem). This is possible simply because, according to the traditional conception, while charge would be reversed only once (by C), spin would also be reversed only once (by T), but as one can show, there is a kind of equivalence, at least for fermions, between a reversal of the polarization state associated with spin and a reversal of charge and this is why, even as it is traditionally defined, the combined *PTC* symmetry operation would have to leave physical states invariant (although it would seem to alter the direction of space and time coordinates, which could turn out to be physically significant under particular circumstances). It is also interesting to observe that, in the context of my revised definitions of the discrete symmetry operations, any two operations applied together is *explicitly* equivalent to the remaining operation, so that applying PT , for example, is totally equivalent to applying C , which again demonstrates that charge conjugation must really be conceived as a space- and time-reversal operation and that time reversal must involve a reversal of charge from a certain viewpoint. What those relationships really show is that the discrete symmetry operations, as they are now defined, are all necessary and together sufficient to provide a complete account of the possible transformations involving a reversal of any of the fundamental properties of matter, aside from the sign of action (in fact, as I will explain in section 4.3, charge can also be reversed independently from any space- and time-related attribute, but

the states of matter so obtained usually do not interfere with the processes involving ordinary matter and antimatter particles).

In this regard, I must also mention that it is not possible to assume that applying either P or T alone, but twice, should necessarily produce invariance (in the sense that it would leave any system with no discernible change that could be related to unchanged physical parameters), despite the fact that it would appear to leave all parameters unchanged, because such a combined transformation may not leave the quantum phase associated with fermions unchanged, given that it would only be equivalent to a rotation in space by 2π radiant (as a single space reversal introduces a π radiant rotation and a single time reversal introduces an equivalent additional π radiant rotation in space) and only twice such a complete rotation would necessarily produce invariance in the presence of fermions. Of course, applying P or T alone, twice, would already be more likely to produce invariance than applying P alone or P combined with T only once, because at least some of the effects of applying P or T once would indeed be neutralized by a second application of the same operation, but the point is that, in such a case, invariance would not *necessarily* follow. The case of C is different, however, given that this operation involves a reversal of both space and time parameters, all at once, which produces an equivalent 2π radiant rotation with only one application (therefore allowing the changes involved to be related to the incomplete transformation of fermion wave functions), so that applying C twice reverses *all* parameters twice and introduces twice a 2π rotation, that must leave even the quantum phase of fermions invariant. The C operation, as I redefined it, is thus unique, because it is the only one of the three relationally distinct discrete symmetry operations that reverses both space- and time-related parameters together and from its alternative definition it can be seen that applying C is actually and explicitly equivalent to applying a combined PT operation. In this context applying PTC could be considered equivalent to applying PT twice, which clearly shows that the PTC operation involves a reversal of all parameters twice and is also equivalent to two complete rotations, which can only produce invariance.

In fact, any one of the three basic discrete symmetry operations can be considered as equivalent to a combination of the other two, so that T , for example, would here be equivalent to CP and P would be equivalent to the combined CT . Therefore, applying T twice would be equivalent to applying CP twice, which would amount to reverse both space- and time-related parameters twice (which considered alone would have to produce invariance)

and then also reverse space-related parameters twice (the order of application of the discrete symmetry operations in a combined operation has no importance and only the number of times a parameter is reversed is significant). But such a combined operation would not leave fermion wave functions invariant, for the same reason that applying P alone twice should not be expected to necessarily leave things invariant. It remains, however, that the fact that some combinations of basic discrete symmetry operations, which are not required to necessarily produce invariance do involve twice a reversal of some specific physical parameters, allows one to expect that an invariance which was lost when one of those fundamental operations was applied alone, can sometimes be regained by application of such combined operations. This should indeed be expected to occur given that, as I mentioned above, reversing one physical parameter twice, even if it is not guaranteed to leave all processes invariant, still allows the possibility of neutralizing some asymmetries which would occur as a consequence of the reversal of this single parameter.

What must be retained here is that there may be a difference between applying a symmetry operation twice and applying the *outcome* of this combined operation only once (which would amount to effect no change), even if in certain cases, as when the operation considered is the C symmetry operation, we would necessarily observe no change when the same operation is applied twice. This particularity of the C operation is merely a consequence of the fact that it reverses more individual parameters all at once, so that applying it in combination with itself actually allows to leave no parameter unchanged, relative to which an asymmetry could be properly defined. It must be understood, however, that despite their equivalence with combinations of distinct operations, the three basic operations defined above are all essential to a description of the allowed discrete transformations of physical parameters and none is more fundamental than any other. Indeed, two operations are distinct from a relational viewpoint, when one of them reverses one category of parameter, say space, relative to the other category, say time, while the other reverses another category of parameter, say time, relative to the previous one, say space, and each one of those operations is relationally distinct from yet another one that reverses both categories of parameters together and which constitutes the necessary complement to the other two operations.

3.8 The significance of classical equations

We can now return to the problem of understanding how it is possible for the momentum p to be left unchanged by a properly defined time-reversal operation T which, from the most fundamental viewpoint, must be assumed to reverse time intervals dt , but to leave space intervals dx unchanged. A problem would, in effect, appear to arise from the fact that, according to the classical equation that defines the momentum of a particle with mass m , we should have $p = m dx/dt$, which would clearly imply that if dt is reversed or negative, while dx is invariant or positive, then p should be negative, which is contrary to my proposal that both space intervals and momentum are unaffected by a reversal of time. But I believe that this contradiction is only apparent and a result of the fact that the classical equation for momentum is actually valid only from a unidirectional-time viewpoint, because it was originally introduced under the implicit assumption that physical properties are always measured in the conventional forward direction of time.

Indeed, what the classical equation is telling us is merely that, from the unidirectional viewpoint of an observer always following events in the unique direction of time associated with entropy increase and providing an account of physical quantities like momentum and space intervals in relation to that unique direction of time, relative to which time intervals dt are, in effect, positive definite, independently from the true direction of propagation in time of the particles involved, some quantities, like dx , which we might assume not to be reversed by T , are actually observed to be reversed, while dt itself is kept unchanged. Thus, if we use the viewpoint relative to which we are allowed to assume that the above equation is valid, then dt would actually remain positive definite, despite the reversal of time, while dx would have to be assumed reversed (for reasons I have already explained), which according to this action-sign-preserving classical equation would imply that momentum is also reversed, a conclusion that agrees with the definitions I provided for the unidirectional-time viewpoint and which is certainly appropriate, given that particles submitted to such a time-reversal operation must have unchanged momentum relative to the apparent (but false) direction of their motion, which is satisfied when both the momenta and the physical space intervals are together reversed.

There is no contradiction here, despite the fact that we must assume that the true signs of conjugate physical parameters, such as the space intervals and the momenta, are together invariant under a reversal of time from the

alternative time-symmetric viewpoint (according to which the sign of time intervals is itself reversed), because in such a case the classical equation no longer applies, simply because, as a conventional formula, it *never* really applied to such situations. The classical relation between momentum and the space and time intervals was deduced on the basis of the validity of a thermodynamic viewpoint of time and therefore does not apply in a context where time intervals are allowed to change sign. The classical equations are logical deductions, dependent on a certain viewpoint of time which must be considered inappropriate at the most fundamental level of description. In other words, it is not the validity of the classical equations in a limited context which implies that the assumptions made from a time-symmetric viewpoint (concerning the sign of physical quantities) are contrary to experimental evidence, but really the limited value of the classical equations which imply that the assumptions associated with a unidirectional viewpoint are not generally valid. We must recognize that the assumptions used in the more appropriate time-symmetric context, regarding the variations of space- and time-related quantities under a reversal of time, are not just theoretically well-founded, but that, under the right interpretation, they are fully supported by observations, while the variations deduced from a unidirectional-time viewpoint are explainable merely in the context where they are assumed to derive from the more fundamental bidirectional description.

It must be clear that, in this context, we would also be unjustified to make use of the classical formula for angular momentum \mathbf{L} , to which the spin of elementary particles is related, to decide what would happen to spin, from a fundamental viewpoint, under a reversal of time generated by a properly defined T or C operation. Indeed, the classical formula defines the angular momentum $\mathbf{L} = \mathbf{r} \times \mathbf{p}$ in terms of the position vector \mathbf{r} and the momentum $\mathbf{p} = m(dx/dt)\mathbf{i}$, and if we assume a reversal of time intervals dt to follow from both a T and a C reversal operation then according to this equation it would seem that \mathbf{L} should reverse under both types of time reversal, because, either dt reverses alone (as under a properly defined T), or else it reverses along with \mathbf{r} and dx (as under a properly defined C). But, as I already mentioned, and for reasons I have previously discussed, it would be incorrect to assume that angular momentum reverses under either T or C , from the bidirectional-time viewpoint relative to which dt does, in effect, reverse. Yet there is no problem here, because the classical formula is only right when we consider things from the unidirectional viewpoint, according to which dt is positive definite, but under such conditions either dx and \mathbf{p} reverse

together with unchanged \mathbf{r} (as occurs when T is applied), or else dx and \mathbf{p} are unchanged and \mathbf{r} is reversed (as occurs when C is applied and only space positions are reversed), so that in both cases spin angular momentum should actually reverse. Again, it must be emphasized that the incompatibility of the classical equation for angular momentum with the proposed definition of time reversal, as it occurs from a fundamental bidirectional viewpoint, cannot be considered to imply that the proposed fundamental definition is inapplicable, because all that it means is that the equation itself is of limited scope, having been developed in the context of a unidirectional perception of the evolution of physical systems, when it had not yet even been realized that there exists a fundamental degree of freedom associated with the direction of propagation in time.

3.9 Reversal of action

The clarification of the situation which was achieved in the preceding sections, regarding the interdependence of fundamental physical properties as they vary under application of any of the three essential discrete symmetry operations, has allowed to establish that that none of the traditionally considered discrete symmetry operations engenders a reversal of the sign of action. This is of course a consequence of the fact that, regardless of the viewpoint we adopt, those symmetry operations always reverse the sign of energy in combination with the sign of time intervals associated with the propagation of particles, just as they always reverse the direction of momentum in combination with the direction of space intervals. Thus, the T operation in particular, despite the ambiguity of its traditional definition, cannot be assumed to reverse the action, because, while it reverses the time coordinates and leaves the sign of energy unchanged from the unidirectional-time viewpoint, it is also implicitly assumed to preserve the sign of time intervals associated with the propagation of elementary particles. The role of inverting the sign of action must therefore be attributed to some symmetry operations distinct from all of those which are usually considered.

I have come to understand that there is not a unique single operation relating positive and negative action states, but that there are basically four different ways by which action can be reversed, which give rise to four different action-sign-reversing symmetry operations, whose four different outcomes are each related to phenomenologically distinct states of negative-action mat-

ter. If any one of those operations is applied independently from the others, it may not necessarily produce invariance. I will collectively denote those operations by the letter M to emphasize the fact that they constitute a different category of reversal transformations which are unlike those already studied. The states produced by those four distinct operations can be transformed into one another by individually applying each of the three action-sign-preserving symmetry operations P , T , and C , and therefore I will denote the various action-sign-reversing operations by applying the appropriate indexes corresponding to the operations which relate the states they generate to the state which is produced by one of those action-sign-reversing operations, chosen arbitrarily as the basic operation, which will itself be denoted M_I . The four discrete symmetry operations so defined are thus the M_I , M_P , M_T , and M_C operations displayed in Table 3.4. It must be clear, however, that the choice of which action-sign-reversing transformation must be associated with the basic operation M_I is completely arbitrary and we could, for example, have defined the operation originally denoted M_C to be the basic operation, which we would instead denote M'_I and we would then obtain the states produced by the other three operations by applying P , T , and C to the state generated by M'_I . That way it would appear that it is the redefined M'_C which would be equivalent to the original M_I , while M'_P would be equivalent to M_T , and of course M'_T would be equivalent to M_P and therefore we see that attribution of the indexes is purely a matter of convention. The letter M was chosen to denote action reversal, because the operations it represents would actually alter the gravitational properties of the matter submitted to such reversals and mass (which is usually denoted m) is the property that was traditionally associated with the gravitational interaction.

From Table 3.4 it is possible to see that there are two different ways by which a given type of fundamental physical parameter, either space- or time-related, can be reversed in such a way that the sign of action is reversed. We can either assume a reversal of the signs of momenta and energies relative to unchanged space and time intervals, or we can assume a reversal of the space and time intervals associated with the propagation of particles that would occur while keeping the signs of momenta and energies invariant. But given that those two different kinds of reversal can be applied differently to space- and time-related parameters (you can apply one kind of reversal to space and the other to time, or vice versa, as long as you do apply any one type of reversal to each type of parameter), it means that there are four different kinds of operations in all which can reverse the sign of action. From those

Bidir.	t	Δt	E	x	Δx	p	q	s	h
$M_I = M'_C$	t	Δt	$-E$	x	Δx	$-p$	q	$-s$	h
$M_P = M'_T$	t	Δt	$-E$	$-x$	$-\Delta x$	p	q	$-s$	$-h$
$M_T = M'_P$	$-t$	$-\Delta t$	E	x	Δx	$-p$	q	$-s$	h
$M_C = M'_I$	$-t$	$-\Delta t$	E	$-x$	$-\Delta x$	p	q	$-s$	$-h$

Table 3.4: Variations of physical parameters under the four relationally distinct action-sign-reversing symmetry operations, as described from the bidirectional-time viewpoint. Here I chose the basic action-reversal operation M_I to be that which reverses energy E independently from time intervals Δt , and momentum p independently from space intervals Δx . Under an equivalent definition it would be the time intervals Δt and the space intervals Δx which would be reversed by the basic action-reversal operation M'_I , while the energy E and the momentum p would be kept invariant.

definitions it is clear that what the M_I , M_P , M_T , and M_C operations really involve is the reversal of an additional degree of freedom, relationally distinct from those already affected by the P , T , and C operations, because, even the state obtained by applying the basic M_I operation actually involves a reversal of action, which means that all possible states related by application of P , T , and C , including the original state obtained by application of the identity operation I , have their counterpart as M -reversed states, and under such conditions we can only conclude that we are actually dealing with a transformation that applies to a distinct property of matter. The illustration of the effects of the various action-sign-reversing operations depicted in Figure 3.3 allows to clearly identify this degree of freedom as the relative orientation of momentum p compared to space intervals Δx or equivalently that of energy E compared to time intervals Δt , which for negative action states is the opposite of what it is for positive action states.

The C , P , and T operations, therefore, do not together operate a reversal of all fundamental physical parameters, because they merely reverse all parameters while leaving the sign of action invariant. The four action-sign-reversing symmetry operations proposed here are then the additional operations which are required to complete the set of discrete space- and time-related symmetry operations, because they perform the only remaining possible changes that the traditional operations do not produce, by actually

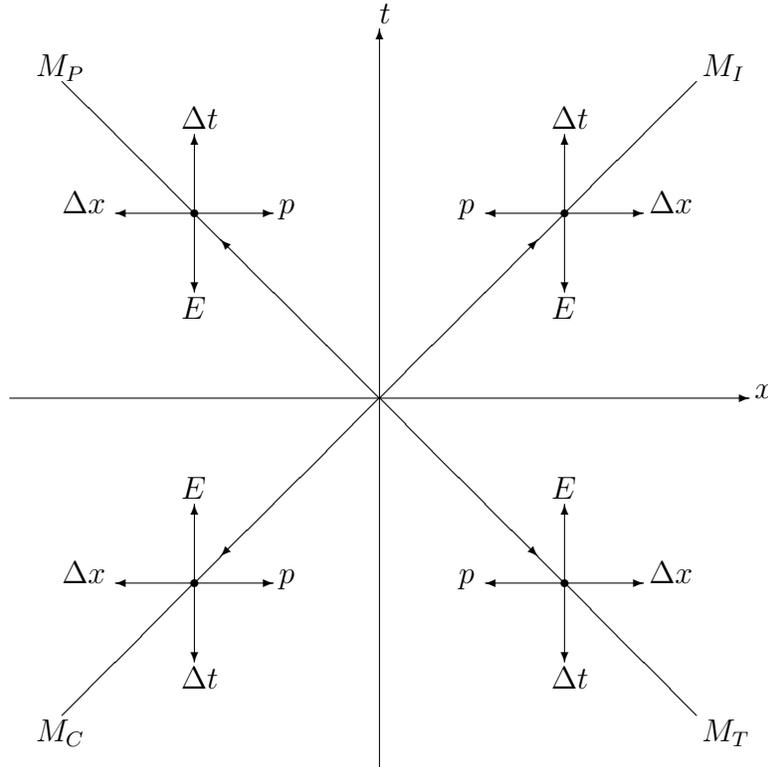


Figure 3.3: Four different outcomes of applying each of the relationally distinct action-reversal symmetry operations, as described from the bidirectional-time viewpoint. Here we notice that the orientation of the vectors which correspond to the signs of space and time intervals is always opposite that corresponding to the signs of momentum and energy, as we should expect to observe when action is indeed negative. If we were to consider a unidirectional-time viewpoint, we would have to reverse all space and time intervals and all momentum and energy signs for the processes obtained by application of both the M_T and M_C operations, which means that all four operations would give rise to the propagation of negative energies forward in time.

reversing the sign of momentum and energy relative to the direction of space and time intervals. From that viewpoint, it appears that even though they are usually ignored, the M_I , M_P , M_T , and M_C operations cannot in fact be avoided. The fact that there are actually four distinct operations that can perform a reversal of action, on the other hand, simply means that it is not possible to associate a unique state of momentum or energy, or of propagation in either space or time, to negative-action matter and that all the different action-sign-preserving variations of the direction of fundamental physical parameters which can be applied to positive-action matter, could also be applied to negative-action matter. We can, thus, actually expect that there would, for example, be a charge-conjugation symmetry operation C applying independently to negative-action matter, which would therefore have its own antimatter particles, distinct from ordinary antiparticles.

In this context, it appears that the distinction that exists between matter and antimatter must be attributed essentially to the true direction of propagation in time of particles, independently from their sign of action. An antiparticle is therefore always just a particle which reversed its energy while changing its direction of propagation in time, which is not very different from the situation of a particle which reverses its momentum by changing its direction of motion in space. Indeed, by reversing its momentum when it changes its direction of propagation in *space*, a particle is allowed to keep the sign of its momentum relative to the direction of its motion unchanged, so that its action sign is also unchanged, just like a positron retains the sign of action of the electron with which it annihilates, because the electron reverses its energy when it starts propagating backward in time (which is viewed as the annihilation process forward in time). But a negative-action particle would be clearly distinct in this respect, as a consequence of the fact that it would not merely carry negative energy forward in time (or positive energy backward in time, which is equivalent from a unidirectional-time viewpoint, when the sign of charge can be ignored), but also negative momentum in the observed direction of its propagation in space (the momentum would point in the direction opposite the observed velocity of the particle), unlike any ordinary matter particle (including antiparticles). It must be clear, however, that according to the proposed definition of action-sign-reversing symmetry operations which is described in Table 3.4, non-gravitational charges are assumed to be unaffected by a reversal of action, just as they were left invariant by the action-sign-preserving reversal operations. Only the practical necessity of a forward-in-time viewpoint would, for negative-action matter, also

imply that charges appear to be reversed when a process is submitted to an action-sign-preserving reversal of time.

Another particularity of the operations of action reversal defined above is that spin is deduced to be reversed under all such relationally distinct operations, when their effects are considered from the bidirectional-time viewpoint. This is certainly just as appropriate as the invariance of spin that is produced by all action-sign-preserving symmetry operations, because, as I previously mentioned, spin has the dimension of an action and should therefore vary in correspondence with the sign of action associated with momentum and energy, from a fundamental viewpoint. The constraint on the variation of the direction of spin is actually the same constraint that requires that, either both space- and time-related parameters are such as characterizing a positive action state, or else that they are both such as characterizing a negative action state, and that it should not be possible for one single particle to propagate, say, positive momentum in the direction of its motion in space and at the same time propagate negative energy forward in time. This is a simple matter of consistency, because a physical system cannot have at once both the gravitational properties associated with positive-action matter and those associated with negative-action matter if, as I suggested in the previous chapter, the attractive or repulsive nature of the gravitational interaction between two particles actually depends on the difference or identity of their action signs. This does not mean, however, that spin cannot vary independently from the sign of action associated with energy and momentum, but merely that while it cannot reverse as a consequence of applying an action-sign-preserving discrete symmetry operation, it also *must* reverse as a consequence of applying a reversal of action.

It may also be noted that just as is the case for the action-sign-preserving discrete symmetry operations, some combinations of two of the four operations describing a reversal of action are equivalent to a combination of the other two operations (in the case of the action-sign-preserving operations, one operation, which is that of charge conjugation C , was equivalent to the other two, but in fact this single operation was implicitly combined with the invariant operation I , which produced no additional change and thus could be ignored). Here, a combination of M_I and M_C or a $M_I M_C$ operation, would be equivalent to a combination of M_P and M_T and this is what allows a combination of all action-sign-reversing symmetry operations (or a $M_I M_P M_T M_C$ operation) to necessarily produce invariance, given that all relevant parameters are actually reversed twice by such a combined operation.

In fact, it turns out that combining any of M_P , M_T , or M_C with M_I produces an operation equivalent to the above defined P , T , or C , respectively (while a combination of M_I with itself produces an operation equivalent to the identity operation I), so that a combination of the other two remaining action-sign-reversing operations would also be equivalent to those action-sign-preserving operations. For example, the combined $M_P M_T$ operation is mathematically equivalent to a C operation, because it reverses both space- and time-related parameters once and reverses the action twice, which is equivalent to leave action unchanged.

One must understand, however, that even though applying any one action-sign-reversing operation twice would be equivalent to applying the identity operation I , such a combined operation would not necessarily produce invariance and this for the same reason that applying P or T twice would not necessarily leave everything invariant, despite the fact that it would also appear to be equivalent to applying the I operation, which produces no change. This is, again, because applying an operation that does not reverse all physical parameters twice, even if it may appear to return a system to its original state, may still produce a change which can be characterized in a relational way, because some parameters would be reversed relative to other parameters which remain unaffected by the transformation and this may not leave the processes involved invariant. Still regarding the conditions necessary to ensure invariance, it should be clear that simply combining a PTC operation with the basic M_I or any other action-sign-reversal operation, as a way to try to regain invariance which may be lost upon reversing the action (in the way we would apply T to a CP violating process), cannot be expected to produce invariance, given that the action-sign degree of freedom would then be reversed only once. Thus, a violation of any of the M symmetries would not imply that there must be a violation of PTC symmetry, as we may understand to be independently required on the basis of the fact that invariance under PTC alone must itself be considered unavoidable. The appropriate generalization of the PTC (or really $IPTC$) symmetry must then be recognized to be the $M_I M_P M_T M_C$ symmetry, which combines all the relationally distinct action-reversal symmetry operations and which must therefore be equivalent to no change at all (because there would remain no unchanged physical parameter relative to which a change could be determined). Indeed, as indicated in Table 3.4, a physical parameter may either not be reversed by any of the action-reversal operations, or else be reversed by two or all of those symmetry operations, which explicitly guarantees invariance under a

combination of the four operations.

Now, in order to avoid confusion, it is important to understand that the action-sign-reversal symmetry operations must be considered as operations distinct from one another that apply to an identical state, rather than as an identical operation that applies to different states. In such a context, it transpires that the fact that the M_I , M_P , M_T , and M_C operations are related to one another through application of the various action-sign-preserving symmetry operations, merely shows that the states obtained by applying the four action-sign-reversing operations are themselves related to one another through the same action-sign-preserving operations that transform unchanged action-sign states into one another. What must be clear, then, is that no action-sign-reversing symmetry operation can be identified as *the* action-reversal operation and under such circumstances it is not possible to avoid having to consider the many operations as distinct from one another, despite the fact that all such operations can be obtained by combining, in turn, each of the action-sign-preserving symmetry operations with just one single action-reversal operation.

Action-reversal symmetry can, therefore, be violated to different degrees when one transforms a state of positive-energy matter into the different states of negative-energy matter which are related to one another by the redefined action-sign-preserving reversal operations P , T , and C , because each of those states is related to a corresponding state of positive-energy matter by a specific action-sign-reversing symmetry operation and these operations do not necessarily produce invariance when applied separately. Thus, the P , T , and C operations can be violated to different degrees by negative-energy matter (compared to how they are violated by positive-energy matter), when applied independently from one another, and this precisely because M_I , M_P , M_T , and M_C can themselves be violated to different degrees in comparison with one another, so that they relate the different asymmetric states of positive-energy matter to corresponding states of negative-energy matter which can be asymmetric in different ways relative to one another. The only requirement is that the different states of negative-energy matter which are related to the different states of positive-energy matter by the various action-sign-reversal symmetry operations be subject to the same invariance under a combined PTC transformation as are states of positive-energy matter, even if P , T , and C are violated to different degrees by negative-energy matter, in comparison with the violations occurring for positive-energy matter. The four action-reversal symmetry operations, therefore, simply allow to relate

all the positive energy states, which are transformed into one another by the action-sign-preserving symmetry operations, to all the negative energy states which are transformed into one another by similar operations. Thus, despite the existence of four distinct action-sign-reversal symmetry operations, action reversal must really be conceived as transforming one single degree of freedom and this means that I'm justified in referring to the action-reversal operations collectively as the M symmetry.

In any case it appears that the commonly met remark, to the effect that gravitation is invariant under a reversal of time, must be nuanced. What I mean is that, while it is certainly true that there would be no change to the attractive or repulsive nature of the gravitational interaction if time was locally reversed for some physical system by a time-reversal operation such as T , we should certainly expect a reversal of time independent from the sign of energy, such as that produced by an M_T operation, to exert a change on the nature of the interaction of the affected system with the rest of the universe. Indeed, such a transformation would reverse the sign of action and as I previously explained, the repulsive or attractive nature of the gravitational force between two bodies depends on the relative value of their action signs (because gravitation is always attractive only for particles with the same sign of action). But, even if we consider the reversal of time produced by an action-sign-reversing operation like M_T to apply to the whole universe (in which case we would have to use negative-energy matter in place of positive-energy matter when testing for invariance), the preceding discussion made clear that we should not necessarily expect to observe phenomena which would be completely similar with those of the original universe, because M_T applied alone could be violated, just as any operation which is not reversing all physical parameters twice. This would also be true of M_P , for example, because, just as the change in the sign of time intervals produced by an M_T operation can be related to an unchanged sign of energy, so the change in the direction of space intervals produced by an M_P operation can be related to an unchanged direction of momentum.

Yet the fact is that there *could*, in effect, be invariance under a reversal of time that does not preserve the sign of action, if the operation is applied to all particles in the universe, because in such a case the difference or the identity of the signs of action of the various particles would not be affected and this is the only aspect that would be significant from a gravitational viewpoint. But this invariance would apply only to the extent that there is, in effect, no violation of symmetry under exchange of positive and negative

action states. It is important to mention, however, that even if one might be tempted to conclude, based on a certain interpretation of the generalized gravitational field equations which were proposed in section 2.15, that the minute imbalance responsible for the observed small, but non-vanishing positive value of the cosmological constant arises from such a violation of M symmetry, this would not be a valid conclusion, because, as I will explain in section 4.2, this imbalance rather develops as a consequence of the fact that the rates of expansion of space experienced by observers of opposite energy signs are allowed to differ as time goes, even if they were initially the same and this can occur even in the absence of a violation of M symmetry. Also, as I have explained in the preceding sections of this chapter, simply reversing the direction of motion of particles cannot be considered to consist in a true time-reversal operation in any meaningful way, so that assuming that such a transformation would leave all processes unaffected, even when gravitation is involved, could not be understood to mean that gravitation is invariant under time reversal.

3.10 Black-hole entropy

We are now entering the realm of a more uncertain domain of scientific inquiry, where classical gravitation theory reaches the limits imposed by quantum indeterminacy. In order for the following discussion to be meaningful it will first be necessary to recognize that the theoretical justifications and the observational evidence for the existence of black holes is sufficiently well established that these objects can be considered legitimate subjects of study. The objective I will try to achieve is then simply to show that it is possible to identify the degrees of freedom of matter which give rise to the exact measure of black-hole entropy derived from the semi-classical theory of black-hole thermodynamics. This explanation will be based on the results achieved in the previous sections while deriving an improved formulation of the discrete symmetry operations, as well as on a better understanding of the implicit assumptions entering the derivation of the semi-classical formula for black-hole entropy. More specifically, I will explain that based on certain plausible hypotheses concerning the constraints that should apply on matter particles approaching a spacetime singularity, it is possible to deduce that a finite number of discrete degrees of freedom characterizes the microscopic state of the elementary particles which were captured by the gravitational field of a

black hole. As a consequence, it becomes possible to actually confirm the existence of an exact relationship between those matter degrees of freedom and the binary measure of missing information or entropy which, according to the semi-classical theory, should be distinctive of those situations in which event horizons are indeed present.

I will be working here under the hypothesis (now usually recognized as appropriate) that the information concerning the matter which produced the gravitational collapse that gave rise to a black hole (or the matter which was later captured by the same object) is not lost, but is rather encoded in the detailed microscopic configuration of certain degrees of freedom associated with microscopic elements of surface on the event horizon of the object. Ignorance of this microscopic configuration, when a black hole is described using the classical macroscopic physical parameters of total mass, angular momentum and charge, is what gives rise to gravitational entropy. What is not fully understood presently is how we can reconcile the fact that matter appears to be characterized by physical parameters that vary in a continuous fashion, while the information contained in the microscopic degrees of freedom on the surface of a black hole must be given in binary units. What is the exact nature of the microscopic degrees of freedom of matter which would correspond with the missing information encoded in the microscopic degrees of freedom present on the event horizon of a black hole? Given the limitations imposed by the Bekenstein bound (according to which the amount of information that can be obtained concerning the microscopic state of the matter contained within any surface is proportional to the finite number of elementary units of area on the surface), it would appear that this question actually applies to the microscopic configuration of matter under *any* condition, regardless of the strength of the gravitational field on the surface through which information about this exact state must be obtained.

It therefore seems that the problem of identifying the fundamental degrees of freedom of matter which are associated with the binary measure of information encoded on a two-dimensional boundary is not one that concerns only those situations in which black holes are present, even though its significance is made more obvious when we are actually dealing with event horizons. I think that the fact that there is a similar limit, regarding the measure of information, for both event horizons and ordinary surfaces means that we must admit the reality of what would be occurring beyond the limits of any event horizon, despite the fact that the processes involved cannot be subject to direct observation (as I will explain below, this is only true in the sense

that a remote observer cannot obtain that information). Thus, regardless of the practical limitations which clearly exist for actually determining the exact state of whatever microscopic degrees of freedom are to be associated with the particular measure of missing information encoded on the surface of a black hole, this problem should nevertheless be considered a tangible one, even if only because, in the case of a normal surface, information about this microscopic state could be obtained (even by a remote observer). In fact, I believe that the constraints imposed by quantum theory, concerning the conservation of information, require that we recognize the reality of the microscopic degrees of freedom which encode all the relevant information about the matter which was captured by the gravitational field of a black hole and whose existence appears to be necessary for the consistency of the semi-classical theory.

Indeed, it is merely the classical nature of the general-relativistic description of the event horizon of a black hole which makes it incompatible with the hypothesis that information must be conserved for matter that is captured by the gravitational field of such an object. But once we recognize that this smooth and uniform description of the gravitational field is no longer appropriate on a microscopic scale, it becomes apparent that there is no basis to the commonly held viewpoint that the process of black-hole evaporation involves fundamental, irreducible irreversibility, or that information is actually lost when a black hole decays through the emission of Hawking radiation. There is no more reason to believe that information is lost when black holes evaporate, than there would be to assume that the information that appears to be lost when a drop of ink spreads into a liquid is fundamentally irretrievable. Later in this section, I will further explain what motivates this conclusion and it will become perfectly clear that there is no rational motive for assuming that the evolution that takes place when matter is captured by the gravitational field of a black hole is fundamentally (rather than statistically) irreversible, even when we recognize that there is something more objective about the growth of entropy that is associated with the formation and the evolution of event horizons.

What the semi-classical theory of black-hole thermodynamics implies is that there does exist information about what lies behind event horizons, but that this information is missing from the description of a black hole in terms of its classical macroscopic parameters and therefore we must assume that it could only be obtained through measurements of the microscopic configuration of some physical parameters associated with the surface delimited by

the event horizon of the object. The fact that a consistent theory of black-hole thermodynamics actually exists means that we have no reason to expect that, when such objects are involved, there could be departures from the rules which govern ordinary physical systems with a large number of degrees of freedom, for which it is already recognized that any apparent information loss merely occurs as a practical limitation. In the context where it is understood that, from a physical viewpoint, information must involve a distinction, this assumption is actually supported by the existence of a relation between the mass of a black hole and its entropy, because any distinctive features must be carried by elementary particles and when the number of particles absorbed by a black hole grows, its mass necessarily becomes larger. This observation would remain significant even if it was determined that the actual microscopic degrees of freedom which are allowed to vary for matter that fell into a black hole do not consist of mere energy differences. Also, if we recognize that information, as a measure of physical distinction, can be conserved without the knowledge of some such distinction being shared by any specific observer, then we are certainly allowed to assume that information persists, even when black holes are involved.

Some well-known results appear to confirm that the information concerning the microscopic state of the matter which was captured by the gravitational field of a black hole may, in effect, be encoded in the detailed configuration of certain degrees of freedom associated with the event horizon of the object⁴. Those conclusions are all dependent, basically, on one assumption, which is that there is a finite, maximum level of accuracy applying to measurements of spatial distances. This limitation would then also apply to the description of surfaces, such as those which are associated with event horizons. Indeed, the still largely uncertain quantum gravitational theories which were used to achieve those results all have as a key characteristic that they involve a discrete description of physical space on the shortest scale. Based on what I have learned concerning this issue, I can safely argue that it is this unique particularity of current quantum gravitation theories which allows to explain that they can predict that black-hole event horizons are characterized by a finite number of microscopic degrees of freedom which vary as binary parameters and which appear to encode the information about the unknown

⁴One of the most significant results of loop quantum gravity is that it allowed to identify those microscopic degrees of freedom as being given by the quantum numbers associated with the intersections of the edges of the spin network that describes the discrete structure of space on the quantum gravitational scale with the event horizon of the black hole.

microscopic state of the matter contained within the objects. Current quantum theories of gravitation would therefore have succeeded in unveiling at least one distinctive aspect of the structure of space and its associated gravitational field, in the context where quantum indefiniteness can no longer be ignored.

What was learned, more exactly, is that two events must be considered indiscernible, from the viewpoint of any measurement, when they would occur within intervals of space and time smaller than the natural scale of quantum gravitational phenomena. We now understand that trying to describe the state of matter and energy at a level of definition of spatial distances and during time intervals more precise than those provided by the Planck scale would constitute a superfluous characterization of physical reality. But, despite the fact that this constraint now appears clearly inescapable, it is still often ignored, as when someone is talking about what may have happened at a time shorter than the Planck time after the Big Bang. Here I will assume that the limitations imposed by quantum indeterminacy, which imply the existence of a smallest meaningful spatial distance, constitute a fact which will gradually become as well established as the existence of elementary particles of matter and on which further insights can, therefore, be based. In such a context, it would appear that, if the microscopic quantum gravitational degrees of freedom associated with the event horizon of a black hole are to be encoded in some discrete elements of space on its surface, then under no circumstances could two particles actually be present at the same moment in such a unit of surface. It would, therefore, be impossible for any physical attribute of some matter or gravitational field particle present within such a unit of area to be attributed more than one value at any particular time (although it remains to establish what is the exact size of this fundamental unit of area and therefore it is still possible that what may now appear to be a fundamental unit of area, would actually allow to encode more than one fundamental degree of freedom, as I will suggest in section 4.3).

Thus, it seems that it is from discrete elements of structure with a size of the order of the Planck interval that a proper description of the exact configuration of the microscopic degrees of freedom associated with the event horizon of a black hole can be formulated, that may also be valid to some extent in the case of ordinary surfaces. What is remarkable is that it appears that the physical parameters associated with those microscopic elements of surface also vary in a discrete way, which means that they actually provide a binary measure for the entropy, or missing information, which character-

izes those objects. Indeed, the relevant microscopic degrees of freedom on a surface can only be this or that, or yes or no, rather than assume any value from a continuous spectrum of possibilities as we go from one discrete surface element to the next. It appears that not only must we accept that space is divided in elementary units on the shortest scale, but we must also recognize that the values taken by the physical parameters associated with those discrete elements of surface can only be either one thing or another and nothing in between. Therefore, at the most fundamental level of description, it would appear that the physical properties of a surface must be described using discrete elements of structure corresponding to the smallest physically meaningful measures of area, to which are associated only two possible states of some microscopic degree of freedom. In such a context, the entropy of a black hole would derive merely from an absence of knowledge of the detailed configuration of those microscopic degrees of freedom (characterizing elements of surface on its event horizon) which arises as a consequence of the existence of insurmountable practical limitations, which effectively prevent any remote observer from obtaining experimental data about what is actually occurring at this level of precision of measurement, whenever an event horizon is, in effect, present.

In a semi-classical context one may assume that the elementary quantum gravitational degrees of freedom present on the event horizon of an isolated macroscopic black hole in empty space would be contained in the microscopic configuration of its surface gravitational field. It is necessary, however, to distinguish between the degrees of freedom characterizing the states of the particles which were captured by the gravitational field of a black hole and the degrees of freedom on the event horizon of the object, which merely reflect the microscopic state of the matter and which may be of a different nature from a physical viewpoint. But despite the ambiguous nature of the relationship between the physical degrees of freedom which allow information to be encoded on the event horizon of a black hole and the exact state of the matter from which an observer has become separated as a consequence of the presence of this theoretical boundary, it must be assumed that such a unique correspondence exists. What's more, given the size of the elementary units of surface on which the information concerning the microscopic state of the matter contained inside an event horizon is encoded, it appears that we would be justified to assume that the degrees of freedom of matter which we must identify are those which would apply to a description of matter at the Planck scale.

In any case, I think that the existence of such a correspondence between the microscopic degrees of freedom associated with an event horizon and those of the matter it contains should be considered unavoidable, even if only because we can never get more information concerning what is located beyond any surface than is obtainable by observing through this surface. But if there does, in effect, exist a limit to the accuracy of measurements that can be effected on a surface (due to the existence of a smallest meaningful spatial distance), then it necessarily follows that there must be a limit to the amount of information that could be obtained through a detailed probing of the processes actually occurring on that surface and this limit should naturally be expected to be proportional to the number of discrete surface elements, through which the information must flow. It should not come as a surprise, therefore, that it is the total area of a black hole which actually provides a measure of the number of elementary units of missing information which should ultimately be related to the exact microscopic state of the matter which is located past the event horizon of the object. What's more difficult to explain is why this constraint does, in effect, appear to be relevant to what is actually taking place beyond event horizons, rather than merely to what we can *tell* about what is going on there. Despite the enduring uncertainty associated with this question, I believe that the following discussion will help clarify the nature of the relationship between the microscopic degrees of freedom on a surface and the microscopic state of the matter located within that surface.

Before I undertake the task of explaining why it is that the states of the elementary particles which have been absorbed by a black hole can become so constrained that they are allowed to match the required binary measure of missing information which is encoded on the event horizon of the object, it would be appropriate to first recall what the semi-classical analysis of black-hole thermodynamics has revealed. What we know, in effect, is that for a non-spinning black hole of mass m with an event horizon of area $A_{BH} = 4\pi R_S^2$, where $R_S = 2mG/c^2$ is the Schwarzschild radius of the black hole, the entropy is given by $S_{BH} = \frac{1}{4}A_{BH}/A_P$, where $A_P = l_P^2$ is the Planck area given in terms of the Planck length, which is defined as $l_P = (\hbar G/c^3)^{1/2}$ and the units are chosen so that Boltzmann's constant k is equal to unity. In general, a black hole would therefore have an entropy that is determined by the value of the area of its event horizon in Planck units of surface, divided by a factor of four. Given that entropy is simply a measure of the information that

is missing from the description of a black hole in terms of its macroscopic parameters of mass, radius, or area, it seems that the amount of information encoded in the unobserved microscopic degrees of freedom characterizing the surface of the object is equal to one fourth its area in natural units. It was pointed out by Gerard 't Hooft, before the previously mentioned results obtained from quantum gravity were derived, that this actually means that information appears to be encoded on the surface of the black hole in binary units corresponding to an area equal to four Planck areas.

Now, if we are willing to accept that the Planck unit of area may actually be given as equal to $A_P = 4\pi l_P^2$ (following the traditional formula for the area of a sphere in terms of its radius) when the mass of a black hole approaches the Planck mass (from higher values associated with macroscopic event horizons), then an interesting result can be shown to follow. Indeed, using the above equation for A_{BH} , I can deduce that the event horizon of what we may call an *elementary black hole*, with a mass equal to exactly one Planck mass $m_P = (\hbar c/G)^{1/2}$, should have an area that is actually equal to four such Planck areas⁵. Using the formula for the entropy of a conventional black hole I would thus be allowed to conclude that the detailed configuration of the microscopic degrees of freedom on the surface of a Planck-mass black hole must carry one single binary unit of information. I think that the outcome of this simple derivation is extremely significant, because, on the basis of the hypothesis that there can be no significance in attributing existence to a particle which would occupy a volume smaller than that which is associated with the most elementary unit of area (as current quantum gravitational theories appear to require), it seems necessary to assume that such an elementary black hole, would contain at most one single elementary particle and in such a case we have no choice but to attribute the information encoded in the microscopic degrees of freedom on the surface of the black hole to its matter content.

But if, in the case of an elementary black hole at least, the missing information encoded on the event horizon of the object must definitely be associated with the single Planck energy particle it contains (which need not necessarily have a large rest mass), then, even for a black hole of larger mass,

⁵In section 4.3 I will argue that there need not be a problem with the fact that the value obtained from such a semi-classical approximation for the elementary unit of area is not exactly the same as the expression one currently obtains for the parameter known as the *area gap*, based on purely gravitational aspects of Planck-scale physics, using the mathematical framework of loop quantum gravity.

it should be possible to associate this binary information with the states of matter particles contained within the surface, despite the fact that, according to the above equations, the entropy of a black hole S_{BH} is not in general proportional to its mass m , but rather to its mass squared (so that entropy rises faster than the matter content). The fact that no simple relationship between entropy and matter content appears to exist in the more general case of a macroscopic black hole is simply due to the fact that the gravitational field must itself carry a portion of the entropy when large accumulations of matter are involved. However, in the context where the particles mediating the gravitational field are to ultimately also be understood as being a form of matter (in the most general sense), we would have no choice but to associate the entire amount of missing information associated with a black hole's event horizon with the 'matter' content of the object, which would then include gravitons. In any case, if an elementary Planck-mass black hole, containing a most elementary particle, with an energy of the order of the Planck energy, can be associated with the smallest unit of information, then it is necessary to recognize that the binary nature of the microscopic degrees of freedom on the event horizon of any black hole is a reflection of the existence of states of matter which can only vary in a discrete way.

So, what are exactly those degrees of freedom prevailing for matter trapped by the gravitational field of a black hole? When we ask this question in the context where the information associated with an *elementary* black hole is understood to provide a complete description of the state of a most elementary particle in the conditions where an event horizon is constraining the motion of this particle, it appears necessary to assume that this state must be completely definable by one single binary unit of information. It may therefore appear that we should be seeking to identify a unique physical parameter, that reverses under a given discrete symmetry operation, as being the binary degree of freedom related to the information encoded on the event horizon of our elementary black hole. But if we are to assume that the same fundamental parameters characterize the spacetime-related properties of matter under all conditions, then it rather seems that all the truly independent discrete symmetry operations, like the previously defined T , P , and M operations, should have their counterpart in the information necessary to characterize the state of a particle constrained by the gravitational field of an elementary black hole. Indeed, all of those reversal operations allow to distinguish the sign or the direction of some physically significant property of elementary particles and there is no *a priori* reason why only a subset of

those variable properties should need to be taken into account in the characterization of the discrete degrees of freedom applying at a fundamental level in the presence of an event horizon.

It must be clear that if all of the independent discrete symmetry operations were considered to determine one distinct degree of freedom of a particle confined by the event horizon of an elementary black hole then we would need not one binary unit of information or one bit to be encoded in the microscopic configuration of the gravitational field on the surface of the object, but rather three bits. Indeed, with two yes or no questions we can determine the action-sign-preserving direction of time intervals (reversed by T or not reversed) and the action-sign-preserving direction of space intervals (reversed by P or not reversed), which already allows to distinguish four states of matter (identity being the state where neither space nor time is reversed). The distinctions which exist between each of those four states as they appear from the bidirectional- and the unidirectional-time viewpoints are illustrated in figures 3.1 and 3.2. With an additional yes or no question we can then determine the sign of action (reversed by M or not reversed), which doubles the number of states of matter that can be distinguished, so that we can differentiate between the eight possible states of matter related by the discrete symmetry operations defined in the preceding sections. The C symmetry operation, being a combination of T and P , does not provide an additional distinct degree of freedom and therefore need not be considered here (even though we may as well consider only T and C or only P and C to provide the relevant discrete degrees of freedom and then it would be P or T which could be ignored). But three bits is not equal to one bit and so it may seem that there is a problem with associating the missing information encoded on the surface of a black hole with the degrees of freedom transformed by the discrete symmetry operations, despite the fact that those parameters should, in effect, characterize the states of elementary particles under all circumstances.

However, I believe that this discrepancy cannot be assumed to rule out the validity of the theoretically unavoidable conclusion that any binary distinction between the states of the matter particles that crossed the event horizon of a black hole must be a reflection of the structure underlying the previously defined discrete symmetry operations, which together allow to operate a reversal of all space- and time-related attributes of matter particles. I will show that very restrictive constraints actually limit the variability of certain microscopic physical parameters, whenever black holes are involved. Those

limitations imply that some microscopic physical parameters are restricted to a subset of the values they would otherwise be allowed to take. This actually contributes to reduce the amount of information needed to specify the microscopic states of particles trapped by the gravitational field of a black hole. Further insights will be needed, however, to allow the number of binary units of information required for achieving this complete description of the state of gravitationally collapsed matter to be made entirely compatible with the measure of black-hole entropy derived from the semi-classical theory of black hole thermodynamics.

In order to clarify the situation regarding what variations are allowed for the various physical properties of elementary particles when matter has become confined by the gravitational field of a black hole, we may first recall that the three macroscopic physical parameters characterizing a black hole are its total mass m , its total charge q , and its overall angular momentum j . To those three parameters, I would add the momentum p , which is not usually considered to define the macroscopic state of a black hole (given that the object can always be described in the reference system relative to which it is at rest), but which I believe provides essential information required to identify the parameters which must be taken into account in defining the *microscopic* state of the matter that form such an object. It must be clear that each of those macroscopic parameters must be allowed to vary not just in magnitude, but also in sign or in direction. The total mass m , in particular, must be conceived as being either positive or negative, depending on the overall sign of energy of the black hole. This is also an aspect that is usually not taken into consideration in the conventional treatment of black-hole thermodynamics, but which must be recognized as a necessary assumption in the context where the existence of negative-energy matter is theoretically unavoidable.

A different question would be to ask whether the sign of energy or action is a variable parameter for the particles forming a black hole. Given that I have already argued that negative-energy matter cannot be absorbed by a positive-energy black hole, it would seem that only positive energy states need to be taken into account in describing the microscopic configuration of the matter that was captured by the gravitational field of a positive-mass black hole. One cannot assume, however, that all black holes with a given mass sign must, at all times, be formed only of particles with the same mass sign as that of the object itself, because even if no particle of energy sign opposite that of a given black hole can cross its event horizon from the outside, it

is indisputable that a positive-energy black hole with a very large radius and a rather low density could form, despite the initial presence of some comparatively small amount of negative-energy matter inside the surface that is to become its event horizon. Thus, it is not strictly forbidden, for a macroscopic positive-energy black hole, to contain negative-energy matter, even though this matter would only be allowed to be present inside the event horizon associated with such a black hole if it was already contained inside the surface that became this event horizon before the gravitational collapse occurred.

But, even if a positive-energy black hole was to contain negative-energy matter, this matter would not remain in this situation for very long, because it would rapidly be expelled by a gravitationally repulsive force equivalent in strength to that which is attracting the rest of the matter toward the central singularity, so that the black hole would actually end up containing exclusively particles having the same sign of energy as its own. Thus, even if the sign of energy of the particles contained within any macroscopic surface was to constitute a relevant microscopic degree of freedom (transformed by the M symmetry operation) which, under particular circumstances, could contribute to the measure of information encoded on this surface, it is nevertheless appropriate to assume that positive energy black holes which have reached a stable state are composed exclusively of positive-energy particles. In fact, if I want to explain the results of the semi-classical theory of black-hole thermodynamics, I have no choice but to assume that the energy sign of every matter particle forming a black hole is the same as the energy sign of the object itself, because the conventional theory is based on the implicit hypothesis that positive-energy black holes exist in a stable state and are not in the process of releasing negative-energy particles, which means that they must be formed exclusively of positive-energy matter. I will, therefore, assume that a knowledge of the sign of mass of a macroscopic black hole allows to determine the energy signs of *all* the matter particles whose states are reflected in the detailed configuration of the microscopic degrees of freedom on the event horizon of the object.

It should be clear that under such conditions it cannot be assumed that the energy sign of particles, which is transformed by the action-sign-reversal symmetry operation M , constitutes the one binary degree of freedom per elementary unit of area which is associated with the measure of black-hole entropy provided by the semi-classical theory, because if that was the case, then, in the most common of situations, nearly all the microscopic physical

parameters of a black hole would be fixed by a knowledge of the sign of mass of the object and no information would be missing from the macroscopic description. It would thus follow that entropy would always be minimum, which is certainly not desirable, given that the semi-classical theory rather requires entropy to be maximum when matter collapses into a black hole. The constraint imposed by the sign of mass of a black hole on the energy sign of its constituent particles may not be so significant, however, given that even if we ignore any additional degree of freedom which could be transformed by the other discrete symmetry operations, a determination of the sign of energy cannot alone be considered to exhaust the requirements of a complete description of the state of the matter particles forming a black hole, because, in principle, energy must also be allowed to vary in magnitude.

For now, we may choose to leave aside that difficulty, but then we are still left with having to explain how it can be that the other two independent discrete symmetry operations, which should also characterize the states of matter under all conditions, provide at most only one single binary unit of information, even though they together transform two degrees of freedom. As I suggested above, those two symmetry operations may be chosen to be the action-sign-preserving time-reversal operation T and the action-sign-preserving space reversal operation P . You may recall that in the context of the redefinition of the discrete symmetry operations which I proposed in a previous section of this chapter, the T symmetry operation must be assumed to reverse all momenta, as well as all angular momenta and all non-gravitational charges, even if merely from the unidirectional-time viewpoint. The P operation, on the other hand, has absolutely no effect on the direction of angular momentum or the sign of charge, from any viewpoint, but must be considered to reverse the direction of momentum and the handedness of particles (as indicated in Table 3.3). Thus, taken together, the T and P symmetry operations would transform all the physical attributes of elementary particles which add up to produce the total momentum p , angular momentum j , and charge q parameters that characterize the macroscopic state of a black hole with a given sign of mass.

Yet this does not necessarily mean that all that must be specified in order to determine all of those macroscopic physical parameters are the signs of the microscopic parameters transformed by the T and P operations (reversed or not reversed for each of the two symmetry operations) for every elementary particle that forms a given black hole. There is, in effect, no *a priori* reason to assume that the momentum of elementary particles (like their energy) can

vary only in sign and it would rather seem that, not only must the magnitude of this parameter be allowed to vary like that of energy, but its orientation must also be allowed to vary, not in a binary way like the sign of energy, but as a continuous two-dimensional angular variable, which would forbid its complete determination through a knowledge of the value that would be taken by one single binary degree of freedom. What's more, even if, under ordinary circumstances, an action-sign-preserving reversal of time intervals generated by T would affect the direction of angular momentum (because it would reverse momentum independently from position), it would not affect the handedness of particles, and in the context where we are trying to identify the microscopic configuration associated with the states of elementary particles present on the Planck scale, it appears necessary to restrict our account of the spin state of matter particles to handedness. But while an action-sign-preserving reversal of space intervals obtained by applying P would actually reverse the handedness (because it would reverse the momentum of particles without also reversing their spin), we have no reason to assume that the spin could not itself reverse independently from momentum, thereby also reversing the handedness. It would then appear necessary to specify the handedness of particles independently from the other degrees of freedom which are reversed by those two symmetry operations. As a consequence, only the sign of charge of a given particle can be assumed to be entirely determined by its dependence on the redefined time-reversal symmetry operation T , when the effects of such a transformation are considered from the unidirectional-time viewpoint, which usually applies in a classical context.

Now, despite the fact that the T operation reverses both momentum and charge, it certainly seems appropriate to assume that, as far as those microscopic physical parameters are concerned, we are actually dealing with two distinct degrees of freedom, because momentum can also be independently reversed by the P operation. But even though it may appear obvious that the sign of charge should be independent from the direction of momentum, it is reassuring to observe that, from a bidirectional-time viewpoint, this hypothesis is unavoidable given that the variation of the sign of charge only occurs from a unidirectional viewpoint and is actually the consequence of a reversal of time intervals obtained while leaving the sign of action invariant, which would, in effect, reverse the sign of energy, but leave invariant the direction of momentum. In any case, the outcome of this reflection is that we have to accommodate three independent microscopic degrees of freedom which are the sign of charge (or equivalently the sign of time intervals),

which is reversed by T , the direction of momentum (or equivalently the sign of space intervals), which is reversed by P , and the handedness of elementary particles, which depends on the direction of spin and which can be reversed independently from charges and momenta. The state of all the relevant, macroscopic physical parameters of a black hole (except for its sign of mass or action) can then be derived from a knowledge of the microscopic state of those three independent parameters.

It is also important to mention that despite the fact that what I'm seeking to determine are the degrees of freedom which would apply on a very small scale, at which the fundamental interactions would presumably be unified, I'm nevertheless assuming that the sign of any non-gravitational charge would remain a parameter distinct from the sign of action (the sign of gravitational charge), because the variation of the sign of charge would here occur merely as a secondary consequence of a reversal of the direction of propagation in time, which must still be considered a significant change, clearly distinct from a reversal of action (which also involves a reversal of time, but which leaves the sign of energy unchanged), even under such conditions⁶.

If we recognize the appropriateness of those remarks, it would then seem that the situation may be even more problematic than I indicated above, because, despite the fact that we are considering as relevant only those parameters which are affected by the redefined discrete symmetry operations, the degree of freedom associated with handedness would provide an additional contribution (dependent on the direction of spin, which varies independently, like the sign of action) to the measure of missing information concerning the microscopic state of the matter that crossed the event horizon of a black hole. This contribution would add to those provided by the degrees of freedom associated with the sign of charge and the momentum of a particle (which are dependent on the sign of time intervals and the sign of space intervals, respectively). It would then seem that we still need at least three binary units of information per elementary particle to completely determine even just the

⁶I hope that the reader will forgive me for implicitly assuming that specifying the sign of one unified charge would allow to completely describe the state of the non-gravitational attributes of an elementary particle on the quantum gravitational scale, but it seems to me that this hypothesis is unavoidable given the current trend observed in particle physics, which clearly indicates that the distinctions we observe between the various charges characterizing matter particles in their low energy states are not fundamental, thereby implying that the particles themselves should not be distinguishable other than by the sign of their unified non-gravitational charge, on the energy scale at which the known forces unify.

signs of all the relevant physical parameters characterizing the microscopic state of matter under such conditions.

It is while I was trying to visualize what would happen to a negative-mass body which would find itself inside some surface that was about to become the event horizon of a positive-mass black hole that I realized that both positive- and negative-mass particles would actually be submitted to very restrictive constraints when experiencing the effects of the gravitational field which exists inside the region delimited by the event horizon of a black hole. Indeed, a negative-energy particle which would happen to be located near the center of a positive-mass black hole at the time of its formation would soon be repelled outward by a force as large as that it would experience inside the most powerful of particle accelerators. While it is being ejected outside the event horizon, the negative-energy particle would reach an arbitrarily high (negative) energy and its negative momentum would also become arbitrarily large in the direction of the forming central singularity (considered as the point where the density of the dominant form of energy reaches its maximum theoretical limit), regardless of what its initial state of motion was. The nearer to the center of the object the particle would initially be, the larger its final negative energy would be when it would emerge from the event horizon of the positive-mass black hole. But given that, in the case a non-rotating black hole at least, the force which accelerates the particle is always directed away from the forming singularity, it follows that the lateral components of its momentum would become completely negligible in comparison with the component of its negative momentum directed toward the singularity. Thus, if we were to consider only negative energy particles literally emerging from a positive-mass singularity, we would end up (for a non-rotating black hole, in the absence of collisions with infalling positive-energy matter) obtaining particles reaching the event horizon with a maximum (negative) energy and a negative momentum invariably directed opposite the positive normal to the surface of the black hole. In other words, we would always obtain (in the absence of interferences) particles in a very specific state of energy and momentum.

The process would be even more constraining for a positive-mass body in the gravitational field of a positive-mass black hole, given the rising tidal effect which, in this case, compresses the object laterally and stretches it vertically, as it is accelerated in the direction of the singularity. In such a case, we would necessarily end up with a very focused beam of particles whose

lateral motions would again be completely negligible (in the reference system in which the black hole is not rotating). Indeed, the force attracting the particles toward the singularity of the black hole would grow with time from the moment they cross the event horizon, eventually becoming so large that the energy of the particles would become as high as it can be, while the horizontal components of their momenta would become completely negligible in comparison with the vertical component of their positive momenta oriented toward the central singularity. Any residual, lateral motion would simply contribute to increase or decrease the total angular momentum of the black hole, whose rotation is shared by all the particles that fell into the singularity (as a consequence of collisions and relativistic frame dragging), and should not, therefore, contribute to entropy (as a measure of missing information concerning microscopic degrees of freedom). Thus, when a positive-energy particle reaches the singularity of a positive-mass black hole, its momentum (in the reference system relative to which the object is not rotating) has basically become a unidirectional variable. In fact, space itself must be considered to become analogous to unidirectional time for a positive-energy particle that crosses the event horizon of a positive-mass black hole, but what I came to understand is that this actually means that momentum (along with the sign of space intervals) would then become a fixed parameter, with a unique direction and a maximum magnitude. As a result, we once again obtain a unique final state of maximum energy and invariant momentum.

The crucial assumption in the present context is that a maximum energy must actually exist. I believe that this conjecture is appropriate, given that, from a quantum gravitational viewpoint, the existence of a minimum meaningful time interval or spatial distance implies the existence of a maximum limit for the magnitude of energy. Indeed, when we add energy to an elementary particle we can more accurately determine its position, which means that its energy becomes concentrated in a smaller region of space. Therefore, if the energy of a particle was to reach the point where it is so large, in a region of space so small, that it would form an elementary black hole (which occurs precisely when the energy of the particle is of the order of the Planck energy) it would follow that adding even more energy to it would simply amount to raise the mass of the elementary black hole in which it would be confined and to increase the area of its event horizon, which could only mean that a larger number of elementary particles are actually present, in a larger number of elementary units of area. Hence, there would be no sense in attributing one single elementary particle at the Planck scale an energy larger

than the Planck energy. The situation we encounter here is somewhat similar to that which we have in quantum chromodynamics, where, beyond a certain threshold, the energy spent at trying to separate two oppositely charged quarks in a meson no longer contributes to increase the distance-dependent attractive force between the two quarks, but merely ends up splitting the original particle into two new mesons, thereby neutralizing the force that existed between the original two quarks.

I would therefore suggest that we assume that the elementary particles that reach a singularity, after having been accelerated by its gravitational field, must be in a state of maximum energy, which we must recognize to be the Planck energy. Given that it is not that difficult to visualize what would happen to a positive-energy particle which would cross the event horizon of a positive-mass black hole, it is surprising that it had never been fully realized what the outcome of such a process would mean for a description of the final states of matter particles submitted to a gravitational collapse. But I believe that it is crucial to recognize, in order to clarify the whole question of black-hole entropy, that what happens when positive energy matter collapses into a black hole singularity is that it invariably reaches a state in which every elementary particle has a (positive) Planck energy and a correspondingly large momentum, characterized by a unique invariant direction (to which is associated a unique sign of space intervals) which is straight toward the singularity, regardless of the initial state of motion of the particles at the time they crossed the event horizon of the black hole.

Here it must be understood that, despite the fact that the wavelength of the light emitted by a positive-energy particle which is about to be absorbed by a positive-mass black hole would be infinitely redshifted (from the viewpoint of a remote observer not moving with respect to the event horizon of the object) and would show time as standing still, we are nevertheless allowed to assume that the events occurring after the particle crosses the event horizon of the black hole can be characterized in a physically meaningful way. It is certainly appropriate to consider, in particular, that a particle's momentum will keep increasing in the direction toward the singularity, as I'm suggesting would be the case, because time dilation does not mean that the particle itself would become motionless, but merely that the signals it emits are infinitely redshifted by the gravitational field of the black hole (still from the viewpoint of a remotely located observer). Thus, despite the fact that signals would show the particle as apparently immobilized on the event horizon, we must still assume that, from its own perspective, this particle

actually crosses the event horizon of the black hole and in a finite amount of time acquires an energy which, relative to a motionless outside observer, would be maximum.

Also, the idea put forward by certain authors that, from the viewpoint of an external observer, a positive-energy particle could in fact acquire a negative energy after it crosses the event horizon of a black hole, would only be appropriate if we were to consider that the negative gravitational potential energy reduces the energy of the particle itself into negative territory. But, in fact, this is not an appropriate approach to defining the energy of matter (especially in the context where the true properties of negative-energy matter are understood to make such a notion implausible), because as far as this potential energy is concerned we are actually dealing with a distinct contribution to the total measure of energy, which is that of the gravitational field. The truth is that the kinetic energy of the particle itself would keep increasing to arbitrarily large values, even if this energy is compensated by a growing negative contribution to the energy of the gravitational field associated with the interaction of this particle with the rest of the matter in the black hole.

However, one may perhaps question the conclusion that momentum would have a fixed magnitude for any positive-energy particle that reaches the singularity of a positive-mass black hole, in the context where the rest mass itself may be a variable parameter. It is true, in effect, that the magnitude of this momentum would depend on the rest mass of the particle which is accelerated in the gravitational field of the black hole, given that all masses have the same acceleration and are therefore subjected to the same velocity increase. But in the context where we are dealing with final kinetic energies which are so large, it appears appropriate to assume that the energy associated with the rest mass of the particles which are reaching a singularity, after having been absorbed by a black hole, would be negligible or null⁷, so that, if the total energy of those particles is the Planck energy (the maximum physically meaningful measure of energy), we are allowed to conclude that the final magnitude of their momentum would always be what we may call a Planck momentum, understood as the maximum theoretically meaningful value of momentum which can be carried by a massless particle (and

⁷In fact, even if that was not the case, from the viewpoint of a unified theory of interactions it should probably be the case that only one magnitude of mass would exist for the most elementary particles which compose the currently known matter and radiation particles and this would have consequences similar to those discussed here.

associated quantum mechanically with the smallest meaningful measure of spatial distance). Under such conditions, we would have no choice but to recognize that the final magnitude of the momentum of all particles reaching a black-hole singularity actually constitutes an invariant property, just like the direction of this momentum.

Now, I initially thought that it would be appropriate to assume that if space actually comes to an end for matter that reaches a singularity, then momentum, as the conjugate attribute to space, simply cannot continue to evolve after the final stages of a gravitational collapse and would remain in the final state I have identified above for the whole lifetime of the black hole. But some relatively recent results from loop quantum gravity appear to show that the final state of a gravitational collapse is not a singularity, but merely a state of maximum matter density, which would immediately be submitted to a ‘quantum bounce’ that would turn the collapse into a process of outward expansion. It is sometimes argued that this might be problematic, given that, if a black hole was to expel matter, it seems that entropy could decrease in the process. However, given that black-hole evaporation does involve a *local* decrease of entropy for the black hole itself (independently from its environment) over its entire lifetime, then the prediction that the singularity would decay may not be as paradoxical as one might assume. In fact, I think that if black holes do evaporate, then something like the quantum bounce must occur, so that there remains no singularity in the final state, when the mass of the black hole itself has become minimal. The perceived problem merely arises when we fail to recognize that the near infinite time dilation that is attributable to the enormous gravitational field of a black hole implies that the process of gravitational collapse and the following quantum bounce that would take place over a finite and relatively short time from the viewpoint of the matter that falls toward the singularity, actually appear to occur over the entire lifetime of the object from the viewpoint of an external observer⁸.

Thus, from the viewpoint of an external observer, the whole collapse process, as well as the quantum bounce that follows it, would take place over the arbitrarily long period of time during which the black hole would exist. The quantum bounce, if it could be observed from outside the event horizon,

⁸After I first wrote those lines I realized that Carlo Rovelli, one of the original contributors to the theory of loop quantum gravity, has more recently arrived at the same conclusion.

would, therefore, appear as a very slow process during which all the matter particles that ever fell toward the singularity would reverse their collapsing motion and begin expand outward, eventually reaching the horizon at which point they would be released in the form of macroscopically thermal radiation, as the black hole slowly evaporates. Given that the energy originally contained in the objects which were absorbed by the black hole can only be released in high-entropy form as the black hole decays through the emission of such macroscopically thermal radiation, it follows that no violation of the second law of thermodynamics would be observed.

Despite what is sometimes suggested, therefore, the process that takes place following a generic quantum bounce is different from a white hole (conceived as the time-reverse of a black-hole gravitational collapse), because the matter which is released following the bounce has high entropy and does not consist in the same macroscopic objects that originally fell through the event horizon. Yet it must also be the case that it is as a result of the quantum bounce that the information concerning the microscopic state of the particles which were captured by the gravitational field of a black hole can be released in higher entropy form as Hawking radiation. It may well only be the widespread ignorance of the unavoidable character of this conclusion that prevents us from acknowledging the fact that this information is not really lost as a result of the evaporation process, but is actually contained in the detailed microscopic state of the emitted radiation.

It is simply the fact that the prediction that a black hole must emit radiation was originally derived from a semi-classical theory and is, therefore, dependent on the hypothesis that the gravitational field is microscopically uniform and itself devoid of small-scale structure that explains that Hawking's original approach cannot account for the presence of information in the flow of energy that emerges from a black hole as it decays⁹. But it must be clear that, for an observer outside a black hole, the information about the matter that was captured by its gravitational field does not vanish from reality, but actually remains encoded in the microscopic degrees of freedom on

⁹What is happening according to the most knowledgeable experts [27] is that quantum correlations with the discrete, high-energy degrees of freedom present in the radiation remain unobservable on a macroscopic scale and could only be obtained by performing observations on the Planck scale (at very high energy) and this is why the radiation appears to remain thermal from a semi-classical perspective, despite the fact that it does contain the information that was encoded in the microscopic degrees of freedom on the event horizon of the decaying black hole.

the event horizon at all times, because, due to time dilation and space contraction (in the radial direction), the matter particles would appear to spend all their time on the event horizon itself, right until the moment when their energy is released as thermal radiation, after a period equivalent to the entire lifetime of the black hole (from the viewpoint of an observer outside a black hole, everything that happens to the particles as they reach the singularity would seem to take place on the event horizon itself).

As a result, any entanglement of a particle outside a black hole with a particle crossing its event horizon remains in effect right until the time when the energy of this very same particle is released, after an arbitrary long time corresponding to the lifetime of the object. Therefore, it is not appropriate or necessary to assume that the outside particle becomes entangled with the whole thermal radiation or that this entanglement is lost along with the information about the state of the particle that fell toward the singularity. But all current difficulties we encounter while trying to account for the conservation of information in the presence of black holes arise from assuming one or another of those two incorrect assumptions, which means that the only real problems are actually the result of a misunderstanding. Indeed, from the viewpoint of loop quantum gravity it can even be conveniently explained how it is, exactly, that this information is encoded, in the quantum gravitational degrees of freedom on the surface of a black hole and therefore, to the extent that the currently favored quantum theory of gravitation can be assumed to provide an accurate representation of what is happening on the Planck scale, we may consider that the black-hole information-loss ‘paradox’ has been solved already. But it must be clear that there never was any paradox, because, as I mentioned above, we do know that classical general relativity and the hypothesis that the spacetime structure is smooth and uniform down to the smallest scale is not valid in a quantum-mechanical context, while it is also clear that a certain measure of information must be associated with any microscopic structure that exists on such a scale.

In any case, if we are willing to recognize that the description provided by current quantum theories of gravitation constitutes the most accurate account of the process of black-hole gravitational collapse that one can presently derive, it would follow from the preceding analysis that over the entire lifetime of a black hole the particles with the same energy sign as that of the object would spend most of their time either collapsing, with maximum momenta directed toward the singularity, or expanding, with maximum momenta directed in the exact opposite direction (as would occur after the quantum

bounce takes place). This is because, the time dilation effect is maximum when the particles are near the singularity and are either collapsing with maximum energy or expanding with maximum energy, so that from an external viewpoint they would, in effect, appear to spend most of their time in either one of those two states. Those discrete states would therefore be the ones that need to be reflected in the configuration of the microscopic degrees of freedom on the event horizon of the black hole, before the object actually evaporates to nothing. More specifically, the detailed configuration of those microscopic degrees of freedom must be considered to reflect the state of motion of every particle contained in the black hole at the time immediately before or immediately after they reach the singularity of the object. It must be clear, however, that it cannot be assumed that all the matter that ever crosses the event horizon of a black hole reaches the maximum-density state at the same time and reverses its direction of motion from collapsing to expanding all at once. The crucial point, therefore, is that it is not possible for an observer outside the black hole to tell in which of those two different states any given particle is at any given time, as the precise time at which a given particle reverses its motion is not known and from the viewpoint of such an observer this could happen at any moment during the whole lifetime of the object.

If we agree on the plausibility of the above conclusions concerning the state of any particle involved in a gravitational collapse, then we need to recognize that it is not only the sign of energy of the constituent particles that would be constrained by the macroscopic properties of a stable-state black hole. Indeed, it now appears that not only must the signs of energy of the particles in the final stages of a gravitational collapse be considered to be completely determined by a knowledge of the sign of mass of the object, but the magnitude of those energies is also to be considered an invariant parameter, which therefore cannot contribute to the entropy of the black hole. What's more, it would appear that the momenta of the particles which are constrained by the presence of an event horizon can only point in one of two directions and only vary as binary degrees of freedom (along with the sign of space intervals which is transformed by the space reversal operation P). It should be clear, therefore, that while the energy state of the particles which were captured by the gravitational field of a black hole cannot contribute to the measure of entropy or missing information which determines the temperature of the object (once its own sign of energy is known), the momentum state of each of those particles does contribute one binary unit to the mea-

sure of missing information encoded in the quantum gravitational degrees of freedom on the surface of a black hole.

I'm therefore allowed to conclude that three binary degrees of freedom are allowed to freely vary for every elementary particle in the final stages of a gravitational collapse. The first is the direction of momentum (or the sign of space intervals) and the other two are the ones I previously identified as being the sign of charge (or the direction of time intervals) and the handedness of matter particles. All three attributes could potentially contribute to the measure of missing information concerning the microscopic state of a black hole. Once the momentum direction of a given particle is known it appears, in effect, that the handedness of the particle would still be allowed to vary and should, therefore, contribute one additional binary degree of freedom, which would vary upon a reversal of the direction of spin relative to this momentum direction. But it also appears that there should be a contribution to the measure of entropy derived from the semi-classical theory of black-hole thermodynamics by the sign of charge of the particle, that varies upon a reversal of its direction of propagation in time. Even if we assume that there exists only one type of charge, with one given polarity, for the elementary particles present at the unification scale, certainly information should be needed to specify whether this charge is positive or negative from the viewpoint of unidirectional time.

I have explained why three binary units of information would be enough to account for all the fundamental degrees of freedom of any positive-energy particle present in the final stages of the gravitational collapse of a positive-energy black hole. But what appears to be required by the semi-classical theory of black-hole thermodynamics is that every elementary unit of area encodes only one bit of information, which is problematic because, for elementary black holes at least, all the information that exists concerning a given particle must be obtained from only one elementary unit of area. Given that what I'm seeking to achieve is a complete determination of all the physical properties of the particles present within the event horizon of a black hole from a knowledge of the value of all the relevant discrete degrees of freedom, it would seem that I have fell short of this objective. I believe, however, that, in fact, the problem we seem to have encountered is not real.

The truth is that there is no contradiction between my account of the quantity of information required to completely describe the state of an elementary particle which was captured by the gravitational field of a black hole and the measure of missing information encoded in the microscopic state of

the gravitational field on the event horizon of such an object. To understand what motivates this conclusion, we must first acknowledge that the formula for black-hole entropy was derived from arguments related merely to the thermodynamic properties of the gravitational field and only in the context where the source of this field is the conventional stress-energy tensor and the measure of missing information involved allows to set the relationship between the surface gravitational field and the temperature of the thermal radiation emitted by a black hole (even though it is usually assumed that the area of the event horizon of a black hole is also determined in part by the electric charge of the object). But, if it is indeed the case that only one out of three bits of information, concerning the state of matter particles contained within the event horizon of a black hole, is encoded in the detailed configuration of the gravitational field on its surface, I think that this is because there is more information encoded in some other physical properties of black holes that do not contribute to the measure of temperature associated merely with their conventional surface gravitational fields.

Once we have recognized that there must be more information, concerning the microscopic state of matter contained within a surface, than is provided by the detailed configuration of the gravitational field on that surface, what becomes crucial to understand is that there is no reason to assume that the gravitational field should provide information about the microscopic distribution of some non-gravitational charge, because that information must actually be contained in the detailed microscopic configuration of the field of interaction associated with this particular charge. It is surprising in fact that this requirement was never considered before, because when one carefully thinks about this question, it is hard to arrive at a different conclusion. If the information that is potentially missing about energy and momentum (as the physical properties of particles which constitute the source of gravitational fields) is to be associated with the microscopic state of the gravitational field, then it is also quite unavoidable that the missing information about, say, the electric charge is to be associated with similar microscopic aspects of the electromagnetic field. There is, in fact, absolutely no reason to assume that the detailed configuration of the electric charges which are the source of the electromagnetic field should be determined from information contained in a different force field, which would here be the gravitational field. It must be clear, however, that any information, associated with the electromagnetic field on the surface of a black hole, that would encode the details of the configuration of electric charges inside the object, would have to be contained

in the *microscopic* (Planck scale) degrees of freedom of the electromagnetic field and would not be reflected in the classical, macroscopic parameters of this force field, which means that the hypothesis that black holes have no ‘hair’ would still be valid.

Now, if the missing information concerning the microscopic distribution of electric charges or electric charge signs inside a given surface (whether or not this surface is that of a black hole) can only be encoded in the detailed configuration of the electric field on the boundary delimited by that surface (even when the total charge inside the surface would be null), rather than in the configuration of the gravitational field on the same boundary, then it means that a theory that would seek to derive a measure of the amount of information necessary to determine the state of the matter contained inside this surface based only on features of the gravitational field present on the surface (which in the case of black holes would be the event horizon) would necessarily fall short of providing the accurate value. Therefore, the results derived from the semi-classical theory of black-hole thermodynamics, concerning the relationship between the entropy of a black hole and the area of its event horizon, would not rule out the existence of an additional amount of missing information associated with the exact microscopic state of the matter trapped within such a surface.

I believe, in effect, that the missing information concerning the sign of charge of every particle forming a black hole, which is transformed by the T symmetry operation (from a unidirectional-time viewpoint), cannot be reflected in the microscopic state of the gravitational field on the surface of the object and this explains that it need not be taken into consideration when deriving the statistical mechanical properties of black holes associated with their surface gravitational fields¹⁰. This is why we were allowed to ignore the existence of this information when elaborating the semi-classical theory of black-hole thermodynamics from which the conventional measure of black-hole entropy was derived. It thus appears that one additional binary unit of information (distinct from those which are associated with the momentum direction and the handedness of matter particles) is indeed missing concerning

¹⁰This conclusion is especially appropriate in the context where one recognizes that, from the bidirectional viewpoint of time, it is really the sign of energy that reverses under application of a T symmetry operation, while this reversal is combined with a reversal of physical time intervals, which means that it is not significant from a gravitational viewpoint and therefore should not be reflected in the microscopic properties of the gravitational field on the event horizon of a black hole.

the state of every elementary particle in the final stages of a gravitational collapse. This information would allow to determine the sign of charge of each and every particle which contributes to fix the total charge q of a black hole, or more specifically the direction of time intervals along which those particles are propagating and from which depend the sign of their charges from a unidirectional-time viewpoint. We are then allowed to assume that this is the binary unit of information which is actually associated with the T symmetry operation (or alternatively the C symmetry operation) defined in a previous section.

It would, therefore, seem that there is, in effect, more information associated with the microscopic state of the matter contained in a black hole than is encoded in the detailed configuration of the discrete degrees of freedom of the gravitational field on the event horizon of the object. But I have explained why we should not expect the missing information about the distribution of non-gravitational charges, at least, to contribute to the conventionally derived measure of black-hole entropy. Instead, this additional missing information should be encoded in the microscopic state of the interaction fields associated with those non-gravitational charges, which would give rise to its own independent contribution to the temperature of a black hole. In this context it is important to note that there actually exists an analogue to the Hawking radiation process associated with the gravitational field of black holes and which involves the electromagnetic field. It is a known fact, indeed, that, past a certain magnitude, the electrostatic field surrounding a charged nucleus would induce processes of real particle-antiparticle pair creation similar to those responsible for the emission of radiation by the gravitational field of a black hole and I believe that this phenomenon would allow a similar treatment of the thermodynamic properties which, according to the above proposal, should be associated with any distribution of non-gravitational charge. Only, in the case of non-gravitational charge we are usually dealing with situations where the total charge is indeed null, even when large amounts of positive and negative charges are present inside a surface. Such situations are therefore more analogous to that which is occurring when the measure of gravitational entropy is constrained merely by the Bekenstein bound and both positive- and negative-energy matter are present together inside an ordinary surface.

This is all well, but then one may ask: what about the contribution by the handedness of particles to the measure of information encoded on the surface of a black hole? Shouldn't that information be contained in the microscopic

state of the gravitational field? Given that angular momentum couples to the gravitational field, one may be tempted to assume that this must be the case. But under such conditions it seems that there should be twice as much missing information contained in the gravitational field on the surface of a black hole than what the semi-classical theory of black hole thermodynamics requires, because the information about the handedness or the spin direction of elementary particles cannot be contained in the microscopic state of the electromagnetic field or that of any other non-metric force field. So, how can one account for the information that is missing from the macroscopic state of a black hole concerning the handedness of every elementary particle it contains? Faced with this difficulty, I once thought that what one would have to conclude is that the momentum direction of particles must, in fact, be fixed, at all times, to the value it had in the very last stages of a gravitational collapse, just before a particle reaches the spacetime singularity, so that only the handedness of particles would contribute to the measure of missing information contained in the gravitational field. But when I came to recognize the inadequacy of this hypothesis, in the context where a gravitational collapse must be understood to give rise to a quantum bounce, I was forced to admit that the mistake I had made was to ignore the fact (of which I was already aware) that, while momentum and energy are the source of spacetime curvature, spin appears to be the source of a torsion of spacetime.

Indeed, a generalization of relativity theory is known to exist according to which spacetime may not only possess a curvature, described by the Einstein tensor $G_{\mu\nu}$ (that must be proportional to the stress-energy tensor of matter $T_{\mu\nu}$), but may also be subjected to a torsion, which can be described by the Cartan tensor $C_{\mu\nu}^\gamma$ and which would be proportional to the spin angular momentum tensor of matter $M_{\mu\nu}^\gamma$ (through the exact same proportionality constant $\chi = -8\pi G/c^4$). When those last two tensors are not null (contrarily to what Einstein originally assumed), an attractive force $F_T = \pi G/c^2 \nabla(\sigma_{jk}\sigma^{ik})$ exists that depends on the spin density of matter σ_{ij} and that tends to concentrate spin. But, in the present context, it only makes sense that this force would be repulsive for particles with opposite spins or handedness, because particles with opposite spins carry opposite measures of action and given the gravitational nature of this interaction, it can be expected that there is an analogy with the situation we have when negative-energy matter is present, except that in the present case particles with opposite spins and opposite measures of action (associated with those angular momenta) can become trapped in the gravitational field of the same macroscopic, stable-state black

hole, as long as they have the same sign of energy (propagated in the same direction of time). As a result, it becomes possible for information to be contained in the microscopic (quantum gravitational) configuration of the metric field associated with the torsion of spacetime about the microscopic state of handedness of each and every particle that crossed the event horizon of a black hole, which is fixed by the spin direction of those particles. Given that the semi-classical theory of black hole thermodynamics is based on the implicit assumption that spacetime torsion does not contribute any measure of missing information, then it is possible to understand why the theory does not account for the entropy which must, according to my analysis, be associated with the handedness of particles.

If you have understood the essence of my argument, then there should be no doubt that the only missing information which is actually encoded in the microscopic configuration of the degrees of freedom of the gravitational field on the event horizon of a positive-energy black hole associated with spacetime curvature is that which allows to determine the randomly variable momentum direction of every positive-energy particle it contains, using merely one single bit of information for every elementary particle. This conclusion should perhaps have been expected, given that energy and momentum are the only attributes of elementary particles which are the source of the curvature of spacetime, which determines the surface gravitational field of a black hole, while momentum is the only one of those two attributes whose polarity is not fixed by a knowledge of the sign of mass of the object. In any case, if we are willing to accept the validity of the arguments on which this deduction is based, it would then follow that we now have an explanation, not only for the fact that the states of matter particles which are trapped by the gravitational field of a black hole vary as discrete variables, but also for why it is that only one such variable (instead of three or four) actually contributes to the measure of missing information which must be taken into account in determining the thermodynamic properties of such an object associated with its surface gravitational field. As a result, the measure of information associated with the matter content of an elementary black hole is allowed to match the value of entropy derived from the semi-classical theory of black-hole thermodynamics, which requires each elementary unit of surface (equal to four Planck areas) to encode one binary unit of information.

Therefore, it is now actually possible to at least confirm the existence of a definite relationship between the microscopic state of the quantized gravitational field on the surface of a black hole and actual states of the matter it

contains. What held the key to a better understanding of the exact nature of the degrees of freedom characteristic of the states of matter submitted to a gravitational collapse was the recognition that, for matter particles reaching a black-hole singularity, the only relevant variables are the signs of all those physical parameters which are transformed by the previously discussed discrete symmetry operations. It is remarkable that the direction of momentum should be one of the only fundamental parameters of elementary particles (along with the handedness and the sign of charge) that is not constrained to any specific value by the conditions prevailing in the final stages of collapse into a spacetime singularity (or in the moments immediately preceding or following a quantum bounce) and that it must, therefore, alone, contribute to the measure of entropy associated with the gravitational field of a black hole. This is certainly the most significant outcome which has emerged from my re-examination of the question of discrete symmetries as it arises in a semi-classical context.

If we now return to the more general case, for which the density of matter is not large enough to produce an event horizon and the possibility for positive- and negative-action matter to be present together inside a surface cannot be ignored, it transpires that this is a situation in which more information would be required to describe the microscopic configuration of matter, because more states of motion are allowed for the particles in the period before such a configuration reaches a stable state. Indeed, even when an event horizon associated with a positive-mass black hole is present, it is clear that while a positive-energy particle would be drawn toward the center of mass of the object during the collapsing phase, a negative-energy particle which would be present in the same place at the same moment, would be repelled outward by a force of similar magnitude (to the extent that the average cosmic density of positive-energy matter can be neglected). Thus, in such a case, we would need to take into account at least one additional binary degree of freedom, associated with the sign of energy of the matter particles present inside the surface, which would also determine the directions of the space intervals associated with their states of motion.

But this would actually be the simplest case, as more complex states of motion would be allowed if the matter was not contained within a surface that constitutes a black-hole event horizon, because under such conditions not only would the orientations of the momentum of particles be allowed to vary, but it seems that their magnitudes could also vary significantly. It is

important to understand, however, that the validity of the Bekenstein bound would be preserved even if more information was required to determine the exact microscopic state of matter under those less constraining conditions. This is, again, because while more information may be required to describe the state of matter when the energy magnitudes and the momentum orientations are not fixed, this information growth would be offset by the decrease in gravitational entropy (the amount of missing information required to describe the unknown microscopic state of the gravitational field itself) that would result from the lower (nearer to zero) positive and negative densities of matter energy associated with such configurations, or from a mixture of matter of both energy signs (I will explain in section 4.7 why it is, exactly, that a local diminution in the magnitude of matter energy density is associated with a lower measure of missing information concerning the microscopic state of the gravitational field).

Now, it may appear contradictory that under ordinary circumstances, when no macroscopic event horizon is present and the distribution of matter energy is smoother, it is more difficult to tell the exact states of the particles present within a surface. How could it be more difficult, in effect, to determine the microscopic state of the matter inside an ordinary surface, when it seems that you can more easily observe what is going on inside such a surface? But, in fact, all I have argued so far is that there is information encoded on the surface of a black hole about the state of the matter that was captured by the gravitational field of the object, not that this information can actually be obtained by any observer under all circumstances. First of all, it must be clear that from the viewpoint of an observer standing outside a black hole, away from the event horizon, the only information that is readily available about the object is contained in the value of its macroscopic parameters of mass, momentum, angular momentum, and charge. Due to the microscopic nature of the event-horizon degrees of freedom which encode the information about the state of the matter that is trapped inside a black hole, the only way this information could be obtained is by performing very precise measurements right on the surface of the object. But while the state of the various particles contained inside a given surface must always be reflected in the exact microscopic state of the quantum gravitational degrees of freedom associated with that surface, when the surface in question is the event horizon of a black hole, additional difficulties arise that would limit the capacity of a remote observer to obtain knowledge about this exact state.

Indeed, as I explained above, from the viewpoint of an observer outside

a black hole, what happens to the particles that reach its singularity would appear to take place right on the surface of the object, due to time dilation and the contraction of distances in the direction of the singularity. Thus, in principle, information about the state of any elementary particle could be obtained by directly measuring the state of the relevant quantum gravitational degrees of freedom on the surface of the black hole in which they fell. The problem, however, is that doing so would require you to approach the surface of the object with the appropriate measuring device to the point where your distance from it would be no larger than the scale of quantum gravitational phenomena. You may then be able to perform the required measurements, but given the enormous difference between the gravitational potential on the surface of the object and that just above it, if you tried to send back a signal encoding the information you have been able to obtain, this signal could not be received by a remote observer before the black hole itself evaporates through the emission of Hawking radiation, at which point the information would actually be contained in the radiation itself and could be determined by examining the quantum gravitational degrees of freedom on the ordinary surface enclosing it.

What happens, therefore, is that due to time dilation there is an *absolute* limitation that prevents any observer outside a black hole from obtaining the information that does exist, right on the event horizon, about the state of the matter inside the object and it is the unavoidable character of this practical limitation that allows one to consider that the information missing from a description of the state of a black hole in terms of its macroscopic parameters of mass or surface area must remain unknown, as a matter of principle, despite the fact that it does exist. As a result, we are justified to assume that those macroscopic parameters provide a natural definition of coarse-graining that does not exist in the case of a general surface, whose information content is limited merely by the Bekenstein bound. Even though the information about the state of all the particles that crossed the event horizon of a black hole is encoded in the microscopic degrees of freedom on the event horizon of such an object, from a practical viewpoint this information cannot be obtained before it is made irrelevant. With the appropriate experimental means, the microscopic state of the quantum gravitational degrees of freedom on a surface can be determined down to the most intricate details, but when this surface is the event horizon of a black hole, this information remains unknown to the outside world and the related measure of entropy becomes objectively defined, as there is no alternative choice of coarse-graining that

could provide a more accurate description of the state of the object.

Thus, even though, on the basis of the arguments I provided above, it must be assumed that the information about the state of all the particles that crossed the event horizon of a black hole is encoded in the microscopic quantum gravitational degrees of freedom on the event horizon of such an object, from a practical viewpoint this information cannot be obtained before the black hole releases it in the form of macroscopically thermal radiation, which means that a non-subjective measure of entropy exists, to which can be associated a non-subjective concept of irreversibility. This is a very significant constraint, because it is ultimately the non-subjective character of that portion of entropy variation which is attributable to the gravitational field that enables one to conceive of the irreversibility that characterizes the evolution of certain macroscopic physical systems as being an objective property, even under conditions where gravitation does not appear to be involved, given that, as I will emphasize in section 4.8, all the entropy growth that is taking place in our universe must ultimately be attributed to the initial conditions of low gravitational entropy that existed in the remote past. The real difficulty here consists in recognizing that it is possible for information to remain absolutely unknown, concerning what takes place on the event horizon of a black hole, despite the fact that this information does not to vanish from reality. What allowed me to realize that this conclusion is not self-contradictory is my profound conviction that information must indeed be conserved under all circumstances, even when we do not have direct knowledge of the reality it describes.

Of course, this conclusion would also apply, from the viewpoint of a negative-energy observer, for the state of all those negative-energy particles which are under the influence of a negative-energy black hole. No information about this negative-energy matter can be communicated by a negative-energy observer located just outside the event horizon of the object (using negative-energy photons) to an observer located farther away, so that this information must be considered to be missing from the viewpoint of negative-energy observers. Now, one may ask whether it would be possible for a positive-energy observer to obtain more information about the microscopic state of negative-energy matter under the same circumstances? But that does not appear possible, because, due to gravitational repulsion, a positive-energy observer would never be able to reach the surface of a negative-energy black hole, on which the information about the microscopic state of the negative-energy matter inside the object would be located. Indeed, an observer must ap-

proach the event horizon of a black hole to within a quantum gravitational unit of distance in order to determine the state of the microscopic degrees of freedom on its event horizon, which means that, even if it would be possible for a positive-energy observer to communicate information to observers outside a negative-energy black hole (using positive-energy photons), the fact that such an observer cannot get to within a quantum gravitational unit of distance of the object means that this information cannot even be obtained, so that it must remain unknown, as a matter of principle, for positive-energy observers as well, even though it does exist. A certain measure of entropy must therefore be associated with such a black hole and indeed with any overdensity in the negative-energy matter distribution, even from the viewpoint of a positive-energy observer.

The validity of the idea that, in the case of a normal surface at least, it is always possible to obtain detailed information about the exact microscopic state of the elementary particles it contains, by examining the microscopic quantum gravitational degrees of freedom on that surface, can perhaps only be appreciated when it is recognized that the classical gravitational field, as it is usually described in a general-relativistic context, does not provide a complete account of the degrees of freedom present in the curvature of spacetime on a microscopic scale, which actually depends on the small-scale distribution of matter and radiation energy (more arguments will be provided in support of this conclusion in section 4.7). But if there are local variations in the curvature of spacetime, above those described by the smooth macroscopic configuration of the gravitational field, then it is only natural to expect that if some property of the gravitational field was to be measured in a very precise location, this usually unobserved substructure would become apparent and the information associated with it would no longer constitute missing information¹¹.

In any case, what's most significant regarding those situations where the entropy associated with the gravitational field is not maximum is that we are necessarily dealing with transitional states which will, in general, continue to evolve until the configuration described in the preceding paragraphs is

¹¹In the concluding section of chapter 5 I will explain that it can also be expected that there exist *unobservable* random fluctuations in the classical gravitational field, but it must be clear that, due to their fundamentally unobservable and random nature, such fluctuations would differ from the small-scale variations discussed here, which arise from the presence of microscopic inhomogeneities in the distribution of matter and radiation energy that are in principle observable, even though they are usually ignored.

reached. Thus, the negative-energy matter which may be present inside a surface containing mostly positive-energy matter will eventually be expelled from that surface, while the positive-energy matter will keep collapsing on itself until it forms a black hole. When all the negative-energy matter is released from a surface containing a larger proportion of positive-energy matter, the total mass contained within the surface actually *increases* and this means that its gravitational entropy grows larger in the process. We are therefore in a situation where a surface containing less matter (but not less mass) can have a larger entropy. This counter-intuitive outcome is allowed because when negative-energy matter is released outside such a surface, the total amount of information required to describe both the microscopic state of the matter particles still contained within the surface and their associated gravitational field grows larger. A negative-energy particle inside a surface containing more positive-energy matter does contribute (positively) to the amount of missing information concerning the microscopic state of the matter within the surface, but at the same time it reduces the amount of information attributable to the gravitational field, which happens to be larger than that attributable to the matter, so that, overall, the amount of information which must be encoded in the microscopic state of the gravitational field on the surface is smaller than it would be without the presence of the negative-energy particle.

The more general situation, where only the Bekenstein bound may apply, is therefore not incompatible with the results I have derived from a study of stable-state black holes, from which all matter with an energy sign opposite that of the object has been expelled. In fact, it seems that, from a fundamental viewpoint, there is no real difference between the situation we observe in general, when opposite-energy particles are necessarily allowed to be present within a surface, and that which arises when we are considering the surface delimited by the event horizon of a black hole. Yet, the fact that the presence of negative-energy matter within a positive-energy black hole would only be temporary (even from the viewpoint of an external observer, given that negative-energy matter does not experience the metric properties of space and time shared by positive-energy observers) and would always give way to a more stable state in which only positive-energy matter would remain inside the surface delimited by the event horizon of the object, appears to suggest that such end states play a role in gravitational physics which is analogous to that which is played by thermal equilibrium states in statistical mechanics. But the real question, regarding the Bekenstein bound, is how it

can be that, under the more general conditions in which it applies, the energy and the momentum states of matter particles located within a surface are allowed to vary in a continuous way, not just in magnitude, but (in the case of momentum) also in spatial orientation, while the measure of information encoded on the surface must still be provided in binary form.

What my investigations have led me to understand is that, in fact, this freedom is only apparent. It turns out that even under the more general circumstances discussed here, the magnitudes of the energies and the orientations of the momenta of elementary particles are restricted to vary as binary parameters. What allows me to draw such a bold conclusion is that I have recognized the consequences of the fact that event horizons are actually always present on the shortest distances, where quantum fluctuations in the energy of the gravitational field continuously give rise to the formation of ephemeral Planck-mass black holes. It is clear that the fluctuations in energy occurring on the Planck scale do not, all by themselves, imply that the energy of particles must be fixed to some maximum value, but the fact that such fluctuations are omnipresent when we reach this scale means that elementary black holes are actually the substance of physical space and time at this level of precision of measurement and if that is the case, then it means that matter is always shrouded in the event horizons of those microscopic black holes and therefore we can only conclude that, locally, it is submitted to the same constraints that would apply in the presence of a macroscopic black hole.

Thus, the energies that could be measured locally would always be of the order of the Planck energy, because the particles trapped within those microscopic black holes would be accelerated to arbitrarily high energies by the gravitational fields present on their surfaces. Indeed, the surface gravitational fields produced by black holes with such small masses would be extremely large, therefore compensating for the short time intervals during which they would actually be allowed to accelerate the particles which are submitted to their influence. It must be clear, however, that there can still occur variations of energy in units smaller than the Planck energy on larger scales, where only average values of the energy of matter and its associated gravitational field are significant and most contributions can be expected to cancel out. The Planck energy must not, therefore, be conceived as a minimum unit of energy (in a more general context), because, to the contrary, it constitutes a maximum level of energy, which must nevertheless be the only possible measure of energy *magnitude* concerning the states of matter

at the maximum level of precision of spatial distances and time intervals set by current quantum gravitational theories.

The case of momentum orientation is a little more complex, because we are dealing here with a scale at which quantum indefiniteness in position cannot be ignored. This is reflected in the fact that the same elementary unit of surface would actually correspond to every possible orientation normal to the surface of an elementary black hole. But even if it is not possible to associate a classically well-defined orientation to the momentum of a particle submitted to the gravitational field of such a microscopic black hole, it remains that, quantum mechanically, there would exist a definite (but superposed) state of momentum orientation, even for particles in such a situation, and this state would still be constrained by the configuration of the local gravitational field. In other words, there would still be a constraint on momentum orientation to be fixed by the presence of the gravitational field. Thus, I believe that when we are considering the states of particles on the scale of an elementary unit of volume, corresponding to an elementary unit of area (equal to four Planck areas), momentum orientation would still be a fixed parameter, so that a particle's state of motion would only be allowed to vary in a discrete way (given that the sign of space intervals could be either positive or negative), even when the particle is not under the influence of a *macroscopic* event horizon.

Indeed, as a result of quantum indeterminacy, it is impossible to specify the orientation of the local momenta any more precisely than there are elementary units of surface associated with the microscopic black hole in which a particle is trapped. So, each elementary unit of area on the event horizon of a local microscopic black hole still contains the same amount of information as would an elementary unit of area associated with a macroscopic black hole. This is true even if it would be possible to define the orientation of the elementary units of microscopic black-hole surface in a very large number of ways, because the orientation of the momentum of a particle under its influence cannot be determined any more precisely than by selecting an elementary unit of area on the surface of the object. The orientation of the elementary surface elements of the microscopic black hole could vary in a near continuous way, but given that the momenta of the particles constrained by the event horizon of this black hole are in a state of quantum superposition, then their directions cannot be identified any more precisely than by specifying the value of a discrete degree of freedom associated with a particular one of the surface elements, regardless of the exact orientation

of those units of area. Thus, on a local scale, there would be a finite number of possibilities (associated with the finite number of surface elements on a microscopic black-hole event horizon) for the momentum orientation of a particle, which can therefore be specified exactly (relative to that of other particles constrained by the same event horizon) using a minimum number of binary units of information.

Now, given that there appears to exist a precise correspondence between the state of a matter particle reaching a black-hole singularity (conceived as being merely a maximum-density state with finite volume) and a given elementary unit of surface on the event horizon of the object, then, in the context where the Bekenstein bound is assumed to apply, it would seem appropriate to consider that a specific unit of area on a macroscopic surface that is not an event horizon should, in general, also correspond with the state of a specific matter particle inside that surface. In such a context, it should be possible to associate the information which would allow to identify the state of motion of a particle contained in a microscopic black hole present inside a macroscopic surface with some precise element of area (or perhaps with a precise group of such elements) on that surface. Thus, if all the matter particles present inside some surface can be considered to be locally constrained by a microscopic event horizon, then, even in the absence of a macroscopic event horizon, we would be allowed to assume that the information about the exact state of those particles must be provided in binary units corresponding to specific elements of area on the surface enclosing the volume in which the particles are located. But this actually occurs only when we assume that event horizons must always be present locally, on the Planck scale, so as to constrain the magnitudes of the energies and the orientations of the momenta of matter particles, while leaving undetermined the signs of the space intervals associated with their motion.

Of course, under such conditions, more binary units of information would have to be encoded on the macroscopic surface to specify the exact microscopic state of each of the matter particles it contains, because, in addition to specifying the sign of the space interval associated with the motion of a particle, we would now need to determine the orientation of its momentum as well as the sign of its energy (which also determines the sign of energy of the microscopic black hole which is constraining the motion of the particle locally). Therefore, the amount of information associated with the microscopic state of *matter* inside an ordinary surface would be larger than it would be if this surface was the event horizon of a black hole. In fact, the

configurations for which the entropy associated only with the signs of energy and the momentum orientations of elementary particles would be the highest are those where macroscopic gravitational fields would be absent and the areas of the local event horizons associated with the presence of microscopic (Planck scale) black holes would be the smallest and would be found in the largest number. But given that this occurs when positive- and negative-energy matter are as smoothly distributed as they can be in the available space, then it follows that there would be a compensation between the larger amount of information required to specify the energy signs and the momentum orientations of matter particles and the smaller quantity of information required to describe the exact microscopic state of the gravitational field itself, which would be a consequence of the reduction of its overall strength on the boundary of the region considered, and this is what would allow the Bekenstein bound to continue to apply.

If this account of the physical degrees of freedom of matter associated with the information encoded in the microscopic configuration of the gravitational field on a surface is accurate, it means that we would not be justified to assume that there is no longer anything physically significant going on *at* the Planck scale, because the same degrees of freedom of matter which are reflected in the discrete quantum gravitational degrees of freedom on the event horizon of a macroscopic black hole would also characterize the state of matter particles present on such a scale in the absence of macroscopic event horizons. I was able to draw this conclusion only at a relatively late stage of my research program, because for a long period I had assumed, without much thinking, that the possibility that matter could exist in a negative energy state would imply a cancellation of all quantum fluctuations in energy at the Planck scale, which would not allow for the presence of microscopic black holes on such a scale. But in fact, all that is truly implied by the possibility that negative energy states can be occupied is that the fluctuations in energy can occur in both positive and negative territory. Thus, not only do fluctuations associated with positive and negative energy states not compensate one another out at the smallest physically significant scale of space and time, but it seems that their basic distinction actually provides one of the only significant degrees of freedom characterizing the state of matter on such a scale.

The fact that the proposed description of the constraints imposed on the microscopic state of matter by the gravitational field of a stable-state black

hole can be generalized, in the particular manner described above, to situations in which the density of matter is lower and more homogeneously distributed and particles of opposite energy signs are present together inside a surface, strengthens the argument for the existence of a correspondence between the semi-classical theory of black hole thermodynamics and conventional statistical mechanics (the discussion featuring in the following section will add weight to this conclusion). Indeed, I have already pointed out that the situation we have, in the presence of a macroscopic black hole containing only matter with one energy sign, is analogous, from the viewpoint of gravitational entropy, to a state of thermal equilibrium such as we might encounter in conventional statistical mechanics. But if we are justified to assume that the proposed description of the microscopic degrees of freedom characterizing stable-state black holes can be generalized, by assuming the existence of states (the microscopic black holes) which are similar, locally, to those equivalent thermodynamic equilibrium states, then the analogy could be carried over to the field of non-equilibrium thermodynamics. This is because, in effect, the basic assumption of the thermodynamic theory of irreversible processes is that, even systems evolving irreversibly are to be conceived as being locally and momentarily in a state of near thermal equilibrium. What we have, then, is an ensemble of subsystems in a state of near equilibrium, exchanging energy and evolving in such a way that static equilibrium is not required at the level of the system as a whole (which in the current analogy would be any matter-enclosing surface), like it would be in equilibrium thermodynamics.

It is true that, in the present case, the stability of the configurations occurring on the shortest scale would be limited, because the gravitational fields of the microscopic black holes are continuously fluctuating, but then, the local subsystems in the theory of near-equilibrium thermodynamics are also not in states of perfect equilibrium. What is reflected in this particularity is merely the fact that we are here actually dealing with statistical laws, applying to randomly fluctuating systems, for which deviations away from thermal equilibrium continuously occur locally, even when a system is in a state of overall equilibrium. In fact, the situation we would be dealing with in general would be one where a relatively large number of black holes of various sizes and variable stability (including macroscopic black holes) are present inside a surface and exchange energy with one another. In this context, the microscopic states of matter would be locally constrained to maximum entropy configurations that would frequently fluctuate to lower

entropy configurations (whenever a microscopic black hole would evaporate), as in the local subsystems of the theory of non-equilibrium thermodynamics. But the system as a whole would be allowed to evolve irreversibly toward a state of maximum entropy through the merger of smaller mass black holes into ever more massive ones with larger event horizons, whose areas fluctuate much more slowly. One could hardly think of a more perfect analogy between two theories and I believe that this is not a coincidence, but rather a clear indication that the proposed application of the insights derived while studying the problem of discrete symmetries in the context of the existence of negative-energy matter allows a better understanding of the problem of black-hole entropy as a pure thermodynamic phenomenon in the quantum gravitational regime. In any case, it is clear to me that whatever explanation of the discrete nature of the microscopic degrees of freedom of matter particles would be more accurate than the one provided above, would have to be derived from the mathematical framework of a quantum theory of gravitation that would accommodate the developments introduced so far in this report.

3.11 Negative temperatures

It is not a widely known fact that while temperatures are usually confined to positive values, it is nevertheless unavoidable that some physical systems be attributed negative temperatures under certain conditions. Those who have considered the issue have recognized, in effect, that negative measures of temperature must necessarily occur when we are dealing with certain macroscopic systems with a finite number of energy levels. What happens is that, as temperature rises it must in general be assumed that more energy states become available for the constituent particles, so that the amount of missing information or entropy is itself rising. Therefore, entropy must be assumed to be minimum when a system is at zero temperature. But for systems with a finite number of energy levels, it turns out that, as temperature increases, we may reach the point where entropy is maximum and temperature therefore must be considered infinite. This may occur, for example, in the case of a spin system in a magnetic field, where the number of orientation states of each nuclei is finite. For such a system, the lowest energy configuration is that where all the spins are in the direction of the magnetic field, while the highest energy configuration is that which occurs when all the spins are oriented in the direction opposite that of the magnetic field. At infinite tem-

perature all spins would be oriented in the most random way, with as many spins oriented in the direction of the magnetic field as there would be in the opposite direction. If we were to add more energy to a system in such a state, we would witness a decrease of its entropy, as more spins would become oriented in the direction opposite the magnetic field and less information would be required to describe the unknown microscopic state of the system.

Given that temperature merely defines the relationship which exists between energy and entropy, if an increase of energy produces a decrease of entropy, then it must necessarily be assumed that the temperature has become negative. But if adding more energy decreases the entropy only slightly when it reaches its maximum point, at which the temperature is infinite, then it means that the temperature is not ‘minus zero’, but actually ‘minus infinity’. Thus, as even more energy is added to the system, the entropy would gradually decrease back to a minimum, at which point the negative temperature would actually reach the zero value again. In the case of the spin system, this point would be reached when *all* the spins would be oriented in the direction opposite that of the magnetic field and no further change could occur. I may also mention that it was found that when we combine two such systems which happen to have opposite temperatures of equal magnitude, the outcome must be a system with infinite temperature. It must be understood that, despite common expectation to the effect that temperature is a positive-definite quantity, the conclusion that negative temperatures may occur in nature is not just a consequence of adopting some particular definition for what temperature should be, or of choosing a particular reference scale for this quantity¹². Specialists are unequivocal concerning the fact that negative temperatures cannot be avoided in a general context, because they are associated with actual states of any system with a finite number of energy levels.

Now, what I would like to point out is that black holes are somewhat similar, from a thermodynamic viewpoint, to those more conventional systems for which negative temperatures are allowed. Indeed, I have already explained in section 2.13 that it appears necessary to attribute to ordinary negative-energy systems a measure of temperature that is itself negative. This negative temperature becomes necessary in the context where negative-

¹²Of course this statement is only significant under the assumption that we are dealing with ‘absolute’ measures of temperature, like those provided by the Kelvin scale and not with measures of temperature where the zero level is arbitrarily fixed to a non-thermodynamically significant quantity, as is the case with the Celsius or Fahrenheit scales.

energy systems absorb thermal energy or heat as a negative-definite quantity, even though the entropy of their matter can be expected to rise when such a change takes place, as is the case for positive-energy systems. But there is no reason to believe that the situation is any different when black holes are involved and therefore it would appear necessary to attribute to negative-energy black holes a negative temperature. This is certainly appropriate, given that it would seem that if a positive-energy black hole has a positive value of surface gravitational field, then a negative-energy black hole would have a negative value of surface gravitational field, and knowing that the surface gravitational field is the quantity which is associated with the temperature of a black hole in the semi-classical theory, I'm led to conclude that this temperature itself needs to be allowed to vary, not just in magnitude, but also in sign. Actually, this can be considered an absolute requirement in the context where a negative-mass black hole would radiate particles with an energy sign opposite that of the particles radiated by a positive-mass black hole, while the same changes to entropy would be required to take place as a consequence of the decay process. Thus, if negative-energy matter exists, it would seem that some black holes could, in effect, be attributed negative temperatures, which would be made conspicuous by the reversal of their surface gravitational fields.

It should not come as unexpected, therefore, that there exists a certain correspondence between the thermodynamics of black holes and the above described thermodynamic phenomenon involving spin systems. Of particular significance is the fact that, as a positive-energy black hole evaporates through the emission of thermal radiation and its mass decreases toward zero (in positive territory), its temperature would rise until it becomes infinite (which would occur when the object reaches the Planck mass), at which point, if we were to continue to remove energy from it (by actually adding negative energy) its mass would start to increase into *negative* territory with an initial temperature that would be infinite, but negative, and which would decrease in magnitude (toward zero) as the negative mass of the object increases. Of course, the dependence of temperature on total energy is not exactly the same here as in the case of spin systems, given that a larger-mass black hole would have a lower temperature. But if we consider only the relationships between thermodynamic properties, then the analogy is valid. Also, if we were able to combine a positive-energy black hole (to which is associated a positive temperature) with a similar negative-energy black hole (to which is associated a negative temperature), then what we would obtain

is not a zero temperature object, but an object with a larger and possibly infinite temperature (just like when we combine two opposite-temperature systems in the conventional theory), because the mass of the resulting black hole would be smaller, which means that it would radiate energy at a higher rate. Of course, it may not be possible, from a practical viewpoint, to combine opposite-energy black holes, so as to produce a lower-mass object, but mathematically the correspondence between the quantities involved is valid and matches the expectations derived from conventional thermodynamics theory.

The fact that the existence of such a beautiful correspondence between the semi-classical theory of black-hole thermodynamics and the classical thermodynamics of systems with a finite number of microscopic levels of energy is allowed to occur, under the hypothesis that two signs of mass are relevant for a description of the thermodynamics of black holes, constitutes an additional argument for recognizing the legitimacy of this theoretically motivated insight. In fact, I'm surprised that the conclusion drawn by specialists, concerning the unavoidable character of the concept of negative temperature, was never considered to imply that energy itself should be allowed to vary in sign rather than only in magnitude. But as I have always believed that the true motivation behind the widespread idea that energy can only be positive originates from the thermodynamic conception of energy as a measure of heat (which is itself a positive-definite quantity from a conventional viewpoint), I was quite satisfied when I learned that this most thermodynamic concept of all, the temperature, must itself vary in sign. If there is no reason to assume that negative temperatures cannot have a clear significance in physical theory, and if it turns out that they must ultimately be associated with the state of objects whose energy is predominantly negative, then we have one less argument for assuming that the concept of negative energy itself cannot be given clear meaning.

Chapter 4

Cosmology and Irreversibility

4.1 The outstanding problems of cosmology

The situation we face today in the field of theoretical cosmology can be resumed by mentioning two broad categories of problems. The first issue has to do with dark energies in general and the consequences of the existence of invisible forms of matter and energy on the gravitational dynamics of visible matter. One of the main difficulties regarding dark energies has to do with explaining how it is possible for the density of vacuum energy to be as low as one observes it to be, while not being exactly null. Indeed, with the discovery that the expansion of space is accelerating, it has become necessary to recognize that some invisible form of positive energy with negative pressure is present in empty space and in the present theoretical context the only plausible explanation we have for this phenomenon is that it is a consequence of zero-point vacuum fluctuations. But such a small value for the cosmological constant is unexpected and therefore one is encouraged in seeking alternative and more exotic interpretations for this dark energy. In the first portion of the present chapter I will explain that it is, in fact, still possible to assume that dark energy is attributable to the existence of a non-vanishing average value for the density of vacuum energy and I will show that this hypothesis is not invalidated by the otherwise inexplicably small, but non-zero value of the cosmological constant.

Another aspect of the problem of dark energies has to do with the phenomenon of missing mass which arises because it appears that the visible material that is present in galaxies and clusters of galaxies does not provide

enough gravitational force to explain the motion of the astronomical objects that compose those large-scale structures. Here, one of the main objectives usually consists in trying to determine the exact nature of the dark-matter particles which are assumed to contribute additional gravitational attraction around visible structures in the positive-energy matter distribution. Despite all the efforts which were devoted to this task, this is a problem which has remained unsolved. But as I will soon demonstrate, it is possible, in the context of the developments which were introduced in the second chapter of this report, to explain most of the missing-mass effects observed around galaxies and clusters as being another, perhaps more unexpected, consequence of the existence of zero-point vacuum fluctuations. However, the presence of underdensities in a uniform distribution of negative-energy matter can also be expected to contribute to the missing-mass effect experienced by positive-energy objects under particular circumstances and therefore I will examine the consequences of such a phenomenon on the formation of large-scale structures.

The other broad category of issues we are currently dealing with in cosmology could be called the inflation problem. This may sound paradoxical, as inflation presently constitutes a dominant paradigm for theoretical cosmology and is still believed to offer solutions to many serious problems in the field. If I'm allowed to speak about a problem concerning inflation, it is because there does exist a series of issues which were most accurately described by the originators of inflation theory and which have long been considered to be appropriately solved by one or another instance of such a model, until it became clear that the theory actually offers so much predictive freedom that it is nearly unfalsifiable. As the following discussion progresses, it will become clear that what made the inflation paradigm so successful is mainly an absence of alternative solution to the various problems it was originally proposed to address. Given that I believe that the most important contribution of the originators of inflation theory was to show that there does remain decisive, unresolved issues in cosmology, which could perhaps be solved using their theory, then I will not refrain from discussing those issues as a genuine category of problem to which new solutions can be proposed, even in the context where we do not reject the basic idea that there may have occurred a short period of exponentially accelerated expansion in the first instants of the Big Bang.

Two different aspects of the inflation problem will be discussed in this chapter. The first aspect has to do mainly with the problem of flatness, or

the fact that the present rate of expansion of matter on the cosmological scale appears to be set to some unnatural value, which requires an extremely precise adjustment of parameters in the initial state at the Big Bang. I will explain that in the context of the progress I have achieved while solving the cosmological-constant problem, this difficulty occurs merely as a consequence of our failure to appropriately recognize that the constraint of relational definition of physical attributes must also apply to the energy of the universe. The other aspect of the inflation problem which I will address is the horizon problem, which has to do with the fact that it is not possible to explain the uniformity of the very-large-scale distribution of matter energy as being a consequence of smoothing processes that would obey the principle of local causality. Two further issues actually constitute particular aspects of the horizon problem. They are the smoothness problem and the problem of topological defects. Actually, the smoothness problem would not exist if it was not for the fact that it is usually assumed that a solution to the horizon problem would have for consequence to leave the universe perfectly homogeneous, therefore requiring an independent explanation for the fact that some inhomogeneities nevertheless remained in the primordial matter distribution, which gave rise to present-day structures. It will be shown that inflation is not required to solve this problem and perhaps also that which is associated with the rarity of topological defects, given that those difficulties arise merely as a consequence of the inappropriateness of inflation theory as a solution to the horizon problem.

The one truly amazing consequence of the particular approach I followed in dealing with the horizon problem, however, is that it allowed me to gain a new perspective on another decisive problem which is not always recognized as a problem for cosmology, despite the fact that it can be traced back to the particular boundary conditions which were in effect at the Big Bang. This is the problem of the origin of the arrow of time, which is probably the most serious difficulty currently facing cosmology. It is merely the fact that the problem is so old, and has remained unsolved for so long, that explains that it is often not recognized as a problem for cosmology, as if we had long ago given up trying to resolve it. But the developments which have been introduced in the preceding two chapters and those which will be discussed in the second portion of the current one will allow to confirm the cosmological nature of the issue and will culminate in providing the first-ever plausible explanation of how it can be that a fully time-symmetric fundamental theory can conspire to enforce boundary conditions which give rise to irreversible evolution and

the second law of thermodynamics.

We therefore have two broad categories of problem in cosmology, which are the problem of dark energies and the inflation problem and which each involve several different aspects. I will first discuss the cosmological-constant problem, along with the problem of missing mass, as particular aspects of the problem of dark energies, which will then allow me to approach the problem of structure formation from a new perspective. The progress achieved while solving the cosmological-constant problem will then enable me to provide a satisfactory solution to the flatness problem as one particular aspect of the inflation problem. Then I will discuss the horizon problem as another aspect of the inflation problem, but while addressing this issue and the related problem of the origin of primordial inhomogeneities I will contribute significant insights into the nature of gravitational entropy that will provide the necessary means to formulate a definitive solution to the problem of the origin of time asymmetry.

4.2 The cosmological-constant problem

One of the key parameters of the standard model of cosmology that remains unexplained is certainly that which we call the cosmological constant. If there is often reticence to assume that the cosmological constant is a manifestation of the energy contained in zero-point vacuum fluctuations, it is certainly because it is normally expected that the density of energy contained in the vacuum at the present epoch should be either precisely null (due to some unknown symmetry principle) or much larger than the energy density we may associate with the observed, current value of the cosmological constant. It appears much more natural, therefore, to assume that we are rather dealing with some dark energy of unknown nature whose density could vary with the expansion of space, like that of matter. If dark energy is merely a material substance with negative pressure, then it would appear natural to assume that it should now have a density similar to that of matter, while it seems rather unlikely that vacuum energy would simply happen to have a density comparable to that of matter (visible and dark), given that the density of vacuum energy is usually assumed to be unaffected by expansion. Thus, either dark energy is not vacuum energy, in which case we have no idea what its material nature is, or we restrict ourselves to known phenomena and we recognize that it must be vacuum energy, in which case we need

an explanation for the observed similarity between the current value of the energy density of matter and that of vacuum fluctuations, that is to say, we need to explain how it can be that the vacuum contains so little energy, and yet does not provide a null contribution to the universe's energy budget, as we usually assume should have been the case if some symmetry principle was responsible for the fact that this energy is much smaller than the natural value associated with the quantum gravitational scale, which is more than 120 orders of magnitude larger than the observed value.

I find it significant that the problem associated with the small value of the cosmological constant is usually recognized to be a disagreement between the viewpoint of experimentalists and that of theoreticians, because, from that perspective, it becomes apparent that resolving the issue will necessarily require reconsidering the validity of certain hypotheses we take for granted in the current theoretical context. First of all, it must be acknowledged that despite the fact that the empirical determination of a positive value for the cosmological constant contributed to reinforce the traditional belief that any energy density that could be associated with this parameter should probably be positive, this restriction would be totally unjustified in the context of the progress achieved in the second chapter of this report. Thus, vacuum energy, in particular, could certainly have been negative and the only thing we can be certain about is that it is the observer independent sum of all positive and negative contributions to vacuum energy density which would have an effect on the expansion rates experienced by positive- and negative-energy observers, unlike would be the case with a material substance like quintessence with pressure opposite its energy sign, which would only influence the expansion rate measured by a positive-energy observer through its positive energy component, as any smooth matter distribution with both a positive- and a negative-energy component. Therefore, in the context of the developments discussed in section 2.6, it may perhaps look like quintessence has an advantage over vacuum energy as a candidate for dark energy, in that it could produce the desired effect even when the material contains just as much positive energy as it contains negative energy. But I will show that this is not really the case and that the advantage rather goes to vacuum energy for at least originating from known physical principles applying to known forms of matter, or forms of matter whose existence can be deduced from known principles.

There is a certain similarity between the prediction of an arbitrarily large magnitude of energy in zero-point vacuum fluctuations and the old problem

of the ultraviolet divergence of black body radiation which was solved by the creation of quantum theory. I believe that the commonly met suggestion that a cut-off may come about in the calculation of the density of vacuum energy, which would be associated with the quantized nature of space at the most elementary level is certainly appropriate, but it is also insufficient to solve the cosmological-constant problem. Indeed, such a cut-off would simply decrease the energy contributions from their potentially infinite values to very large values associated with the scale of quantum gravitational phenomena and those various energy contributions would still need to cancel out in order to produce the much smaller observed value. This is precisely the problem we face right now: the required cancellation must occur by chance out of a myriad of potentially enormous, independent contributions to the energy of the vacuum. The validity of the hypothesis that space itself must be submitted to quantization (so that there must exist a maximum theoretical value of vacuum energy density) is certainly quite inevitable, especially in the context of the developments introduced in section 3.10 concerning black-hole entropy and the relationship between discrete symmetry operations and the microscopic states of the matter that crosses the event horizon of such an object. But even if this assumption is well-founded, it is simply inadequate all by itself to reconcile the theoretically derived and observationally inferred values of vacuum energy density.

In fact, I believe that we have no choice but to assume that some symmetry principle must be responsible for the almost perfect cancellation that gives rise to the observed small value of vacuum energy density, because under current assumptions there would be virtually no limit to the expected value of this parameter, which would then be more likely to have a relatively high positive or negative value. However, I also share Feynman's opinion that it may not be quantum field theory, or the preferred Grand Unified Theory, which needs to be modified in order to accommodate such a requirement, but rather our current theory of gravitation. Indeed, the generalized gravitation theory I have introduced in chapter 2 has allowed me to identify a new category of matter particles with negative action sign, with which we may naturally expect to be associated a contribution to the energy of zero-point vacuum fluctuations which would be opposite that associated with positive-action matter particles.

It is true that there are already both positive and negative contributions to the energy of the vacuum in the context of traditional theories, but it is simply too unlikely that the required outcome could arise by chance from an

extremely precise cancellation of the countless, independently varying, positive and negative contributions which are normally taken into account. What I'm suggesting is that there exists a whole new class of contributions whose total energy must necessarily compensate the sum of all currently considered contributions to the energy of the vacuum. Indeed, in the context where there must be a symmetry under exchange of positive and negative energy states, we are allowed to expect that the energy of the vacuum should actually be null, because negative-energy observers would necessarily experience vacuum fluctuation processes which contribute energies that are the exact opposite of those contributed by the vacuum fluctuation processes which are experienced by positive-energy observers and which are the only type of vacuum fluctuations currently taken into account from the viewpoint of conventional quantum field theory. This is a consequence of the fact that, while only one category of positive and negative energy fluctuations directly interacts with positive-energy matter, both categories of contributions exert a *gravitational* influence on positive-energy matter and must be taken into account in determining the current value of the cosmological constant measured by a positive-energy observer.

From my viewpoint, the presently considered negative contributions provided by certain particles present as zero-point vacuum fluctuations would become the positive contributions of those same particles in the negative-action sector of quantum field theory (that which describes the processes which directly affect negative-energy matter other than through their gravitational influence) and the currently considered positive contributions provided by other particles, also present as zero-point vacuum fluctuations, would become the negative contributions of the same particles in the negative-action sector of quantum field theory. This would be true despite the fact that, as I explained in section 3.9, there are actually four distinct action-reversal symmetry operations, which can be violated in different proportions, because, when we are considering all possible processes occurring in the vacuum, we are actually dealing with the outcome of all those operations combined and as I explained in the same section, there must be invariance under such a combination of all action-reversal symmetry operations that relate positive-energy matter to negative-energy matter.

Thus, all currently considered contributions to the energy density of the vacuum, whether they are positive or negative, must have a counterpart of equal magnitude and opposite sign, which guarantees a cancellation of all contributions, regardless of the details of the Grand Unified Theory chosen to

describe elementary particles and their interactions. It is not the conclusion that there are no unexpected cancellations among the multiple independent terms which add up to produce the total energy density of that portion of vacuum fluctuations experienced by positive-energy observers which is wrong, but the ignorance of the fact that there is a corresponding set of contributions, experienced only by negative-energy observers, whose distinguishing feature is that all of its terms contribute energies which are naturally the opposite of those which are already taken into account, as a consequence of the requirement of symmetry under exchange of positive and negative energy states. It is merely the fact that no fully consistent theory incorporating the concept of negative-energy matter had ever been formulated that justified the implicit assumption that no contributions of the kind proposed here needed to be taken into account, because from that perspective the whole idea that virtual processes could take place in the vacuum that would interact merely with negative-energy matter appeared meaningless, as no such matter would exist in our universe.

The usual remark to the effect that it is highly unlikely that all contributions to the energy of the vacuum could conspire to produce a vanishing density is justified, but only in the context where the sole class of contributions which is recognized to exist is that which is associated with those zero-point fluctuations and virtual particles which exert a direct influence on positive-energy matter. However, if we recognize the unavoidable character of the assumption that negative action states are not forbidden, then it would seem that we can now predict a vanishing value for the energy of the vacuum. It is no longer necessary to assume that there occurs a miraculous conspiracy, that results in the numerous, currently envisaged, independent contributions to vacuum energy density adding up to produce a number several orders of magnitude smaller than those individual terms. It is also no longer required that the details of some Grand Unified Theory be invoked that would allow to derive the existence of such a precisely adjusted set of independent contributions in order for the right outcome to be derived. We are not really looking for compensations among multiple unconstrained parameters, but for an overall cancellation among two identical sets of parameters, whose corresponding elements have equal magnitudes and opposite signs, even on the low-energy scale at which the symmetries associated with the unified theory are spontaneously broken. This does not mean that there must be a cancellation of energy fluctuations locally on the Planck scale, however, because, as I mentioned in section 3.10, even the sign of energy must be considered a vari-

able parameter on such a scale (in the absence of a macroscopic event horizon to constrain the states of matter particles) and it is merely on the scale at which classical gravitation theory applies that a cancellation of positive and negative contributions is allowed to occur.

What's surprising, therefore, is not that the cosmological constant is presently so small, but rather that it is, in effect, not perfectly null. But even if this may not be as serious a problem as that of the discrepancy between current estimates of vacuum energy density and the actual value of this parameter provided by astronomical observations (given that in the present case the amplitude of the required adjustment is much smaller than that which would have to occur in the context of a traditional model), it would not be appropriate to assume that the progress achieved so far in this section constitutes a complete solution to the cosmological-constant problem. What I will now explain is that, despite the fact that it is natural to expect that there should be a perfect compensation between the currently considered contributions to vacuum energy density and the additional contributions arising from the presence of those virtual particles which directly interact only with negative-energy matter, it is nevertheless possible, in principle, for the cosmological constant to take on arbitrarily large values, even though it does appear that, for some reason, the magnitude of vacuum energy density was negligible compared to the magnitude of positive and negative matter energy densities in the very first instants of the Big Bang.

Faced with the dilemma presented here, I must acknowledge that I initially tried to explain how it can be that we appear to measure a small but non-vanishing value for the cosmological constant by assuming that, in fact, the cosmological constant is actually null, while the effects we attribute to it, instead of being the consequence of a non-zero density of vacuum energy, are rather a consequence of the presence of a very-large-scale inhomogeneity in the invisible negative-energy matter distribution. Indeed, as I explained in section 2.8, an overdensity of negative-energy matter should produce an outward-directed (repulsive) gravitational force on positive-energy matter. Thus, if we happen to be located near the center of such a very-large-scale overdensity of negative-energy matter we should expect to observe a 'local' acceleration of the rate of expansion that would merely be a consequence of the presence of this inhomogeneity in the invisible distribution of negative-energy matter. In fact, it was also suggested by others that just the opposite might be occurring and that we may be located inside an *underdensity* in the distribution of *positive-energy* dark matter, which would exert a similar

outward directed gravitational force on positive-energy matter.

But it is precisely here that a problem occurs with my own original hypothesis, because it was later shown that the accelerated expansion which was revealed by observations of high-redshift type Ia supernovae is incompatible with any such explanation of the acceleration of expansion. In fact, in the context where there is a constraint on the amplitude of density fluctuations arising from the uniformity of the cosmic microwave background, it appears that there simply could not have existed inhomogeneities of sufficiently large magnitude to provide an alternative explanation of the acceleration of expansion. What's more, if we recognize the observational and theoretical necessity of a critical density of positive energy, then we have an additional argument to reject such an explanation for the acceleration of the rate of expansion, because we actually need the additional positive energy that would be contained in the vacuum in order to reach the critical density, which cannot be provided by positive-energy matter alone¹.

It must be acknowledged, therefore, that despite the fact that, in the context of the developments proposed in the preceding chapters, we may expect the natural value of vacuum energy density to be zero, there must nevertheless exist an imbalance between the positive and negative contributions to vacuum energy density. What must be understood is that this imbalance cannot be attributed to a violation of the symmetry under exchange of positive and negative energy states, which is a necessary requirement of the constraint of relational definition of physical properties. At this point it is necessary to recall the definition of the vacuum-energy term that emerged from the generalized gravitational field equations developed in section 2.15. There, I proposed that the value of vacuum energy density which, on a global scale, would be associated with the cosmological constant measured by a positive-energy observer be defined as the sum of the natural vacuum-stress-energy tensors $\gamma^{-+}\mathbf{V}_P^+$ and $-\mathbf{V}_P^-$, which provide the maximum positive and negative values of energy density contributed by those portions of zero-point vacuum fluctuations that directly interact only with negative- and positive-

¹The same argument can also be used to rule out the possibility that dark energy could actually consist of gravitationally-repulsive negative-energy matter of the traditional kind, which would repel both positive-energy matter and negative-energy matter itself, because such material would contribute negatively to the energy budget and while it would not form local structures, it would interfere with current estimates concerning the initial rate of expansion of matter at the Big Bang (which allow to successfully predict the observed abundance of light chemical elements), when its density would be much larger.

energy matter (respectively), but which both exert a gravitational influence on positive-energy matter:

$$\mathbf{T}_V^+ = \gamma^{-+}\mathbf{V}_P^+ - \mathbf{V}_P^- \quad (4.1)$$

From that particular viewpoint it would appear clearly inappropriate to consider the existence of a ‘bare’ cosmological constant, distinct from that which would be associated with the energy contained in zero-point vacuum fluctuations, because the cosmological term $\mathbf{T}_\Lambda = -\Lambda\mathbf{g}$ that enters the original form of the gravitational field equations associated with a positive-energy observer (with Λ as the cosmological constant) must now be understood to consist of the globally uniform portion of the locally variable, vacuum-energy term \mathbf{T}_V^+ measured at a given epoch of cosmic time and this means that it must be considered a particular form of vacuum energy, even in a purely classical context.

Now, what is significant in the above equation is the appearance of the metric conversion factor γ^{-+} in front of the positive contribution to vacuum energy density, which becomes necessary once we recognize that the portion of vacuum fluctuations that cannot be directly experienced by a positive-energy observer (other than through its gravitational influence) is, in effect, the one that provides a maximum positive contribution $\gamma^{-+}\mathbf{V}_P^+$ to the density of vacuum energy, while the portion that can be directly experienced only by a positive-energy observer would be the one that provides a maximum negative contribution $-\mathbf{V}_P^-$. This is what justifies submitting the maximum positive contribution to the same metric conversion factor as applies to the measures of negative-energy matter density effected by positive-energy observers, because, in the absence of direct interactions, it cannot be assumed that the metric properties of space which govern this portion of vacuum fluctuations are necessarily the same as those experienced by a positive-energy observer. In section 2.15 I mentioned, in effect, that the γ^{-+} conversion factor is the mathematical object that allows to map the metric properties of space experienced by negative-energy observers onto those experienced by positive-energy observers. But if the portion of zero-point fluctuations that provides a maximum positive contribution to the density of vacuum energy, is directly experienced only by negative-energy observers, then from the viewpoint of positive-energy observers the measure of energy density involved must be submitted to the same metric conversion factor as applies to measures of negative-energy matter density.

It must be understood, therefore, that the maximum positive contribution to the energy of the vacuum is not the sum of all positive contributions directly experienced by both positive- and negative-energy observers, but really the sum of all contributions, positive and negative, which are directly experienced only by a negative-energy observer and which must necessarily produce a maximum positive outcome (given that negative-energy matter must by definition consist of voids in positive vacuum energy and such voids cannot be considered not to interact with that very portion of vacuum energy in which they develop and propagate). Thus, while the hypothesis that the sum of all contributions to the density of vacuum energy experienced by a positive-energy observer produces a negative number (while the sum of all such contributions which are directly experienced by a negative-energy observer produces a positive number) may at first, perhaps, appear arbitrary, it is actually unavoidable in the context where it cannot be assumed that the density of negative matter-energy itself could be directly determined (other than through its gravitational influence) by a positive-energy observer, while that would be allowed if such an observer could directly measure the actual value of the maximum positive contribution to the density of vacuum energy that would be reduced by the presence of negative-energy matter, as I explained in section 2.15. Thus, if one considers the measure of vacuum energy density that is contributed by the maximum positive-energy term as it is perceived by a positive-energy observer, then it must be submitted to metric conversion. But even though the necessity of such a mapping is justified by the absence of direct interaction between positive- and negative-energy matter, its legitimacy can only be understood based on considerations of a cosmological nature.

First of all, it must be noted that the magnitude of positive vacuum energy density which would be associated with the natural vacuum-stress-energy tensor \mathbf{V}_p^+ that is directly experienced by an observer made of negative-energy matter is an invariant quantity, which according to the requirement of symmetry under exchange of positive and negative energy states should be the same as that which is provided by the magnitude of negative vacuum energy density associated with the natural vacuum-stress-energy tensor \mathbf{V}_p^- that is directly experienced by a positive-energy observer. Thus, if, in the context where the vacuum-energy term does not vanish to zero, there must be a difference between the maximum positive and negative contributions to vacuum energy density, it can only arise because the metric properties of space that determine the magnitude of the positive contribution, as it is

perceived by a positive-energy observer (through its gravitational influence), are not the same as those that determine the magnitude of the same positive contribution, as it is perceived by a negative-energy observer (and similarly for the maximum negative energy contribution). What I'm suggesting is that this means that the appearance of the metric conversion factors in the definition of the net values of vacuum energy density is a consequence of the fact that the volume of space contained within a given boundary may vary depending on whether this volume is measured by a positive- or a negative-energy observer, so that the same invariant magnitude of vacuum energy densities can provide different contributions for observers of opposite energy signs.

Now, when I introduced the notion of observer-dependent gravitational fields, which gives rise to observer-dependent metric properties, I emphasized that it must be recognized that there is still a correspondence between the local topology of space associated with positive-energy observers and that which is associated with observers of opposite energy sign. Thus, the set of events occurring in spacetime must be the same regardless of the way the metric properties of space are perceived, which also means that every particle that is present inside a surface parameterized using the metric properties of space associated with a negative-energy observer must also be present within a corresponding surface parameterized using the metric properties of space associated with a positive-energy observer, even when the volume contained inside the surface varies as a function of the sign of energy of the observer. In such a context, even when the ratio of the average densities of positive- and negative-energy matter is fixed from the viewpoint of any given observer (as would be the case before the early annihilation of matter with antimatter, for reasons I will explain later), the average densities of both positive- and negative-energy matter could be different for observers with opposite energy signs, which do not share the same metric properties. The crucial point, here, is that those observer-dependent metric properties may not only differ locally, as a consequence of the curvature of space attributable to the presence of positive and negative matter energy, but may also be different on a cosmological scale, when a discrepancy emerges between the scale factors experienced by opposite-energy observers.

To visualize the nature of the relationships between the measures of energy density perceived by positive- and negative-energy observers on a cosmological scale, it may help to consider the analogy provided by the case of a universe with bi-dimensional space and closed geometry. More specifi-

cally, we may imagine two spherical surfaces centered on the same point (in three-dimensional space) which would represent the entire volumes of space experienced by opposite-energy observers². It would then be appropriate to assume (for reasons that will be discussed later) that initially, in the very first instants of the Big Bang, the two surfaces both have minimum areas which correspond to a state of maximum positive and negative energy densities. Under such conditions, the average densities of positive- and negative-energy matter particles determined using the metric properties of space associated with one of the surfaces would initially be exactly the same as those which are determined using the metric properties of space associated with the other surface. But even if we assume that the average matter densities only vary as a result of expansion (additional variations can be expected to arise as a result of the early annihilation of matter with antimatter), as space expands and the two closed surfaces grow in size, any difference in their expansion rates would make their respective areas to differ. Yet, even if such a divergence was to develop, to each position of a particle on the smaller surface would still correspond a unique position on the larger surface associated with observers of opposite energy sign and to each boundary on the smaller surface would correspond one larger boundary on the other surface. In the absence of any *local* variations in the metric properties of space experienced by opposite-energy observers, the only difference which would characterize the matter distributions observed on the two surfaces would therefore be the difference between the magnitudes of their average densities, which would follow from the fact that the same particles occupy spherical surfaces with different total areas.

Even in the absence of local space curvature, therefore, it seems that the metric properties of space could differ for observers of opposite energy signs, because, on the cosmic scale, regions of space delimited by corresponding boundaries (associated with observers of opposite energy signs) could have different volumes depending on the sign of energy of the observer that determines this volume, if the scale factor determined by positive-energy observers is different from that which is determined by negative-energy observers. This is due to the fact that the present average densities of positive- and negative-energy matter measured by a positive-energy observer are allowed to differ

²It must be clear that the situation described here is only valid as an analogy, because, as I will explain in section 4.5, in a more realist context it is not even possible for space to be closed from both the viewpoint of positive-energy observers and that of negative-energy observers.

from those measured by a negative-energy observer, even when there was no difference, initially, between the scale factors experienced by observers with opposite energy signs, given that it is possible for the rate of expansion measured by positive-energy observers to differ, or to have differed at some point, from that which is measured by negative-energy observers.

I believe that what is implied by the appearance of the metric conversion factors in the proposed definitions of the density of vacuum energy, therefore, is that the invariant, maximum, negative and positive contributions (\mathbf{V}_P^- and \mathbf{V}_P^+) to the energy density of the vacuum can be made to differ not only in sign, but also in magnitude, as a consequence of the fact that opposite-energy observers do not necessarily share the same metric properties of space, even on the global scale, where it can be expected that matter is homogeneously distributed. The rule would be that when the scale factor is measured as being proportionately smaller by a positive-energy observer, the density of the maximum positive contributions to the energy of the vacuum (which cannot be directly measured by such an observer) would be increased from the viewpoint of this observer, in comparison with the density of the maximum negative contributions to the energy of the vacuum which is measured by the same observer, so that according to equation (4.1) above, the average density of vacuum energy would be positive and our positive-energy observer would measure a positive cosmological constant Λ . This would be due to the fact that, from the viewpoint of an observer that measures a smaller volume of space on the cosmological scale, those vacuum fluctuations whose invariant maximum energy density can only be directly measured by an observer of opposite energy sign would actually take place within a comparatively larger volume and would therefore appear to have a higher positive or negative energy density (when projected on the smaller volume of space perceived by the first observer) which means that they would provide a larger contribution than the vacuum fluctuations whose invariant maximum energy density our observer can directly measure.

A definite relationship would therefore exist between the net value of average vacuum energy density or the cosmological constant and the difference between the scale factors determined by observers with opposite energy signs, which is made even more significant by the fact that the cosmological constant must itself modify the rates of expansion experienced by positive- and negative-energy observers which determine those scale factors. Thus, if the current value of the cosmological constant is positive, it means that any volume of space, enclosed by a sufficiently large boundary, that would be de-

terminated using the metric properties of space experienced by positive-energy observers must presently be smaller than the corresponding volume which would be determined based on the metric properties of space experienced by negative-energy observers. I'm therefore allowed to predict that if those volumes were exactly the same in the initial Big Bang state, as I will propose in section 4.5, then space must have expanded at a smaller rate, from the viewpoint of positive-energy observers, during a certain portion of the universe's history, in comparison with the rate at which it expanded from the viewpoint of negative-energy observers.

The problem that may seem to arise, under such conditions, is that the smaller scale factor presently experienced by positive-energy observers can be expected to produce a positive cosmological constant that would actually contribute to accelerate the rate of expansion of space determined by positive-energy observers and to reduce any difference between this measure of the scale factor and that which is determined by negative-energy observers. Indeed, while a positive cosmological constant would contribute to accelerate the expansion of space from the viewpoint of a conventional, positive-energy observer (due to the larger contribution of its negative pressure), it would actually contribute to decelerate the expansion rate for a negative-energy observer (again as a result of its negative pressure), which would have for consequence to reduce the divergence between the scale factors associated with observers of opposite energy signs. It may, therefore, appear that the current conditions could only be realized if the observed positive value of the cosmological constant did not arise as a result of positive-energy observers experiencing a smaller scale factor, but rather as a consequence of those same observers experiencing a *larger* scale factor, that would then contribute to further accelerate the rate of expansion of space experienced by those the same observers.

This is the reason why I originally thought that the empirical evidence appeared to support the hypothesis that, contrarily to what I have proposed above, positive-energy observers should directly experience a maximum contribution to vacuum energy fluctuations that happens to be positive (while negative-energy observers should directly experience a maximum contribution that is negative). One must recognize, in effect, that if one was to assume that the sum of all contributions to the energy of the vacuum which are directly experienced by a positive-energy observer actually produces a maximum positive number, then a different form of the generalized gravitational field equations would have to be adopted, such that, from the viewpoint of a

positive-energy observer, the metric conversion factor would rather apply to the negative portion of the maximum contribution to vacuum energy density, thereby giving rise to a modified version of the vacuum-energy term:

$$\mathbf{T}_V^+ = \mathbf{V}_P^+ - \gamma^{-+} \mathbf{V}_P^- \quad (4.2)$$

(this equation is to be compared with equation (4.1) above). If I had originally believed that this alternative form of the vacuum-energy term was a more appropriate choice to model the evolution of the densities of matter energy, it is because I had difficulty seeing how the universe could have evolved in such a way that the scale factor experienced by negative-energy observers could have become so much larger, in comparison with the scale factor experienced by positive-energy observers, that the cosmological constant which results from this divergence could have grown into a positive value that is much larger than the density of positive-energy matter experienced by positive-energy observers (whose magnitude would already be larger than the density of negative-energy matter experienced by negative-energy observers), while according to the definition of the generalized gravitational field equations that gives rise to the vacuum-energy term provided by equation (4.1), one would expect that any difference that develops between the scale factors experienced by opposite-energy observers would rather tend to be reduced by the gravitational force attributable to the pressure of the vacuum. Thus, it appeared desirable to assume that the alternative form of the vacuum-energy term provided by equation (4.2) applies, because that would allow vacuum energy to produce the very conditions which allow it to grow even larger.

What must be clear, first of all, is that what we measure, through astronomical observations, at the present epoch, is an acceleration of the rate of expansion that is experienced only by observers with our own sign of energy, while observers with an opposite sign of energy could measure a different variation of the rate of expansion, not just because the same vacuum energy would exert an opposite gravitational force on negative-energy matter, but because only the average density of positive-energy matter influences the rate of expansion determined by positive-energy observers, while only the average density of negative-energy matter influences the rate of expansion determined by negative-energy observers. But this means that a positive cosmological constant would only be allowed to actually produce an acceleration of the rate of expansion experienced by a positive-energy observer that could reduce the divergence of the scale factors experienced by opposite-energy observers

if the average density of positive *matter* energy determined by a positive-energy observer is sufficiently smaller than the average density of positive *vacuum* energy. Now, while we do observe the present average density of positive matter energy (both visible and dark) to be somewhat smaller than the current average density of vacuum energy and while this can be expected to produce an acceleration of the rate of expansion observed by positive-energy observers, which will reduce the average, positive density of vacuum energy and the magnitude of the cosmological constant to a smaller value, the situation may have been different at a certain epoch in the past.

It is not possible, in effect, to conclude that the magnitude of the positive cosmological constant must have been decreasing at all times in the past, as a result of the effect it exerts on the expansion rates experienced by positive- and negative-energy observers. This is not due merely to the fact that the average value of vacuum energy density must have already been null initially if the scale factors experienced by positive- and negative-energy observers were themselves precisely equal in the first instants of the Big Bang (as I will propose in section 4.5), it is also due to the fact that the average density of matter may not only change as a result of expansion, but also as a consequence of the early annihilation of matter with antimatter. Thus, even if the magnitudes of the average, initial densities of positive- and negative-energy matter (and antimatter) were exactly the same, from the viewpoint of all observers, in the very first instants of the Big Bang, they could come to differ in a relatively large proportion later on, due to the potentially distinct violations of time-reversal symmetry which can be experienced by positive- and negative-energy matter. Indeed, it is under conditions where more matter than antimatter particles with a given sign of action (or vice versa) are produced, as a result of a violation of time-reversal symmetry, during the very first instants of the Big Bang (by processes I will describe in section 4.3) that some baryonic matter is allowed to survive the later processes of pair annihilation. But if the violation of time-reversal symmetry that gives rise to those violations of matter-antimatter symmetry is not as substantial for negative-action matter as it is for positive-action matter (which is always possible, as I have explained in section 3.9) then the magnitude of the density of negative-energy matter could become much smaller than that of positive-energy matter following the early annihilation of matter with antimatter, even if those two magnitudes were identical initially.

What's significant here is that the energy that is released in the course of those matter-antimatter annihilation processes is radiation energy, while the

density of this energy decreases more rapidly than that of matter energy as a result of expansion, which means that it does not contribute to decelerate the rate of expansion experienced by an observer with the same sign of energy as much as matter itself over time. Therefore, the rates of expansion experienced by positive- and negative-energy observers may be influenced differently by the presence of matter, during the matter-dominated era, even if the magnitudes of the densities of matter with opposite energy signs were exactly the same initially, from the viewpoint of any observer (which is unavoidable, as I will explain in section 4.5). As a consequence, the deceleration of the rate of expansion experienced by positive-energy observers could have been almost exactly the same as that experienced by negative-energy observers in the radiation-dominated era (while the cosmological constant would have had a null value) and the rate of expansion determined later on (during the matter-dominated era) by the same positive-energy observers could have become smaller than that which is determined by negative-energy observers, if the contribution of positive matter energy to the deceleration of the rate of expansion experienced by positive-energy observers had become larger than the contribution of negative matter energy to the deceleration of the rate of expansion experienced by negative-energy observers. This is what one can actually expect to happen whenever significantly more positive- than negative-action matter is allowed to survive the early annihilation of matter with antimatter. But under such conditions it can be expected that the cosmological constant would grow, from its initial zero value toward larger positive values, during the matter-dominated era³.

Now, astronomical observations do appear to show that the average density of baryonic negative-energy matter measured by a positive-energy observer is now much smaller than the average density of baryonic positive-energy matter measured by the same kind of observer (for reasons I will discuss in section 4.3), and this means that the expansion of the universe must have been slowed down to a greater extent during the matter-dominated era, from the viewpoint of a positive-energy observer, compared with what

³It is important to note that the initial annihilation of baryons with antibaryons would not make the rates of expansion experienced by positive- and negative-energy observers to vary much during the radiation-dominated era, because the rate of expansion was then determined by the density of radiation and all particles are relativistic during the radiation-dominated era, so that if there is no overall change to the energy density of matter plus radiation, then the rates of expansion should continue to decelerate as if all the matter was still present.

a negative-energy observer experienced. The only conclusion we can draw, therefore, is that the cosmological constant did grow to a larger positive value during a certain portion of the matter-dominated era, despite the fact that a positive, average value of vacuum energy density would contribute to accelerate the rate of expansion determined by positive-energy observers and to decelerate that which is determined by negative-energy observers, thereby moderating its own growth rate.

Thus, it is not impossible for the average density of vacuum energy to become dominant over that of positive matter energy, despite the fact that an average positive density of vacuum energy has a tendency to accelerate the rate at which matter is expanding from the viewpoint of a positive-energy observer, so that if it does become dominant, then it should reverse the very tendency that allows it to grow. Even though the average positive density of vacuum energy was smaller than the average density of positive matter energy until relatively recently, it must have grown in magnitude while matter energy was dominant over the uniform portion of vacuum energy. But if there was reason to expect that vacuum energy would eventually become dominant over positive matter energy on the global scale, it is not merely due to the fact that the cosmological constant was growing, but also because, the average matter density was decreasing as a result of the expansion of space. Therefore, despite the fact that the average density of vacuum energy must have been null in the very first instants of the Big Bang, it was allowed to grow into a positive value larger than that of positive matter energy (from the viewpoint of a positive-energy observer), at which point it finally began to accelerate the rate of expansion experienced by positive-energy observers and to diminish its own magnitude. Those deductions would appear to agree with astronomical observations, which indicate that the dominance of the average density of vacuum energy over that of matter energy occurred only recently on the cosmic time scale, given that the expansion of space is observed to be decelerating immediately before it began accelerating, in the most recent period of the universe's history.

What I originally failed to understand is that the difficulty one faces while trying to explain how it could be that the average positive density of vacuum energy was allowed to become dominant over the average density of positive matter energy can be avoided in the context where the early annihilation of matter with antimatter allows the divergence between the scale factors experienced by opposite-energy observers to rise continuously during the matter-dominated era, despite the fact that the positive value

of average vacuum energy density so produced contributes to accelerate the expansion rate experienced by positive-energy observers and to decelerate that which is experienced by negative-energy observers. Therefore, the fact that it wouldn't be appropriate to assume that the portion of zero point vacuum fluctuations which produces a maximum negative contribution to the energy of the vacuum cannot directly interact (other than gravitationally) with matter of positive energy sign, even though the presence of such matter is assumed to be equivalent to an absence of energy in this very portion of vacuum fluctuations, constitute sufficiently strong a motive to conclude that the final form of the gravitational field equations, from which is derived the vacuum-energy term provided by equation (4.1) above, is the one that must be retained. Thus, contrarily to what I had originally envisaged, the constraint that allows to decide whether the maximum value of the density of energy associated with those vacuum fluctuations that interact with positive-energy matter is positive or negative is not *purely* empirical, but also constitutes an unavoidable consistency requirement.

To avoid confusion, however, it must be understood that a positive cosmological constant contributes to accelerate the rate of expansion of space measured by a positive-energy observer and not merely to accelerate the rate of expansion of positive-energy matter, because the same metric conversion factor that is involved in determining the net value of vacuum energy density also affects the measure of negative-energy matter density determined by a positive-energy observer, as is made perfectly clear in the formulation of the generalized gravitational field equations introduced in section 2.15. Thus, what we may call the *specific density* of negative-energy matter (that which is measured by a negative-energy observer) is allowed to become smaller or larger than the specific density of positive-energy matter (that which is measured by a positive-energy observer), even in a universe in which the densities of positive- and negative-energy matter are of equal magnitudes initially (from the viewpoint of any observer), because, despite the fact that any such divergence would increase the magnitude of the average density of vacuum energy in such a way that the resulting cosmological constant would have a tendency to reduce this very divergence, the magnitudes of those specific densities can still be made to differ as a result of the early annihilation of matter with antimatter. However, the presence of the metric conversion factor in the second term of the decomposed generalized gravitational field equations associated with a positive-energy observer produces the same reducing effect on measures of negative-energy matter density as applies to the

positive instance of natural vacuum-stress-energy tensor that gives rise to a net positive value for the energy density of the vacuum.

As a result, even if the average, *specific* density of negative-energy matter had grown comparatively larger than that of positive-energy matter before the early annihilation of most baryons with their antibaryon counterparts, the average density of negative-energy matter which enters the gravitational field equations associated with a positive-energy observer would have remained as similar as it originally was to the specific density of positive-energy matter (that which is measured by a positive-energy observer) right until the annihilation process began. Of course, a similar effect would then occur for the measures of average positive-energy matter density entering the gravitational field equations associated with a negative-energy observer, because, despite the fact that the average, *specific* density of positive-energy matter would have become comparatively smaller than that of negative-energy matter during that very short period, the average density of positive-energy matter that is physically significant for a negative-energy observer would then actually grow in comparison with that measured by a positive-energy observer, along with the average, specific density of negative-energy matter, as a consequence of the presence, in the gravitational field equations, of the metric conversion factor associated with a negative-energy observer, which must give rise to the same unique cosmological constant (so that it must have an effect opposite that which arises from the metric conversion factor associated with a positive-energy observer).

To return to the analogy of the two embedded bi-dimensional spherical surfaces representing the spatial volumes of a closed universe which are experienced by opposite-energy observers, we may determine (through cosmological observations) the average density of negative-energy matter on the smaller surface associated with positive-energy observers in a universe with a positive cosmological constant, in order to predict the future evolution of the distribution of negative-energy matter. But, in doing so, we would have to take into account the fact that the surface on which the negative-energy particles actually evolve has a larger area, so that the distribution of negative-energy matter would appear to be deflated as it is projected on the surface over which positive-energy particles evolve. The average density of negative-energy matter which would be ‘observed’ on that surface would therefore be higher than the ‘real’ density which would be determined based on measures of distances associated with the larger surface on which negative-energy particles evolve. As a consequence, the ratio of the average

density of negative-energy matter to that of positive-energy matter obtained while using the measures of area associated with the smaller surface would remain identical to what it was initially, when the two surfaces had nearly exactly the same minimum area, but again, only as long as the magnitudes of the densities of positive- and negative-energy matter are not made to differ by the annihilation of matter with antimatter that takes place before the end of the radiation-dominated era. This, I believe, is the true significance of the transformation that is accomplished when one considers the stress-energy tensor of negative-energy matter in the form under which it is combined with the appropriate metric conversion factor in the generalized gravitational field equations from section 2.15.

If this interpretation is correct, it would mean that the average density of negative-energy matter over which are measured any density perturbations which may potentially affect the gravitational dynamics of positive-energy matter is not the specific density of negative-energy matter which is measured by negative-energy observers, but a measure of matter density dependent on the metric properties of space specific to positive-energy observers and which varies as a function of the rate of expansion measured by such observers. Thus, the variation of the average density of negative-energy matter which takes place either before or after the early phase of matter-antimatter annihilation is always assessed by a positive-energy observer based on the rate of expansion of space related to his own measures of distance and duration, which on a cosmological scale are influenced only by the average densities of positive-energy matter and vacuum energy and the same is true for the density of positive-energy matter measured by a negative-energy observer. This is why the ratio of the average cosmic densities of positive- and negative-energy matter must be considered an invariant quantity that is not affected by the actual value of the cosmological constant (even though it may still vary for independent reasons during the early phase of matter-antimatter annihilation).

There is no *a priori* motive, therefore, to assume that if space is expanding at a certain rate from the viewpoint of a positive-energy observer, then it should expand at the same rate from the viewpoint of a negative-energy observer, even during the radiation-dominated era (before the early annihilation of matter with antimatter had a sizable effect on the rates of expansion). It remains, however, that during the first instants of the Big Bang the average densities of positive- and negative-energy matter can be expected to have had exactly the same magnitude, so that the early rates of

expansion measured by positive- and negative-energy observers should themselves correspond, to an arbitrarily high degree of precision, as I will explain in section 4.5. Anyhow, it transpires that even the positive cosmological constant must affect positive- and negative-energy matter in the same way from the viewpoint of a positive-energy observer, because any acceleration or deceleration of the rate of expansion would depend merely on the metric properties of space associated with the gravitational field that this positive-energy observer experiences, despite the fact that the same average density of energy of zero-point vacuum fluctuations would influence the rate of expansion of matter in a different way from the viewpoint of a negative-energy observer. On the cosmological scale, the rate of expansion does not differ depending on the sign of energy of the expanding matter, but depending on the sign of energy of the observer who measures the expansion.

It must be emphasized again that the rule invoked above for justifying that the maximum positive contribution to the average density of vacuum energy is predominant when the scale factor determined by negative-energy observers is larger than that which is determined by positive-energy observers, simply follows from the fact that in such a case the metric conversion factor associated with the measurements of negative-energy matter densities effected by a positive-energy observer transforms the magnitude of the specific density of negative-energy matter (measured by a negative-energy observer) to a larger value, while the magnitude of the density of the maximum positive contribution to the energy of vacuum fluctuations that can be directly measured by a negative energy observer is an invariant quantity, so that when it is submitted to the same metric conversion as apply to negative matter energy, it would appear to be increased in comparison with the magnitude of the density of the maximum negative contribution to the energy of vacuum fluctuations that can be directly measured by a positive energy observer, thereby giving rise to a positive cosmological constant. It should be clear, however, that it is really the *specific* value of negative-energy matter density measured by a negative-energy observer that is transformed by the metric conversion factor which enters the gravitational field equations associated with a positive-energy observer and not the measure of negative stress-energy that is observationally determined (through its gravitational influence) by a positive-energy observer.

If such a transformation is necessary, it is merely as a consequence of the impossibility to directly compare the average density of matter energy observed by a negative-energy observer, or the average energy density of

those vacuum fluctuations which are directly experienced by such an observer, with the average density of matter energy experienced by a positive-energy observer, or with the average density of energy of those vacuum fluctuations which are directly experienced only by a positive-energy observer, due to the fact that it is not possible to directly compare the measures of spatial volume effected by such opposite-energy observers on a cosmological scale. This is made unavoidable not just by the absence of direct interactions between positive and negative energy matter, but also as a consequence of the fact that the presence of a smooth distribution of negative-energy matter exerts no influence on the gravitational field experienced by a positive-energy observer. But there is no reason to expect that the density of vacuum energy itself cannot vary with position, because it remains that the metric conversion factors were defined as locally-variable parameters and if that is allowed, then there is no *a priori* motive to assume that variations of vacuum energy density cannot occur, above those directly associated with the presence of ordinary matter itself (as voids in the homogeneous distribution of vacuum energy). In the following section I will explain what the freedom for the vacuum-energy term to vary as a function of position really means and how this property actually becomes an advantage of the particular interpretation of dark energy developed above.

It is also important to mention that when it is recognized that all positive contributions to vacuum energy must have a negative counterpart of equal magnitude, the whole notion of false vacuum with a larger than usual energy density becomes somewhat irrelevant, at least from a gravitational viewpoint, given that, under such circumstances, a non-zero cosmological constant can only arise when there exists a difference between the metric properties of space perceived by observers with opposite energy signs and not as a consequence of the actual nature of the processes taking place in the vacuum. Thus, when we say that a symmetry is broken in a low-energy vacuum state, what we should really mean is that the matter particles in this vacuum interact in a manner that is different from that by which the same particles interact when they are cooled in a different way in the same vacuum, or by which they interact at higher energies. But that does not mean that the vacuum itself is physically different, in particular with regards to its energy content. Of course, given that I have described matter as being equivalent to missing vacuum energy, I must recognize that the fact that matter can behave in different ways depending on how a symmetry is broken may nevertheless justify that we refer to the products of such symmetry breakings as

consisting of different vacuums. In any case, I think that it would no longer be appropriate to argue that, as ordinary baryonic matter contributes less than 5 percent to the average positive density of energy, then 95 percent of all matter must be considered of unknown nature, because if dark energy, which comprises more than 70 percent of the density of positive energy, really is vacuum energy, then a significant portion of it would consist in the exact same matter particles continuously fluctuating in and out of existence in their virtual form.

Now, if one demands an explanation for the smallness of the cosmological constant in the context of the above description of its origin, one would have to explain why it is that the scale factors and the rates of expansion experienced by observers with opposite energy signs (which we may call the *specific expansion rates* of positive- and negative-energy matter) were so similar in the very first instants of the Big Bang that they only began to differ significantly during the matter-dominated era, as a result of the early annihilation of matter and antimatter (which appears to have been more complete for negative-action matter). Indeed, despite the fact that a smaller specific rate of expansion of positive-energy matter would produce a larger positive cosmological constant, which would accelerate the rate of expansion of space measured by a positive-energy observer, thereby reducing the difference between this expansion rate and the specific expansion rate of negative-energy matter, which would eventually allow to reduce the magnitude of the cosmological constant, there is no doubt that the average value of vacuum energy density could have been much larger (into positive or negative territory) initially, even before the early annihilation of baryons and antibaryons had an effect on the specific rates of expansion of positive- and negative-energy matter, in which case its present magnitude would still be much larger than the measured value.

Here it would appear that one may have no choice but to invoke the weak anthropic principle, because it is not sufficient to recognize that the magnitude of the cosmological constant must have been reduced to a certain extent as a result of the negative feedback exerted by the average density of vacuum energy on itself. Indeed, according to Steven Weinberg [28], the current value of the cosmological constant is so close to the maximum limit imposed by the anthropic principle that it would appear that, if it is not much larger, this may simply be a consequence of the fact that a larger value would be incompatible with the existence of an observer. What I will explain in section 4.5 is that, in the context where we impose a requirement of null

energy on the universe as a whole, it becomes possible to assume that it is really anthropic selection which, alone, requires that the density of vacuum energy be as small as it is currently observed to be.

In any case, the validity of the approach advocated here is not compromised by the fact that it once seemed that empirical data was perhaps favorable to the hypothesis that the cosmological constant has not changed much during the recent history of the universe, because, given the current smallness of the observed average value of vacuum energy density (compared to the natural scale of quantum-gravitational phenomena), it is natural to expect that the rate of change of the cosmological constant, which during the most recent portion of the history of the universe was determined by the very magnitude of this density of vacuum energy, would have remained too small to be detected. But even under such conditions, it is certainly fortunate that it would appear that the effect of a non-zero cosmological constant is not to produce an even larger average (positive or negative) density of vacuum energy, as would have been the case if the vacuum-energy term that enters the generalized gravitational field equations had been that which is provided by the alternative equation (4.2) above, because, in such a case, one would have to conclude that despite the fact that the cosmological constant is still relatively small, it should eventually become much larger, while the specific rate of expansion of positive-energy matter should accelerate ever more rapidly, as a result of its own growth. Yet, I cannot be considered guilty of having chosen the right form of the vacuum-energy term based on a desire to avoid the prospect of predicting such an end for the world, given that, as I have explained above, I had originally assumed that the alternative form of the equations was actually the right one, on the basis of what appeared to be an unavoidable empirical constraint and despite the discomfiting outcome of such a choice. It is only in order to achieve greater consistency that I was later forced to recognize that this position is untenable.

To resume the situation, it transpires that the problem of the cosmological constant was complicated by the fact that it no longer appeared possible to explain the value of this parameter as being the outcome of a symmetry principle, when astronomical observations began to show that it is not exactly zero. This is because any violation of symmetry would likely produce a value of average vacuum energy density much closer to the natural scale of energy associated with quantum gravitation. What I have explained is that it is the necessary invariance under exchange of positive and negative energy states (which is justified by the requirement of relational definition of

physical attributes discussed in chapters 2 and 3) that allows one to expect a perfect cancellation of all contributions to the average density of vacuum energy in the absence of a divergence between the scale factors experienced by opposite-energy observers, while it is possible to assume (as I will explain in section 4.5) that it is the weak anthropic principle which, alone, explains that this divergence of the scale factors was not much larger than it could have been initially, thereby allowing the current value of the cosmological constant to be as small as it is observed to be. I believe that the fact that such a relatively simple and efficient solution to what has been called ‘the mother of all physics problem’ had never been seriously considered before is simply a consequence of the preconceived opinion that negative-energy matter cannot exist, which is a consequence of both irrational prejudice and what always appeared to be the insurmountable difficulties preventing a consistent description of gravitationally-repulsive matter.

4.3 Missing mass and dark matter

In this section I would like to discuss the impact of the developments introduced in the earlier portions of this report on our understanding of the phenomenon of missing mass⁴, which is currently believed to always arise solely from the presence of additional, unseen, but normally-gravitating positive-energy matter. What will emerge from those considerations is that additional effects, similar to those we normally attribute to ordinary dark matter, must be taken into account in the context of a cosmological model based on the generalized gravitation theory introduced in chapter 2. As a result, it is no longer necessary to assume that conventional dark matter is responsible for most of the missing-mass effect observed at the present epoch around visible positive-energy galaxies and clusters. Thus, while I will suggest that it is necessary to recognize the existence of an unexpected component of dark, but normally-gravitating baryonic matter, which could be responsible for a small portion of those missing-mass effects, I will also explain that, for the main

⁴It must be clear that what I’m referring to here is the general phenomenon that is usually attributed to the presence of dark matter and not that of voids in a matter distribution (even though I will suggest that those two phenomena may sometimes be related) and if I choose this slightly ambiguous and rarely used denomination it is because the problem I’m referring to is more general than the dark-matter problem itself, which merely consists in identifying a potential candidate for the weakly-interacting massive particles whose presence is usually assumed to explain this missing-mass effect.

part, the phenomenon of missing mass appears to merely be a secondary effect of the presence of energy attributable to zero-point vacuum fluctuations. Before delving into this important issue, however, I will explore another dimension of the dark-matter phenomenon which has been altogether ignored until now and which has to do with the gravitational attraction attributable to the presence of voids in an otherwise uniform matter distribution.

I have already mentioned in section 2.6 that certain forces which could not be distinguished from those traditionally attributed to positive-energy dark matter would arise in the presence of an underdensity in an otherwise uniform distribution of invisible negative-energy matter. This is because the absence of gravitational repulsion by the negative-energy matter that is missing, due to the presence of such an underdensity, would have the same effect on the surrounding positive-energy matter as would the presence of an overdensity of gravitationally attractive matter. If the interaction between positive- and negative-energy matter is governed by the principles enunciated in section 2.14, it would appear that such a phenomenon could in principle occur around positive-energy matter overdensities, given that such structures would repel negative-energy matter and thus create underdensities in this negative-energy matter distribution that could potentially enhance the gravitational attraction of the positive-energy objects, if the average density of negative-energy matter and the magnitude of the overdensities are large enough. In fact the same phenomenon should arise from the presence of voids in the positive-energy matter distribution, which can be expected to exert a gravitational attraction on any negative-energy matter (either visible or dark) that would be present around those voids.

What we can expect to occur, as we approach the center of mass of a sufficiently large overdensity in the positive-energy matter distribution, is that an increasingly smaller density of negative-energy matter would be present, because a larger fraction of it would not be able to overcome the repulsive gravitational force exerted by the overdensity. We can, therefore, expect the reduction in the density of negative-energy matter that is attributable to the gravitational repulsion exerted by the positive-energy matter overdensity to grow, along with the magnitude of the missing mass effect, as we approach the center of the structure. But, clearly, this cannot continue indefinitely, because the average density of negative-energy matter over which the underdensity is measured has a finite magnitude, which, at the present epoch at least, is much smaller than the average density of positive-energy matter in a typical galaxy. When the point is reached at which the magnitude of

the underdensity of negative-energy matter attributable to the gravitational repulsion of the positive-energy matter overdensity equals the magnitude of the average density of negative-energy matter itself, it becomes impossible to further reduce this matter density. This marks the limit beyond which the density of missing mass attributable to the presence of such an underdensity can no longer grow and actually becomes insignificant in comparison with the growing density of positive-energy matter within the structure.

In such a context, it should be clear that it is not possible to conclude that a potential contribution, by negative-energy matter underdensities, to the missing-mass effect around visible positive-energy structures could make contributions of a distinct nature unnecessary, because, for this to happen, the current magnitude of the average density of negative-energy matter (visible and dark) would need to be much larger than the currently inferred average density of positive-energy matter (both visible and dark), so that underdensities of sufficiently high magnitude could exist that would explain all of the missing-mass effects presently attributed to positive-energy dark matter. If the current average density of negative-energy matter determined by a positive-energy observer is merely as large as that of positive-energy matter, or if it is actually even smaller, as would appear to be required by observations (I will return to this question later in this section), then it simply isn't large enough to allow a replication of all the missing-mass effects around visible structures, which are known to involve equivalent matter densities hundreds of times larger than the average density of ordinary baryonic matter.

What's important to understand is that the amount of missing mass that can be attributed to the presence of underdensities in the negative-energy matter distribution is limited, at the present epoch, due to the fact that the current, average cosmic density of negative-energy matter is finite and relatively small compared to the density of matter inside most visible structures. The presence of negative-energy matter underdensities can therefore be expected to have accelerated the process of structure formation in the positive-energy matter distribution only at the epoch, in the remote past, when the average matter density was still relatively large and the matter distribution homogeneous enough on the scale of the structures considered. Indeed, any hypothetical missing-mass effect that could be attributed to the presence of underdensities in the negative-energy matter distribution can only be concentrated around positive-energy structures if negative-energy matter is otherwise smoothly distributed around those structures. But the

problem we face when we try to attribute to the presence of negative-energy matter underdensities most of the missing-mass effect around positive-energy galaxies is that, even though negative-energy matter may be more homogeneously distributed than positive-energy matter, at the present epoch, on the scale of galaxies and clusters, the average density of negative-energy matter is presently much smaller than the density of matter inside those visible structures.

Yet the possibility that negative-energy matter underdensities could have exerted an influence on the gravitational dynamics of positive-energy matter in the early universe, even on the scale of individual stars and galaxies, is real and certainly not undesirable, given that, despite all the progress which was achieved in the last decades to model the formation of large-scale structures, the currently favored theory of structure formation, involving only positive-energy cold dark matter, is still inadequate in certain respects. It is well-known, in effect, that the most recent observations (see in particular Ref. [29]) have kept revealing the presence of well-developed galaxies with masses much larger than expected, at increasingly larger redshifts, corresponding to an epoch when there shouldn't be any such galaxies according to current models. It is my belief that those difficulties will be alleviated once we recognize that, in the remote past, there existed significant contributions to the missing-mass effect which arose from the presence of underdensities in a relatively uniform distribution of negative-energy matter, which could have been produced as a result of the gravitational repulsion of positive-energy matter overdensities.

Indeed, even though we may have reasons to expect that no significant amount of baryonic negative-energy matter survived the early period of matter-antimatter annihilation, there are also good reasons to believe that negative-energy dark matter (dark from the viewpoint of observers made of baryonic negative-energy matter) was present in the primordial universe, with an average density comparable to that of non-baryonic, positive-energy dark matter and it will become clear later in this chapter why most of this matter cannot have been submitted to matter-antimatter annihilation like baryonic matter and must still be present in the universe today. Therefore, the gravitational attraction attributable to the presence of negative-energy matter underdensities could have played an important role (which need not be attributed only to positive-energy cold dark matter) in the formation of the primordial inhomogeneities that gave rise to visible, present day structures.

It is necessary to assume, in effect, that when the distribution of negative-energy matter was still sufficiently dense and uniform, as must have been the case in the primordial universe (for reasons I will explain later in this chapter), the gravitational repulsion of the structures which developed in the positive-energy matter distribution should have triggered the formation of negative-energy matter underdensities concentrated mostly around those developing structures, thereby allowing them to develop more rapidly. What cannot be assumed, however, is that negative-energy matter overdensities were present as well, at the same epoch, on a similar scale, which could have produced stellar- or galactic-size underdensities in the positive-energy matter distribution, that would have similarly accelerated the growth of those negative-energy structures. This is a hypothesis that is both theoretically unnecessary and observationally doubtful, because the presence of overdense negative-energy objects on a galactic scale should have exerted recognizable effects that would have been revealed already, by weak gravitational lensing experiments. But given that I will argue, later in this section, that dark matter overdensities have a tendency to form and to grow where baryonic matter overdensities with the same sign of energy are located, while baryonic negative-energy matter appears to be virtually absent at the present epoch in our universe, then those observations are quite understandable.

But it must be clear that even if a small portion of the missing-mass effect observed around present-day structures could still be attributed to the presence of negative-energy matter underdensities (particularly on larger scales), those contributions would not allow the average density of positive energy in our universe to reach its critical value, because any contributions to the energy budget from inhomogeneities in the negative-energy matter distribution would cancel out on a global scale, if those inhomogeneities developed in an originally smooth distribution of negative-energy matter (which I will argue to be a necessary assumption in section 4.9), given that there would then be as much spread out overdensities as there are localized underdensities in this matter distribution. Thus, even independently from any other considerations, it is necessary to recognize that the presence of negative-energy matter underdensities cannot contribute significantly to the observed missing-mass effect around positive-energy objects at the present epoch, given that the effect is already known to require the contribution of a density of gravitationally-attractive matter energy about as large as that which would bring the total density of positive energy to its theoretically and empirically required critical value.

Now, if one recognizes that the presence of negative-energy matter underdensities would never allow to explain a significant portion of the missing-mass effects which are observed around visible positive-energy structures at the present epoch, then one must admit that there definitely exist additional contributions of unknown origin to positive matter energy in our universe. Faced with the undeniable evidence that a certain form of dark matter must exist, the normal reaction is to seek to identify a weakly-interacting particle, different by necessity from all known particles, that might constitute a viable candidate for this dark matter. But for various reasons, despite the fact that all attempts at detecting and identifying such a particle have failed, it is still believed that dark matter should actually consist of particles that do not interact with ordinary matter *only* through the gravitational interaction. I believe that what really motivates this view is the fact that, if dark matter interacts with the rest of matter only through the very weak gravitational interaction, then it may, in effect, become impossible to determine the nature of those dark-matter particles by experimental means, which justifies that we concentrate, instead, on trying to identify a particle that does interact with ordinary matter through one of the other known forces. But what if we could deduce from certain observable properties of ordinary matter that there *must* exist positive-energy matter particles which can only interact with ordinary matter through gravitational forces?

At this point you may recall the discussion from section 2.9 concerning the possibility I had once contemplated that a certain condition of continuity along the world-lines of elementary particles could perhaps explain the empirically motivated requirement that a particle always reverses its energy sign when it reverses its direction of propagation in time. The idea was that if one allows action to reverse along a particle's trajectory in spacetime, this may give rise to some discontinuity that would be forbidden from a gravitational viewpoint. I eventually recognized, however, that continuity cannot *alone* be invoked for requiring that the action sign of a particle remains invariant under such circumstances and that what forbids the creation and the annihilation of particles with opposite action signs is the absence of any direct interactions between opposite-action particles (which also implies that no interaction boson can decay into opposite-action particles). Indeed, while the observed invariance of the sign of charge under a reversal of the direction of propagation in time is the decisive property that allows the time-direction degree of freedom to be physically significant (from the viewpoint of unidirectional time), the fact that a reversal of the direction of propagation in

time is always accompanied by a reversal of the sign of energy which leaves the sign of action invariant would appear to indicate that it is the sign of non-gravitational charges alone that must remain unchanged as a particle reverses its directions of propagation in space and time.

It became perfectly clear eventually that what a theoretically well-founded condition of continuity of the flow of time requires is merely that there be continuity in the true direction of propagation in time along an elementary particle world-line in spacetime. This restriction becomes relevant in the context where it is recognized that there does exist a fundamental time-direction degree of freedom distinct from the observed direction of motion of elementary particles. Compliance with such a continuity requirement would imply that particle-antiparticle creation and annihilation processes can only occur as the kind of events during which a particle bifurcates in spacetime to start propagating in the opposite direction of time and not as a chance encounter of two opposite-charge particles propagating in the same direction of time. This requirement would then also impose that events cannot occur which, from the unidirectional-time viewpoint, would appear to involve a particle that propagated a given charge forward in time turning into an identical particle that would now propagate an opposite charge backward in time (without an annihilation process taking place), because such processes would imply that the continuous path of a particle in spacetime (the arrow along a particle world-line) could abruptly reverse when, by chance (not as a result of any causal influence), a particle would meet an oppositely charged version of the exact same particle that was propagating backward in time.

Yet we have no choice but to assume that *ordinary* antiparticles (those that routinely take part in interactions involving ordinary matter) are indeed backward-in-time-propagating particles (and not particles propagating opposite charges forward in time), because, as I mentioned in the discussion concerning the time-direction degree of freedom appearing in section 2.2, if we are to view any transformation along a particle world-line as a continuous process, then ordinary antiparticles must always be considered to propagate in the direction of time opposite that in which the corresponding particles are propagating, given that the annihilation of an ordinary particle with an ordinary antiparticle must be allowed to occur with the same probability for all such pairs and cannot only take place for those pairs where the two particles would happen to be propagating in opposite directions of time. Indeed, when a condition of continuity of the flow of time applies along the world-lines of elementary particles, if certain *ordinary* electrons are allowed

to propagate negative charges forward in time, while other *ordinary* electrons would propagate positive charges backward in time, then certain electrons could not annihilate with certain positrons (those that would propagate an opposite charge in the same direction of time) with which they would nevertheless be allowed to interact, while it is known experimentally that no such a restriction to electron-positron annihilation exists (all known electrons can annihilate with all known positrons).

Thus, even if some of the electrons that propagate in a particular direction of time could have a negative charge, while others would have a positive charge, we must consider as empirically forbidden for particles with such opposite *bidirectional charges* (the invariant measures of charge which are not affected by a conventional reversal of the direction of propagation in time of elementary particles) to transform into one another or to interact with one another, at least under ordinary circumstances. No particle of any given kind that propagates some non-gravitational charge forward or backward in time can decay into, or interact with a particle of the same kind that would propagate an opposite charge either forward or backward in time and which would otherwise appear to consist of the exact same kind of particle (two such particles propagating in opposite directions of time, would not merely consist of a particle and its antiparticle, they would actually have the same sign of energy and the same sign of charge, from the viewpoint of unidirectional time).

What I'm suggesting, therefore, is that the fact that particle-antiparticle annihilation processes are allowed to occur for any particle-antiparticle pair does not mean that there can be discontinuities in the direction of the flow of time along a particle world-line in spacetime, but really that particles with opposite bidirectional charges cannot interact with one another or transform into one another, under ordinary circumstances. Thus, I propose that we recognize the existence of a fundamental rule which can be formally expressed using the following definition⁵:

Condition of continuity of the flow of time: There must always be a continuous flow in the true direction of propagation

⁵It will be made clear in the latter portion of chapter 5 that what justifies this rule, from a fundamental viewpoint, is really the principle of local causality as it applies to particle propagation processes from the viewpoint of a time-symmetric quantum-mechanical description of reality, in the context where we require all causes to originate from within the universe in which the processes involved take place.

in time along a particle world-line in spacetime for elementary fermions, even when particle-antiparticle creation and annihilation processes are involved from a unidirectional-time viewpoint.

I will soon explain what justifies (from a theoretical viewpoint) the validity of the empirical rule that particles propagating a given charge forward or backward in time cannot interact (other than gravitationally) with similar particles propagating an opposite bidirectional charge forward or backward in time (even when it may appear that the condition of continuity of the flow of time would not be violated) and therefore cannot transform (under ordinary circumstances) into such particles either, even in the course of particle-antiparticle creation and annihilation processes.

It should already be clear, however, that even when the proposed constraint applies, the charge of a particle (not necessarily the electric charge, but any non-gravitational charge) can still vary on a continuous world-line (as when a blue quark turns into a red quark, or a neutrino turns into an electron), because all that is required is that a particle with a given charge does not change into an identical particle with an opposite *bidirectional* charge (the measure of charge which is independent from the direction of propagation in time) along such a continuous path, particularly in the case of fermions, so that if a particle was initially propagating a given charge forward in time, it can still be assumed to propagate the same charge in the same direction of time, unless a particle-antiparticle annihilation process takes place and the same charge begins propagating backward in time. Of course, a similar conclusion would apply for a particle propagating a given charge backward in time, which must continue to propagate the same charge in the same direction of time unless a particle-antiparticle *creation* process takes place and the same charge begins propagating forward in time.

The difficulty, here, consists in recognizing that there are actually very good reasons to assume that, in the context where an antiparticle must be considered to be an ordinary particle propagating the same charge backward in time with reversed energy, a condition of continuity of the flow of time must be imposed. What we do, from a conventional viewpoint, is that we simply ignore the possibility that an electron, for example, may exist that would propagate a positive charge and a negative energy backward in time, by assuming, as a matter of coordinative definition, that a positive-action electron always propagates a negative charge forward in time, while a positive-action positron always propagates the same negative charge backward in time, as if

there were no other possibilities (which makes the issue of continuity irrelevant). This is similar to what we do when we exclude negative energy states propagating forward in time, or positive energy states propagating backward in time, by assuming that they are nonphysical states. In the present case, however, it is not even understood that in doing so we are deliberately choosing to exclude certain states of matter from our description of reality, because it looks like all that is involved is a definition. But that is not the case, and if the choice of which positive-action electrons propagate a negative charge forward in time and which propagate a positive charge (along with a negative energy) backward in time is, in effect, a simple matter of definition, the decision to exclude as nonphysical those electrons which, according to this definition, would propagate a positive charge forward or backward in time can only be justified on the basis of observational evidence.

One may argue that this distinction is insignificant, because the validity of the traditional approach is in fact empirically confirmed, given that it does provide a theoretical framework whose predictions agree perfectly well with observational constraints. Or does it? We still have a serious problem in theoretical cosmology, because we do not know what most of the matter in our universe is made of. Could it be that there is in fact something wrong with some of the implicit choices which were made a long time ago, while we were trying to make sense of the newly developed mathematical framework of relativistic quantum fields, before everybody even knew about the existence of dark matter? Is it possible that there does exist in our universe positive-action electrons with positive bidirectional charges and positive-action protons with negative bidirectional charges and that those particles actually constitute a non-negligible portion of the normally-gravitating dark matter, along with the positive-action neutrons composed of negatively-charged up quarks and positively-charged down quarks propagating forward or backward in time?

I do recognize that there may be serious difficulties with this idea, because, even if one acknowledges the fact that, from an empirical viewpoint, positively-charged electrons propagating either forward or backward in time should not be allowed to interact with, or to transform into ordinary electrons propagating negative electric charges in any direction of time, or to interact with ordinary protons propagating positive electric charges forward or backward in time, one still needs to explain what justifies this limitation from a theoretical perspective. What's more, even if we could justify the absence of interactions between ordinary matter particles and their dark-matter coun-

terparts with reversed bidirectional charges, then it would remain to explain why it is that those particles do interact through the gravitational interaction. I believe, however, that those difficulties do not decisively rule out the existence of such baryonic dark-matter particles and that it is possible to understand, by making use of the developments already introduced in this report, why reversed-bidirectional-charge particles should, in effect, be dark, despite the fact that they can also be expected to interact gravitationally with the rest of matter, thereby allowing them to contribute to the missing-mass effect around visible positive-energy structures.

What I have come to understand is that the difficulty one may face while trying to explain the absence of electromagnetic interactions between electrons with negative bidirectional charge and electrons with positive bidirectional charge arises merely because one ignores the fact that the previously defined constraint regarding the continuity of the flow of time along a particle world-line must also apply in the case of the particles that mediate the interactions between elementary particles of matter. The problem is that, according to the current interpretation, the world-lines of interaction bosons would appear to abruptly come to an end when they are absorbed by a matter particle, just like they would seem to come into existence discontinuously when they are emitted, either by a fermion or another interaction boson. While this may not appear to violate any principle, a certain tension clearly exists between the traditional description of those absorption and emission processes and the previously discussed constraint regarding the continuity of the flow of time along a particle world-line. But, instead of arguing indefinitely as to why such discontinuities are allowed to occur, despite the fact that they may be at odds with certain rules that seem to apply in the case of fermions, I would suggest that we simply assume that in fact the flow of time along the world-lines of elementary particles is never really interrupted, given that the bosons mediating the interactions between elementary particles of matter, somehow, allow charges to propagate along two opposite directions of time all at once, *as if* the spin-one interaction bosons were composite particles made of a fermion and an anti-fermion which need not carry the same charges.

One important characteristic of such an alternative description is that, if there must, in effect, be a continuity of the flow of time along the world-lines of all elementary particles, then, in the context where the interaction bosons would allow a propagation of charges along two opposite directions of time all at once, it follows that the direction of propagation in time of the interacting

particles would actually be allowed to remain unchanged during any such interaction process, because time flows in and out of the interaction boson at each vertex. This is, of course, in accordance with the previously stated conclusion to the effect that the condition of continuity of the flow of time along the world-lines of elementary particles forbids the transformation of a particle propagating a given charge forward in time into an apparently identical particle propagating an opposite charge backward in time and therefore it seems that it is really the necessary continuity of the flow of time that imposes that the interaction bosons be described as always propagating charges in two opposite directions of time all at once.

From the viewpoint of this equivalent description of interaction processes, it would follow that, for any interaction vertex, time would flow from the incoming fermion into the interaction boson and from the interaction boson into the outgoing fermion (or from the outgoing fermion into the interaction boson and from the boson into the incoming fermion if this particle is propagating backward in time) and the same must be happening at the other vertex of an interaction diagram. An examination of the diagrams describing the interactions between elementary particles, such as those represented in figures 4.1 and 4.2, clearly shows that this hypothesis agrees with the description of all known interaction processes, even those that involve a variation in the charges of the interacting matter particles that must be carried by the interaction bosons, but only when we assume that the bidirectional charge signs of the particles involved (those which are attributed to matter particles which are propagating forward in time, when they are observed from the unidirectional-time viewpoint) must be either both non-reversed or both reversed, regardless of whether the interacting particles are propagating forward or backward in time.

Indeed, while the above defined condition is satisfied for those processes where both of the interacting particles are propagating a certain bidirectional charge with a unique given sign either forward or backward in time, it cannot occur for the same processes where only one of the interacting particles is propagating a reversed bidirectional charge (in any direction of time). In the latter case, either the direction of propagation in time would be reversed along the direction in which the charge is flowing (from a unidirectional-time viewpoint), or else a bidirectional charge would have to transform into an opposite bidirectional charge, in the direction along which time is flowing (from a bidirectional time viewpoint). The crucial point is that, even when a neutral interaction boson is involved, bidirectional time must flow continuously

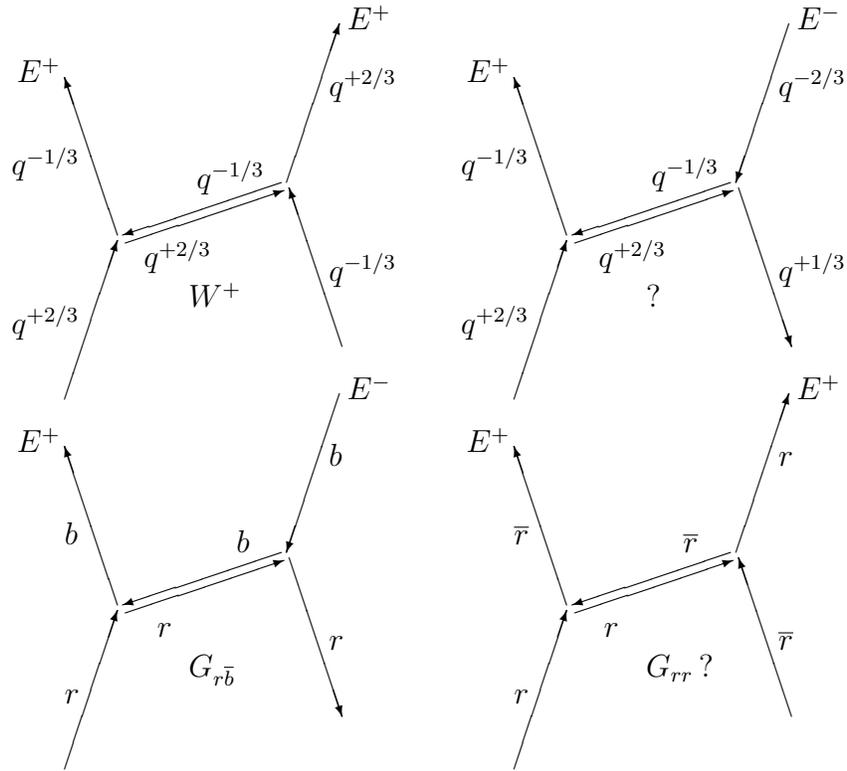


Figure 4.1: Alternative Feynman diagrams for flavor- and color-changing interactions between quarks. Here $q^{+2/3}$ and $q^{-1/3}$ represent the magnitudes and the signs of fractional electric charges, as determined from a bidirectional-time viewpoint, b and \bar{r} represent a quark color and anti-color, while E^+ and E^- are the energy signs relative to the direction of propagation in time, which corresponds to the direction of the arrows relative to the vertical axis. The diagrams on the left represent processes which are allowed to occur, while the diagrams on the right represent processes which are not allowed to occur based merely on the requirement that the bidirectional charge signs be left invariant along the direction of the flow of time.

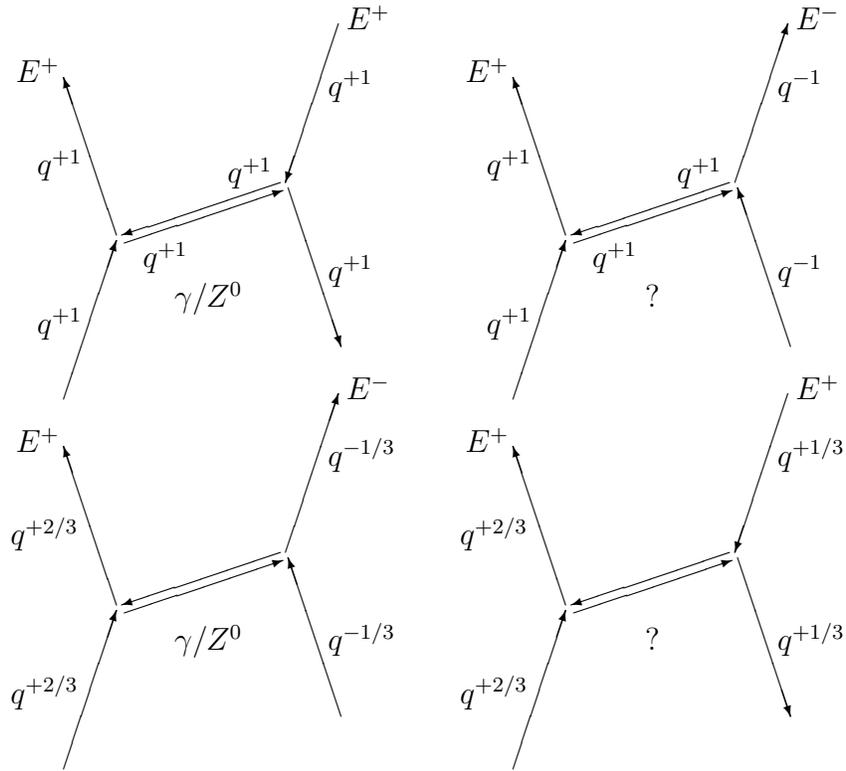


Figure 4.2: Alternative Feynman diagrams for flavor-conserving electroweak interactions between fermions. Here again q^{+1} and q^{-1} represent the magnitudes and the signs of the electric charges, as determined from a bidirectional-time viewpoint, while E^+ and E^- are the energy signs relative to the direction of propagation in time. It is only for processes of the kind described in the diagrams on the left that the sign of the bidirectional charge carried by the interacting matter particles does not vary discontinuously along the direction in which time is flowing and what is observed is that only processes of this kind actually occur in nature, even when there is no actual change in the charges of the interacting particles.

from the original interacting matter particle into the forward-propagating component of the boson and then into the ingoing or outgoing world-line of the other interacting matter particle (depending on its direction of propagation in time), despite the fact that it may appear like no specific charge is being propagated by the interaction boson.

What I'm suggesting is that it is the fact that a physical attribute of elementary particles, which is normally associated with the sign of their non-gravitational charges, would have to vary discontinuously along the direction in which time is flowing in the diagram describing an interaction among elementary particles (that which is indicated by the direction of the arrows) that explains that, from the viewpoint of the above description of interaction processes, the only interactions which are allowed to take place between two particles are those involving identical particles (one of which may be an antiparticle) with the same sign of bidirectional charge, or particles which both have either a reversed or a non-reversed sign of bidirectional charge, although this is only explicitly apparent in the case of an interaction during which there is an exchange of charge that is carried by the interaction boson. Thus, regardless of whether the charges of two interacting particles vary, it is only in those cases where the charges of both particles are either reversed or non-reversed that the interaction is actually allowed to occur.

What's important to understand is that, even though for certain interaction processes (such as that which is illustrated in the lower right diagram of Figure 4.2) we would not merely observe a reversal of charge, when we follow the direction of the flow of time along a particle world-line, or a reversal of time, when we follow the flow of charge (from a unidirectional-time viewpoint) during a hypothetical interaction process between an ordinary particle and a particle with reversed bidirectional charge, because the charges of the particles involved in the process (say an ordinary up quark and a down quark with reversed bidirectional charge) do not have the same magnitudes and do not vary as a result of the interaction (so that, from a traditional viewpoint, no charge would appear to be carried by the interaction boson), such an interaction is still forbidden by the condition of continuity of the flow of time. What I'm proposing, in effect, is that even when two interacting particles do not have the same magnitude of charge and those charges are not altered by the interaction, there still exists a constraint for some property of elementary particles, normally associated with their non-gravitational charge sign, to flow continuously, either forward or backward in time, from an interacting particle with non-reversed bidirectional charge into one that also has a non-

reversed bidirectional charge (or from an interacting particle with a reversed bidirectional charge into one that also has a reversed bidirectional charge), in order precisely that a discontinuity in the direction of the flow of time be avoided at the fundamental bidirectional level.

Thus, if an observable particle (with a non-reversed bidirectional charge sign) has a positive charge of $q = +2/3$ when it propagates forward in time, while another observable particle has a negative charge of $q = -1/3$ when it propagates forward in time, then it must be assumed that a certain attribute of the particles, normally associated with their non-gravitational charge sign, cannot flow continuously from the first particle, with a charge of $q = +2/3$, into an interaction boson and then back into a particle identical to the second one, but whose bidirectional charge $q = +1/3$ would be opposite that of its observable counterpart. What this means is that there is always a clear (even though relationally defined) distinction between what constitutes an ordinary matter particle and what constitutes a particle with reversed bidirectional charge, even when we are dealing with particles which do not carry the same magnitudes of charge, and this means that, even observable particles with different magnitudes and signs of charge must share a certain physical attribute which would only vary if the sign of the bidirectional charge carried by those particles (that which is independent of the direction in which a particle is propagating in time) was allowed to reverse, while it remains unaffected by a mere reversal of the direction of propagation in time that actually leaves the sign of charge invariant (from the viewpoint of bidirectional time). But I must acknowledge that we will probably only be able to fully understand what justifies the rule described here when we obtain a more complete theory of elementary particles which would allow a description of quarks and leptons (and perhaps also of interaction bosons) as composite particles.

In any case, if the constraint of continuity of the flow of time along a particle world-line extends to interaction bosons in the way suggested here, then it would appear that no interaction can occur that would involve two identical matter particles with opposite bidirectional charges (those observed while following the direction of propagation in time of the particles) propagating in any direction of time, or even merely two different particles, when only one of them has a non-reversed sign of bidirectional charge which is propagating in either the past or the future direction of time. I believe that this would be a simple consequence of the fact that no such an interaction could ever be described as a process during which the polarities of all *non-gravitational*

attributes of the elementary particles involved (and this excludes spin and the sign of energy) remain unchanged (are not subject to discontinuous reversal) as we follow the direction of the flow of time along their respective world-lines, from one of the two interacting matter particles into the interaction boson and then back into the other matter particle, either forward or backward in time.

What must be clear is that there is a difference between the description of an interaction process according to which none of two interacting particles has a reversed sign of bidirectional charge and the alternative description according to which one of the particles would have a reversed bidirectional charge sign. From my viewpoint this distinction is such that it forbids the processes from occurring when the sign of bidirectional charge of one of the particles is reversed (not necessarily in the course of the interaction), even if the direction in which the particle is propagating in time is also reversed. Thus, even if it may appear that a quantum-mechanically equivalent description of a certain interaction process could exist that would be obtained by simply reversing both the sign of charge and the direction of propagation in time for one of the interacting particles, we would have to conclude that the description for which the sign of charge is reversed, independently from the direction of its propagation in time, is actually distinct from that which does not involve such a reversal. This distinction would simply be a consequence of the fact that, while the sign of charge that is reversed in the apparently equivalent description of the process would not appear to be reversed from the viewpoint of unidirectional time, along which the particle is observed, from a bidirectional viewpoint this charge would nevertheless be reversed and this constitutes a physically significant change, because even when the two interacting particles do not have the same magnitude of charge (as an up quark and a down quark) and are not transformed by the interaction, they still interact only when they both have non-reversed bidirectional charge signs, or when they both have reversed bidirectional charge signs, given that the same distinction that explains the absence of flavor changing interactions that would violate the condition of continuity of the flow of time along the direction in which the charges are propagating also exists in such a case and must have similar consequences.

Therefore, I'm allowed to conclude that the rule which is implicitly assumed to apply, from a traditional viewpoint, to the effect that no positive-action particle that would be propagating charges opposite those of ordinary particles in the opposite direction of time need be considered to exist, is only

appropriate in the sense that it is not possible for any such particle to interact with ordinary particles, at least through the exchange of interaction bosons associated with non-gravitational forces. But it is also very clear that this does not mean that particles with reversed bidirectional charges cannot exist, because, from a theoretical viewpoint, this conclusion would be as unjustified as that which would amount to argue that ordinary particles themselves cannot exist. Indeed, the distinction between electrons propagating positive charges backward in time and ordinary electrons propagating negative charges forward in time is only a relational distinction, in the sense that a positively-charged electron propagating backward in time can only be distinguished from an ordinary electron through the fact that it actually has a charge that is opposite that of the ordinary electron, even while it propagates in a direction of time opposite that in which an ordinary electron propagates, but those are not absolutely characterized properties and an electron with positive bidirectional charge is only different from an electron with negative bidirectional charge in the exact same way an electron with negative bidirectional charge is different from an electron with positive bidirectional charge and it is not possible to distinguish one from the other, except through those mutual relationships. If there is no intrinsic or absolute distinction between particles in those two different states, however, then it means that none of them can be considered more real than the other. In other words, both kinds of particles must be assumed to exist, even though matter with reversed bidirectional charges must by necessity be dark from the viewpoint of ordinary matter.

Now, obviously, the only way that such a conclusion could come out as not totally meaningless is if the gravitational interaction is not affected by the condition of continuity of the flow of time along the world-lines of elementary particles, because otherwise there should be no interaction at all between ordinary positive-energy particles and positive-energy particles with reversed bidirectional charges. But I believe that this is actually unavoidable, because it is clear from the above discussion that it is merely the non-gravitational attributes of elementary particles that must not be subjected to any discontinuous reversal along their respective world-lines. The gravitational interaction is fundamentally distinct from all other interactions in this respect, given that it is neutral with respect to all non-gravitational charges, which is not really the case with other neutral interactions that couple to charge (even though they are mediated by interaction bosons that do not appear to carry a charge). This essential distinction, which is unique

to the gravitational interaction, appears to be what allows opposite bidirectional charge particles with the same energy sign to interact gravitationally (and attractively) with one another. The fact that gravitons couple only to energy, while the sign of energy or action is not affected by a reversal of bidirectional charge means that gravitation is the only truly neutral interaction, which therefore remains unaffected by the condition of continuity of the flow of time that prevents the existence of other interactions between opposite bidirectional charge particles⁶.

Anyhow, it seems that the only conclusion that can be drawn is that, despite the fact that positive-energy matter with reversed bidirectional charges is dark, it would actually exert attractive gravitational forces on ordinary positive-energy matter particles, as well as indirect repulsive gravitational forces on all negative-energy matter particles, regardless of their bidirectional charge signs. It is necessary to assume, in effect, that there also exist negative-energy (negative-action) particles with reversed bidirectional charges, which would be dark from the viewpoint of ‘ordinary’ negative-energy observers. But due to the requirement of symmetry under exchange of positive- and negative-energy matter, it must also be assumed that any such negative-energy particle is gravitationally attracted to other negative-energy matter particles and is gravitationally repelled by all positive-energy matter particles, regardless of their bidirectional charge signs. What allows the existence of repulsive gravitational interactions between negative-action particles with reversed bidirectional charges and visible positive-action particles is the fact that *all* negative-action particles are equivalent to the presence of voids in the positive-energy portion of the vacuum, while such voids *necessarily* exert indirect gravitational forces on positive-action particles.

To summarize what I have discussed so far, it seems that additional, attractive gravitational forces could, in principle, be exerted on positive-energy matter, as a consequence of the presence in the negative-energy matter distribution of underdensities attributable to the gravitational repulsion of positive-energy matter overdensities. We have no choice, however, but to

⁶The fact that gravitons have a spin that is twice as large as that of other interaction-mediating bosons may suggest that this neutrality is preserved as a result of the fact that the gravitational interaction allows charges to propagate along two opposite directions of time *in two different ways* all at once, for each individual interaction, thereby allowing the changes produced by one of those exchange process to be neutralized by that of the other, in observance of the condition of continuity of the flow of time, *as if* spin-two gravitons were composite particles made of two spin-one *bosons*.

assume that positive-energy dark matter is present in our universe in one form or another, because the magnitude of the underdensities that could exist today in the negative-energy matter distribution is much too small to explain a significant portion of the missing-mass effects observed around visible, present day structures, due to the fact that the average density of negative-energy matter in which such underdensities could develop is itself much too small.

Yet, I must acknowledge that it is not possible to conclude that the missing-mass effect is attributable, for the most part, to the presence of particles identical to those that compose visible matter, but which happen to propagate bidirectional charges opposite those propagated by ordinary matter and antimatter particles. Indeed, if most of the dark matter that is assumed to be responsible for the missing-mass effect was composed of particles which interact with themselves through the same forces by which ordinary baryonic matter particles interact, then it would be more difficult to explain the near spherical shape of dark-matter halos, or certain observations of colliding clusters of galaxies which show that, while the detectable high-energy gas originally present in the clusters is stripped of the galaxies as a result of such a collision, most of the dark matter is unaffected by the process. I initially thought that this difficulty may simply be a consequence of the fact that we ignore the possibility that baryonic dark matter could be more susceptible to collapse into stars and other high-density objects at a very early stage if its average density happens to be larger than that of visible baryonic matter (I will soon explain how I later came to realize that this is not possible and that it is necessary for the average density of normally gravitating baryonic dark matter to be equal that of ordinary baryonic matter). Under such conditions baryonic dark matter would no longer interact with itself on a larger scale, when galaxies begin to form later on, which could perhaps allow to explain the near spherical shape of dark-matter halos. For the same reason, it would have appeared appropriate to assume that the dark matter present inside colliding clusters is mostly unaffected by the collisions, just like the visible stars which are present within the galaxies, despite the fact that the dark matter particles that form those astronomical objects are allowed to interact electromagnetically among themselves.

For this to be a valid hypothesis, however, one would need to assume that a density of baryonic dark matter much larger than that of visible baryonic matter exists in the form of massive compact astronomical objects or MA-CHOs. But even though early studies seemed to indicate that the existence

of a large amount of matter in the form of invisible MACHOs was not completely ruled out, because what really motivated the commonly held opinion that there cannot exist enough MACHOs to provide a sizable portion of the dark matter was merely the impossibility for those objects to be formed of ordinary baryonic matter (whose presence would affect the predictions of Big Bang nucleosynthesis), more recent astronomical observations [30] do confirm that there cannot be such a large portion of normally-gravitating matter in the form of MACHOs (regardless of the nature of their constituent particles). Thus, it is no longer possible to assume that a sufficiently large number of such objects could exist that would be composed of baryonic matter with reversed bidirectional charges⁷ (which would not have affected the predictions of Big Bang nucleosynthesis). As a consequence, it is necessary to recognize that the above discussed difficulties associated with the hypothesis that dark-matter particles may interact with themselves (like ordinary baryonic matter) can only be surmounted if most of the observed missing-mass effect is attributable to a phenomenon distinct from those I have discussed so far. This doesn't mean that none of the dark matter can consist of baryonic matter with reversed bidirectional charges (this is not ruled out by the new observations), but merely that it is not possible to conclude that the necessary existence of such matter provides a valid explanation for most of the missing-mass effect observed around visible galaxies and clusters.

Now, even aside from the question of the origin of the missing mass effect, it would certainly be interesting to know what the average density of positive-energy matter with reversed bidirectional charges actually is in our universe. Does a theoretical constraint exist that would allow one to tell what this density actually is? In fact, it turns out that this value can be determined with great accuracy. To see how this is possible it will be necessary to reexamine the problem of matter-antimatter asymmetry in the light of the progress already achieved in this report. First of all, is it appropriate to conclude that the absence of antimatter in our universe reflects a fundamental lopsidedness with respect to the direction of time? This is an important question, because, while it is usually recognized that thermodynamic time asymmetry is probably not the cause of the violations of T symmetry which have been observed in certain high-energy experiments, the direction of time

⁷This is not to say that there cannot exist any large astronomical objects composed of reversed-bidirectional-charge matter, as it is quite possible, in fact, that invisible stars and planets made of such matter are present in our own region of the universe.

singled out by T violations can be related to the macroscopic arrow of time and this might allow one to conclude that our universe is characterized by a phenomenologically-apparent fundamental lopsidedness. Given that the time-reversal symmetry operation can now be understood to involve a transformation of matter into antimatter, the question of whether there actually exists such a preferred direction in time would, in effect, be equivalent to ask if there really is an absolutely definable asymmetry associated with the predominance of matter over antimatter in our universe?

I initially thought that what would allow symmetry with respect to the direction of time to be regained, despite the observed asymmetry between matter and antimatter, would be the existence of a certain constraint that would require that whenever there is an overabundance of forward-in-time propagating positive-action particles over positive-action particles propagating backward in time, there should be an overabundance of unobservable, negative-action particles propagating backward in time over forward-in-time propagating negative-action particles. It would then merely be the fact that the backward-propagating particles which do exist in large numbers cannot be observed, due to the fact that they would have an opposite sign of action and would therefore be dark, that would make it seem like there is a smaller number of such particles, compared to forward-propagating particles. But, as I realized later, this solution is not valid, because, as I explained in section 3.9, the violation of T or C symmetry that must be involved in giving rise to the observed positive-action, matter-antimatter asymmetry cannot be assumed to directly compensate that which would give rise to the similar, relationally-defined T or C asymmetry that *may* be affecting negative-action matter and antimatter, because T and C can be violated to a different degree for negative-action matter, as long as invariance under PTC is independently preserved by processes involving this type of matter.

Thus, it cannot simply be assumed that what happens is that the matter-antimatter asymmetry is reversed for negative-energy matter and that there is actually the same number of otherwise identical particles propagating the same sign of energy and the same sign of bidirectional charge in opposite directions of time, when we appropriately take into account the contribution of the unseen negative-energy matter. What I will explain in section 4.9 is that once it is recognized that there is a possibility for time to extend past the initial Big Bang singularity, following a hypothetical quantum bounce, then it is no longer necessary to conclude that there is an absolute lopsidedness, that one could attribute to the apparent asymmetry between matter and

antimatter. Yet, even under such conditions, it would seem that a certain asymmetry could persist if there were more matter particles of any given kind with a certain sign of bidirectional charge than there are matter particles of the same kind with an opposite bidirectional charge sign, either before or after the quantum bounce.

Here, I must first mention that it is well understood already that the difference between the density of matter and that of antimatter in the universe today is a consequence of the fact that, due to a certain violation of CP symmetry (which would imply a violation of time-reversal symmetry T), there was a little more matter than antimatter in the primordial Big Bang state, so that some matter was allowed to survive the annihilation of matter with antimatter that took place when the temperature became too low for ordinary pair creation to compensate the related process of pair annihilation. The problem is that it is not known what the exact origin of the asymmetry is that gave rise to this overabundance of matter over antimatter. Now, what I realized is that, given that particles with opposite bidirectional charges do interact gravitationally with one another, then one must conclude that, under conditions where gravitation is strong enough, as was the case in the first instants of the Big Bang (despite the very high homogeneity of the initial matter distribution), it should be possible for pairs of particles with opposite bidirectional charges propagating in the *same* direction of time to be produced out of gravitational radiation energy.

Indeed, while all sorts of particle-antiparticle pairs can be created in the first instants of the Big Bang by all sorts of interactions, most particle pairs which are produced in such a way can annihilate back to radiation and be recreated soon after, just to decay again into radiation, and so on, until relatively late times, as long as the interactions involved are strong enough for their effects to persist despite the reduction of temperature associated with expansion. But once a pair of particles was produced by the gravitational interaction, at the epoch in the far past when space was expanding at an arbitrarily large rate, the particles so produced were allowed to move away from one another rapidly enough that they were no longer able to annihilate back to gravitational radiation (given that, on the scale of quantum gravitational phenomena, the distance between the particles had become too large and the magnitude of their energies too low), which means that the creation process had become permanent. But if those particles had opposite bidirectional charges and if there are more matter particles than antiparticles created in such a way, this imbalance would itself become permanent,

because ordinary matter particles and particles with reversed bidirectional charges cannot interact other than gravitationally.

For such an imbalance to be initiated, a violation of T or C symmetry would need to arise which would affect the production of pairs of opposite-bidirectional-charge particles propagating forward in time in the exact opposite way it would affect the production of pairs of opposite-bidirectional-charge particles propagating backward in time, even though it would not affect the production of pairs of particles with the same bidirectional charge sign propagating in opposite directions of time (precisely because in such a case both directions of time are involved in every process). Thus, the first of the two following graviton decay processes (where $+2/3$ and $-2/3$ are the *bidirectional* electric charges of an up quark and the overline indicates that the quark is propagating backward in time) must have occurred a little more often than the second and the same must be true for all types of particles:

$$g \rightleftharpoons u^{+2/3} + u^{-2/3}$$

$$g \rightleftharpoons \bar{u}^{+2/3} + \bar{u}^{-2/3}$$

I believe that this is what explains that there actually existed a little more ordinary matter than ordinary antimatter in the first instants of the Big Bang, so that there are presently more negatively charged electrons propagating forward in time than there are negatively charged electrons propagating backward in time. But once opposite-bidirectional-charge particles are produced by such processes, then they no longer interact with one another through any non-gravitational interaction, as I mentioned above, which means that matter with reversed bidirectional charges becomes dark and appears to be missing from the viewpoint of observers made of matter with non-reversed bidirectional charge signs.

In such a context it is possible for an exact compensation of the observed asymmetry between the number of electrons with negative bidirectional charge and that of electrons with positive bidirectional charge to arise that would not be immediately apparent, because, while it is not required that the exact same numbers of ordinary electrons and positrons be produced by the processes described above, it is still necessary that the exact same numbers of particles propagating opposite bidirectional charges in the same direction of time be produced by those pair creation events. As a result, if there is an overabundance of protons propagating a positive bidirectional charge forward in time over protons propagating the same charge

backward in time, from the viewpoint of an observer made of such protons, then there should exist an overabundance of protons propagating a negative bidirectional charge forward in time over antiprotons propagating the same bidirectional charge backward in time, from the viewpoint of an observer made of such negative-bidirectional-charge protons, which means that no asymmetry actually exists that would have to do with an overabundance of particles with one bidirectional charge sign over those with an opposite bidirectional charge.

Given that it is possible for the symmetry under a reversal of the direction of time to be violated, even when time direction remains a relationally defined parameter, and given that the gravitational interaction allows opposite-bidirectional charge particles to interact without violating the constraint imposed by the requirement of continuity of the flow of time, I must then conclude that nothing would forbid a violation of matter-antimatter asymmetry from developing as a result of more forward- than backward-in-time-propagating pairs being created (or vice versa). Thus, contrarily to what I had envisaged at a certain point, it is not the number of negatively charged, negative-action electrons propagating backward in time, that must be the same as the number of negatively charged, positive-action electrons propagating forward in time at the present moment in our universe, but really the number of positively charged, positive-action electrons propagating forward in time. As a result, the density of matter with reversed bidirectional charge signs (arising from the presence of those particles that survived the early phase of matter-antimatter annihilation) must be the same as that of matter with non-reversed bidirectional charge signs in our universe, while it is the density of negative-energy matter that is allowed to differ substantially from that of positive-energy matter, at least once the primordial annihilation of matter particles with their antimatter counterparts has taken place.

Therefore, it is possible to assume that the average density of baryonic negative-energy matter is currently much smaller than that of positive-energy matter, even if those two densities must have been equal initially, because the ratio of matter to antimatter in the primordial state may have been closer to unity for negative-energy matter, so that a larger portion of negative-energy matter would have been subjected to annihilation. But if that is the case, then even baryonic negative-energy matter with reversed bidirectional charges would be nearly absent at the present time, as its density must equal that of negative-energy matter with non-reversed bidirectional charges, which means that a significant portion of the total average density

of negative-energy matter, possibly as large as the portion of positive-energy matter density which is attributable to the presence of baryonic matter (both visible and dark), would be missing following the early annihilation of baryons with antibaryons. In other words, this early annihilation of matter and antimatter must have produced a difference between the densities of positive and negative energy matter twice as large as would be the case if only baryonic negative-energy matter with non-reversed bidirectional charges had been submitted to annihilation, thereby allowing the specific rates of expansion to diverge more rapidly and the positive cosmological to grow to a greater extent than would otherwise be possible.

Of course this does not mean that no negative-energy matter would remain in the universe, because there was initially a large portion of negative-energy matter in the form of non-baryonic dark matter (as I will explain below). But, as I came to realize, the possibility that there may no longer exist a significant proportion of negative-energy matter in baryonic form actually constitutes a requirement, from an observational viewpoint, as it implies that there can be no significant, localized overdensities in the negative-energy matter distribution on stellar and galactic scales (for reasons that will be discussed in the following section), in agreement with the limits imposed by astronomical observations (concerning weak gravitational lensing experiments in particular).

It must be clear that the above discussed symmetry between the number of particles with non-reversed bidirectional charge signs and that of particles with reversed bidirectional charges is not just a mere possibility, but that the mechanism responsible for the predominance of ordinary matter over ordinary antimatter implies that there must necessarily be as many particles of a given kind with a positive sign of bidirectional charge as there are particles of the same kind with a negative sign of bidirectional charge propagating in the same direction of time. Thus, a certain equilibrium is recovered that allows one to avoid having to conclude that a preferred, absolutely defined sign of bidirectional charge exists, as would be the case from a traditional perspective. But, given that I have argued (based on independent motives) that it is not possible for absolutely (non-relationally) defined physical attributes to exist at a fundamental level, then we may consider that the solution of the issue discussed here is a confirmation of the validity of the hypothesis that there does exist dark matter with reversed bidirectional charges in our universe and that this matter is as abundant as ordinary matter itself at the present epoch.

One obvious, but nonetheless significant consequence which would emerge, if the above proposed solution to the problem of the origin of the asymmetry between matter and antimatter is valid, is that despite the fact that a condition of continuity of the flow of time must apply that forbids all non-gravitational interactions between particles with opposite bidirectional charges, it should nevertheless be possible for such particles to interact through the gravitational interaction, because this assumption is necessary to derive that solution. Indeed, if we do not allow for opposite-bidirectional-charge particles propagating in the same direction of time to be produced by the decay of a graviton, then no imbalance between the number of baryons and that of antibaryons can develop, because only particle-antiparticle pairs involving particles with the same sign of bidirectional charge propagating in opposite directions of time could be created and in such a case no violation of T or C symmetry could give rise to the required imbalances. Thus, the previously discussed argument, to the effect that it should be expected, based on independent motives, that particles with opposite bidirectional charges can interact gravitationally, without violating the condition of continuity of the flow of time along a particle world-line (as a result of the absolute neutrality of the gravitational interaction), appears to be well-founded.

What I would like to explain, now, is that it is actually another phenomenon, made unavoidable by the existence of negative-energy matter, but not associated with the presence of voids in this matter distribution, that is ultimately responsible for most of the missing-mass effects. You may recall that I mentioned in section 4.2 that from the viewpoint of the particular interpretation of the metric conversion factors I have proposed and which allows the emergence of a non-zero value for the cosmological constant, it should be possible for vacuum energy density to vary with position, in addition to have a non-zero value on the global scale. But, in the context of this particular interpretation, it appears that if local variations of vacuum energy density do arise, then they could only be attributable to the fact that local differences may develop between the metric properties of space experienced by positive-energy observers and those experienced by negative-energy observers. What I have come to understand is that, in fact, such variations are unavoidable, given that the presence of an inhomogeneity in the positive- or negative-energy matter distribution produces a variation of the metric properties of space which, for a positive-energy observer, is opposite that which is experienced by a negative-energy observer.

Indeed, the possibility, for opposite-energy observers, to experience differing metric properties of space as a result of the presence of matter inhomogeneities (which is allowed when it is not possible to directly compare such observer-dependent measures of distance) implies that vacuum energy can vary locally, depending on the strength of local gravitational fields, as long as there is no compensation between the local gravitational field attributable to positive-energy matter and that attributable to negative-energy matter. A positive energy observer, for example, would measure a positive increase in vacuum energy density (really a stronger attractive gravitational field) in the presence of a positive-energy matter overdensity, because the measures of spatial volume around that overdensity would appear smaller for such an observer given that from her viewpoint space is contracted and time dilated by the presence of positive-energy matter, while space is dilated and time contracted by the presence of a positive-energy matter overdensity, from the viewpoint of a negative-energy observer, which means that the maximum positive contribution to the density of vacuum energy which is directly experienced only by a negative-energy observer would be larger than the maximum negative contribution which is directly experienced only by our positive-energy observer, given that the corresponding volume of space in which the zero-point vacuum fluctuations that give rise to this maximum negative contribution would take place would appear smaller to this positive-energy observer.

However, the same positive-energy observer would measure a stronger repulsive gravitational field in the presence of an *underdensity* in the positive-energy matter distribution, because the measures of spatial volume around that underdensity would appear larger to her, given that the curvature of space attributable to such an underdensity, like that which would result from the presence of a negative-energy matter overdensity, would give rise to space dilation and time contraction, while it would give rise to space contraction from the viewpoint of an observer of opposite energy sign, which means that the maximum negative contribution to the density of vacuum energy which is directly experienced by the positive-energy observer would be larger than the maximum positive contribution which is directly experienced by a negative-energy observer, so that more gravitationally repulsive vacuum energy would be present in the volume of space around such an underdensity⁸.

⁸I must warn the reader to be careful in asserting the validity of those conclusions, as it is very easy to make a mistake when evaluating the effects of the curvature of space

Those results are a simple consequence of the fact that the different metric properties of space experienced by opposite-energy observers imply different volumes of space, even locally, and therefore also different measures for the maximum positive and negative contributions to the density of vacuum energy provided by the natural vacuum-stress-energy tensors $\gamma^{-+}\mathbf{V}_P^+$ and $-\mathbf{V}_P^-$ (experienced by positive-energy observers), which add up to those arising from the observer-dependent measures of cosmological scale factor, as if we were dealing with an independent stress-energy tensor, similar to that of ordinary matter and independent from the traditionally considered cosmological term $\mathbf{T}_\Lambda = -\Lambda\mathbf{g}$ associated with the cosmological constant. It would, therefore, appear that if the maximum contribution to the density of vacuum energy which is directly experienced by a positive-energy observer was positive (while that which is directly experienced by a negative-energy observer was negative), so that, for such an observer, the metric conversion factor γ^{-+} would rather apply to the natural vacuum-stress-energy tensor that contributes negatively to the density of vacuum energy, then the curvature of space around a positive-energy object would not increase the mass of the object, but rather decrease it, by providing a negative contribution to its total energy. This means that even from an observational viewpoint, it is preferable to assume that the final form of the generalized gravitational field equations I have proposed in section 2.15 (on the basis of consistency requirements) and which implies that the positive cosmological constant does not contribute to its own growth over time, is the correct one, given that it allows one to predict that dark matter actually contributes to enhance the strength of the gravitational field produced by an astronomical object. It is then merely the fact that the local variations of vacuum energy density involved are correlated, under most circumstances, with the presence of local inhomogeneities in the distribution of baryonic matter, due to the fact that such inhomogeneities are usually required to trigger the development of local variations in the density of vacuum energy, that allows them to provide the long sought explanation of the missing-mass effect as being a particular

produced by various matter configurations on the amount of volume present within the boundary surrounding a matter inhomogeneity. I am myself guilty of having once arrived at the exact opposite conclusion as that stated above and not having immediately realized that it was an error, because it produced the desired outcome in the context where the incorrect form of the vacuum-energy term (entering the generalized gravitational field equations introduced in section 2.15), which I later abandoned, appeared to require the validity of that conclusion.

aspect of the phenomenon of dark energy.

What must be understood is that, even if *local* fluctuations in the density of negative-energy matter can be measured by a positive-energy observer and do have an effect on the gravitational field experienced by such an observer, this does not mean that the gravitational fields associated with the presence of negative-energy matter inhomogeneities cannot give rise to additional effects of a gravitational nature, arising from the response of vacuum energy fluctuations to the presence of those gravitational fields. In fact, even in the absence of any inhomogeneity in the negative-energy matter distribution, there may arise local variations of vacuum energy density as a result of the presence of positive-energy matter inhomogeneities, and the gravitational fields attributable to those local variations of vacuum energy density would actually affect the motion of both positive- and negative-energy bodies. Therefore, I believe that what explains most of the missing-mass effect around visible positive energy structures is the fact that the gravitational fields produced by those inhomogeneities in the matter distribution give rise to such local variations of the density of vacuum energy that must necessarily be concentrated around the visible structures and that must themselves give rise to further variations of vacuum energy density, arising from the gravitational fields produced by those very same concentrations of vacuum energy. The crucial point here, therefore, is that the positive vacuum energy which is produced by the curvature of space attributable to the presence of positive-energy matter must itself contribute to produce additional space curvature, similar to that which would be produced by ordinary positive-energy matter, which in turn produces additional positive vacuum energy.

The problem we would normally face, in such a context, is that it would seem that the mass of an astronomical object would be allowed to increase without limit, as the growth of mass arising from the concentration of vacuum energy would trigger the formation of an even larger concentration of vacuum energy that would further increase the mass of the object. But, in fact, the situation we face here is no more problematic than that which arises as a result of ordinary gravitational instability, because, as I will explain below, the portion of dark matter attributable to local variations of vacuum energy density already existed in a macroscopically homogeneous form on the cosmic scale, before it accumulated around ordinary matter overdensities, yet we already know that gravitational instability cannot produce catastrophic outcomes when the matter distribution is relatively homogeneous, given that the energy of the gravitational field generated by a positive-energy body is

opposite the energy of the source, while the field also interacts with itself. Therefore, the growth of mass attributable to local variations of vacuum energy should be limited, especially since the gravitational interaction itself is very weak. This does not mean that no such an effect would exist, however. In fact, it appears that once one recognizes that baryonic negative-energy matter itself can exist (or at least that it must have existed at a certain epoch), one cannot avoid the conclusion that such local variations of vacuum energy density and a phenomenon which we may call *vacuum dark matter* would arise which would have consequences similar to those we normally attribute to the presence of ordinary dark matter (such as weakly interacting massive particles), as it would actually contribute to significantly increase the mass of any astronomical object present on a sufficiently large scale, without raising the density of baryonic matter.

Now, I must admit that for a long time I, myself, believed that local variations of vacuum energy density could not constitute a solution to the missing-mass problem, because I thought that the equivalent mass attributable to vacuum dark matter would not be allowed to contribute to the total energy of matter that is required to bring the density of positive energy to its critical value, given that there would also be negative contributions to the energy of matter that would arise from those local variations of vacuum energy density attributable to the presence of negative-energy matter overdensities or positive-energy matter underdensities, which I thought would cancel out the additional positive contributions and prevent the density of positive matter and vacuum energy from reaching its critical value, while such negative contributions also appeared unavoidable. In other words, I had forgotten about the idea, because, when I first considered this possibility, I thought that, given that the energies involved were particular instances of vacuum energy, then both the positive and the negative contributions should add up to produce a null density that would not allow to increase the density of positive matter and vacuum energy to its critical value, while this appeared to be required from a theoretical viewpoint, particularly in the context of inflationary cosmology.

Also, when I began seriously considering the possibility that some local variations of vacuum energy density attributable to the gravitational field of large astronomical objects could be responsible for the phenomenon of missing mass, I had actually (but inappropriately) come to believe that voids in the negative-energy matter distribution could provide an alternative explanation to most of the missing-mass effects around visible structures and

therefore I didn't see the need that there was to explain the missing-mass effect as being the outcome of an inhomogeneous distribution of vacuum energy attributable to the presence of matter, even if the existence of such a phenomenon appeared unavoidable. It is only much later that I came to understand that the fact that we are dealing here with *local* variations in the distribution of vacuum energy means that, from a gravitational viewpoint, vacuum dark matter would be equivalent to the presence of ordinary matter, which allows one to expect that the gravitational attraction of positive vacuum-dark-matter energy would not be compensated by a gravitational repulsion attributable to the presence of negative vacuum-dark-matter energy, on the cosmological scale. But for the exact same reason, we should also expect that a local measure of positive vacuum-dark-matter energy, unlike a positive cosmological constant, should not exert a negative pressure that would repel positive-energy matter, either locally or globally (because the average density of positive vacuum-dark-matter energy does not remain constant in an expanding universe).

What is unique, in effect, about the interpretation of the inhomogeneous character of the distribution of vacuum energy discussed here (which is derived from the generalized gravitational field equations introduced in section 2.15) is that, despite the fact that vacuum dark matter is a form of vacuum energy, that must consequently be dark, it nevertheless contributes to the gravitational dynamics of the universe on a global scale in much the same way ordinary matter does. Indeed, if most of the missing mass-effect is attributable to local variations of vacuum energy density, then the gravitational forces exerted by dark matter must be similar to those which are attributable to the presence of voids in the otherwise uniform distribution of positive and negative vacuum energies, while, as I explained in section 2.8, a void of cosmic proportions in the positive portion of the distribution of vacuum energy (arising from the presence of a globally homogeneous distribution of negative vacuum-dark-matter energy) exerts a null gravitational force on positive-energy matter and has no effect on the rate of expansion determined by positive-energy observers.

The one crucial aspect that differentiates vacuum dark matter from ordinary matter, therefore, is the fact that, as I explained in section 2.8, if we are to conceive of negative-energy matter as missing positive vacuum energy, then we have no choice but to assume that what is missing from zero-point fluctuations under such conditions is not just positive energy, but also some positive or negative (non-gravitational) charges, as the missing virtual par-

ticles which are equivalent to the presence of matter also carry charges in addition to energy. Those charges which are missing in the electrically (or non-gravitationally) neutral vacuum are equivalent to the presence of charges of opposite signs, which appear to be carried by the matter particles. The above proposed concept of vacuum dark matter, however, makes it clear that the presence of dark matter does not result from a local absence of virtual particles from zero-point vacuum fluctuations, that would need to be correlated with an absence of charge, but arise from the differing measures of metric properties experienced by opposite-energy observers, which alter the relative densities of the maximum positive and negative portions of vacuum energy without affecting the electrical (or non-gravitational) neutrality of the vacuum. In other words, while ordinary positive-energy matter consists of both missing negative energy and missing positive *or* negative charge, positive vacuum-dark-matter energy is only equivalent to missing negative vacuum *energy* and must remain electrically neutral (it does not carry any non-gravitational charges). This is the only aspect that differentiates vacuum dark matter from ordinary matter.

Dark matter, therefore, appears to be a hybrid form of matter that shares some properties of the uniformly distributed portion of vacuum energy, or the cosmological constant, but that produces gravitational forces which are equivalent to those produced by ordinary matter, due precisely to the fact that its presence is attributable to *local* variations of vacuum energy density and because its density must vary like that of ordinary matter, as a result of expansion. Thus, again, we can expect that as long as it is uniformly distributed on a global scale, negative vacuum-dark-matter energy, just like ordinary negative matter energy, does not, in fact, affect the rate of expansion of matter determined by positive-energy observers and does not contribute to the critical energy density that is relevant to those observers, unlike the negative component of a uniform distribution of vacuum energy. It is only when vacuum-dark-matter energy accumulates around massive astronomical objects, that its presence becomes apparent to both positive and negative energy observers.

Now, it must be clear that the total quantity of energy contained in positive-energy dark matter, like that contained in negative-energy dark matter, does not change with time on the largest scale (even when the amount of positive or negative energy contained in matter itself varies, as must have happened following the early annihilation of matter with antimatter), despite the fact that the portion of missing-mass effects attributable to local varia-

tions in the density of vacuum energy only becomes apparent when macroscopic inhomogeneities develop in the matter distribution and those energies become more concentrated around large astronomical objects. Thus, the additional amount of energy that is present around positive-energy galaxies, but that cannot be accounted for by the presence of baryonic matter, was already present, in more diffuse form, before the formation of those structures, even though it was then exerting a significant gravitational force only on the global scale. What makes this possible is the fact that, even though the distribution of matter and radiation energy in the early universe was very homogeneous from a macroscopic viewpoint (a hypothesis that is theoretically and observationally unavoidable, as I will explain in section 4.9), on the quantum-gravitational scale the magnitude of the initial positive and negative densities of matter and radiation energy was maximum, which means that there existed variations of maximum magnitude in the density of energy attributable to a variation of the sign of energy of matter from one discrete location to another, to which were associated opposite amounts of vacuum dark matter energy.

As the universe expanded and the average density of matter decreased, along with the average kinetic energy of matter and radiation particles, the macroscopically homogeneous distribution of vacuum dark matter (which existed as a consequence of the presence of microscopic inhomogeneities in the primordial distribution of matter and radiation energy) spread into the available space along with the rest of matter and it is only when the small-amplitude inhomogeneities which were present on a macroscopic scale in the matter distribution (including that of vacuum dark matter) began to grow, later on, as a result of gravitational instability, that vacuum dark matter itself began to concentrate in those macroscopic inhomogeneities. But this vacuum-dark-matter energy was not created when those macroscopic inhomogeneities developed, despite the fact that the gravitational forces it exerts are apparent only in those places where matter is inhomogeneously distributed on a macroscopic scale and the curvature of space is more developed. In section 4.5 I will explain that if that was not the case, and the amount of positive-energy dark matter attributable to local variations of vacuum energy density was actually growing in the universe (along with that of negative-energy dark matter), difficulties would arise, even if the total energy of matter (comprising the contributions of both positive- and negative-energy dark matter) was conserved in the process.

What must be understood here is that what we usually describe as a

homogeneous matter distribution actually contains inhomogeneities on a smaller scale and this additional amount of structure gives rise to local variations in the gravitational field or the curvature of space which affect local measures of vacuum energy density (how this is possible will become clearer once the reader learns about certain unexpected microscopic properties of gravitational fields in section 4.7). As a result, even in places where no macroscopic inhomogeneities are present in the matter distribution, there usually exist local variations in the metric properties of space, attributable to the presence of microscopic inhomogeneities in such a macroscopically homogeneous matter distribution and those inhomogeneities result in the presence of both positive and negative vacuum-dark-matter energy. But it is only when the density of positive-energy matter grows larger than its average value on a macroscopic scale, at the expense of the creation of an underdensity of equal magnitude in the surrounding, homogeneous, positive-energy matter distribution, that the uniformly distributed positive-energy portion of vacuum dark matter becomes rarefied in the underdensity and more concentrated around the developing positive-energy structure, where it can begin to exert a gravitational pull on nearby positive-energy matter. At the same time, the uniformly distributed negative vacuum-dark-matter energy becomes rarefied in the positive-energy matter overdensity and more concentrated around the void that formed (as a consequence of the formation of this positive-energy structure) in the surrounding positive-energy matter distribution, where it can now exert a measurable gravitational repulsion on surrounding positive-energy matter.

It must be clear, again, however, that vacuum dark matter is not created by the structure formation process. It was already present in the initial Big Bang state, before macroscopic inhomogeneities began to grow in the matter distribution, only this vacuum dark matter was then more homogeneously distributed in space, from a macroscopic viewpoint and only exerted its full influence on the gravitational dynamics of the universe as a whole. This means that there is no growth in the total amount of positive-energy dark matter when a local overdensity develops in the positive-energy matter distribution, even if it is indeed the curvature of space attributable to this overdensity that is responsible for the presence of vacuum dark matter, because this overdensity formed through the accumulation of both baryonic positive-energy matter and positive vacuum-dark-matter energy and it is only because the initial vacuum-dark-matter distribution was macroscopically very homogeneous that its influence is allowed to grow with time as the

curvature of space itself grows (locally) on a macroscopic scale. Of course, the same is true for the positive vacuum-dark-matter energy that must accumulate, locally, when a void forms (on a large scale) in the macroscopically more homogeneous distribution of negative-energy matter, which already existed before it was made conspicuous by its accumulation in this particular region of space, just like the negative vacuum-dark-matter energy that would accumulate outside this negative-energy matter underdensity.

From an observational perspective, it would appear possible to confirm that dark matter is, for the most part, an outcome of spatial variations in the density of vacuum energy, because currently available data indicates [31] that there is a strong correlation, in general, between the gravitational acceleration attributable to the total amount of matter inside an orbit (say around the center of a galaxy) and the gravitational acceleration attributable to the baryonic matter (and this correlation would probably be even stronger if we were taking into account the presence of baryonic dark matter). Indeed, if the presence of dark matter must be considered to be an effect of the curvature of space (attributable to the matter that is present in a region of space) on the local measures of vacuum energy density, then the more gravitational acceleration that there is as a consequence of the presence of baryonic matter, the more distinct the metric properties of space experienced by opposite-energy observers must be that gave rise to the accumulation of vacuum dark matter around that particular location.

Even though the importance of the empirically derived relationship that allows to confirm the validity of those conclusions is often overlooked, it would certainly be a significant problem if it was to remain unexplained, as would be the case from the viewpoint of a more conventional interpretation of the missing-mass effect (given that in such a context dark matter is simply an additional component of matter whose existence does not depend directly on the presence of ordinary matter). But the conclusion that there must exist a relationship between the amplitude of the gravitational field attributable to visible positive-energy matter overdensities and the amplitude of the missing-mass effect would also imply that, even within galaxies and clusters, the dark matter should be more concentrated around the visible elements of the structure. While this result is certainly unexpected, it does, in fact, agree with some recent observations [32], which indicate that there is a greater than expected concentration of gravitational lensing around individual galaxies within clusters. There is, thus, a strong motive to prefer an interpretation of the missing-mass effect as being a consequence of local variations in the

density of vacuum energy, which must exert gravitational forces proportional to those produced by the baryonic matter inhomogeneities to which their presence is correlated.

It is important to point out, however, that vacuum dark matter would exert its own gravitational field, which would actually allow it to clump together just like conventional dark matter, despite the fact that it really is vacuum energy. This means that the observations which indicate that large overdensities of visible matter can sometimes become separated from their dark matter component (as a result of collisions between galaxy clusters or in the course of galaxy mergers) can be easily explained, unlike would be the case if the currently unexplained correlations discussed above (between the gravitational acceleration attributable to the total amount of matter inside an orbit and the gravitational acceleration attributable to baryonic matter alone) were the result of a more profound modification of the laws that govern the gravitational dynamics of astronomical objects (such as envisaged in the context of the theory known under the MOND acronym). Indeed, once created, such a dark-matter object could continue to exist all by itself, sustained merely by its own gravitational field, just like voids in the matter distribution, while only a minimum measure of vacuum dark matter would be left in the visible structure around which it was originally located⁹. This is a considerable advantage of this original approach which, once again, appears to confirm the validity of the generalized gravitation theory developed in the earlier portions of this report.

To conclude this section, I would like to briefly return to the problem of black-hole information and entropy which was discussed in section 3.10. An important conclusion at which I arrived while trying to determine the nature of the microscopic degrees of freedom of the matter particles captured by the gravitational field of a black hole is that the portion of missing information which is encoded in the microscopic degrees of freedom of the gravitational field on the surface of a stable-state black hole would only allow to determine

⁹From that viewpoint, it would seem that the galaxies which do not appear to contain much dark matter are not galaxies which produce no local variations of vacuum energy density, or which did not grow out of the accumulation of vacuum-dark-matter energy, but merely galaxies which were stripped off of that portion of their dark matter attributable to the presence of inhomogeneities in the distribution of vacuum dark matter itself, as a result of encounters with more massive structures (a conclusion which is supported by the latest computer simulations [33] involving dwarf galaxies, once dark matter is recognized to originate from local variations of vacuum energy density).

the sign of space intervals associated with the momentum direction of each and every matter particle that collapses into its inner singularity and which is then submitted to a quantum bounce. The other physical parameters characterizing the microscopic state of matter particles under the influence of a stable-state black hole which could influence the curvature of spacetime, like the magnitude and the sign of their energies and the magnitude and the orientation of their momenta, are all fixed to common unique values, as a result of the constraints imposed by the gravitational field that is present in the vicinity of the inner singularity. I also explained that, by necessity, the missing information concerning the sign of charge of matter particles (which is transformed by the redefined time-reversal symmetry operation T) would need to be encoded in the microscopic state of the field of interaction associated with this charge and is not reflected in the microscopic configuration of a black hole's surface gravitational field, while the missing information concerning the handedness of particles (which depends on the direction of their spin) would be contained in the microscopic degrees of freedom of the component of the gravitational field associated with the torsion of spacetime. It is only under such conditions that one can obtain the right measure of missing information (that which is determined by the semi-classical theory of black-hole thermodynamics) in the case of elementary black holes containing, at most, one matter particle.

However, in the context where the sign of charge of a most elementary particle may not only differ as a consequence of a reversal of the direction of propagation in time, but may also be different for particles with opposite bidirectional charges propagating in the same direction of time, one may wonder whether it would still be possible to determine the direction of propagation in time of a given particle from information contained in the microscopic state of the field of interaction associated with this charge? This is an important question, because if it is not possible to assess the direction of propagation in time of a particle that was captured by the gravitational field of a black hole, then one would have to conclude that some physically significant aspect of the state of matter particles cannot be uniquely determined from the information that is contained in the microscopic state of the fields of interaction on the event horizon of such an object, which would imply that information is lost when matter is submitted to gravitational collapse. It may, therefore, appear that if reversed-bidirectional-charge particles are allowed to exist, there would be a problem with the fact that the information about the sign of electric charge, that would be provided by the microscopic degrees of freedom

of the electromagnetic field on the surface of a black hole, would not allow to differentiate between a positively-charged electron propagating forward in time and an ordinary positron described as a negatively-charged electron propagating backward in time, while there clearly exists a degree of freedom associated with this physical property, which normally allows to differentiate between matter that is visible and matter that is dark.

The above discussion, however, makes it clear that it need not be the case that information about the direction of propagation in time is lost in the presence of reversed-bidirectional-charge particles, precisely because particles with reversed bidirectional charges would need to be dark from the viewpoint of an observer made of ordinary matter. Indeed, the field that contains the information about the sign of charges, or the direction of propagation in time of ordinary matter particles, is not the exact same field as that which contains the information regarding the sign of charge, or the direction of propagation in time of particles with reversed bidirectional charges. It is the microscopic state of the electromagnetic field with which positive bidirectional charge electrons interact that contains the information about the direction of propagation in time of those particles and given that one can differentiate between this field and that which is produced by electrons with negative bidirectional charges, then it is possible to obtain information about both the bidirectional charge sign of elementary particles and their direction of propagation in time from a determination of the microscopic state of all components of the electromagnetic field on the surface containing those particles. Thus, while the distinction between ordinary electrons and anti-electrons is encoded in the microscopic state of the electromagnetic field that interacts with visible matter, the distinction between baryonic dark-matter electrons and anti-electrons of the same kind is encoded in a different component of the electromagnetic field, which is that with which only baryonic dark-matter particles interact.

This conclusion, which is dependent on the above proposed interpretation of reversed-bidirectional-charge particles, is actually much more unavoidable than one may expect. Indeed, as I already mentioned, it seems that once the sign of energy of an elementary black hole (present on the quantum gravitational scale) is fixed (which also determines the energy sign of the particle submitted to its gravitational field), three discrete variables remain undetermined concerning the microscopic state of the particle under its influence, which are the direction of the particle's momentum, or the sign of space intervals associated with its motion (downward or upward), the parti-

cle's handedness (spin-up or spin-down), and the sign of time intervals from which are determined its particle or antiparticle nature from a unidirectional time perspective. But, if a quantum theory of gravitation is to eventually constitute a unified theory of all interactions, it can be expected that additional information would need to be provided to specify the sign of the bidirectional charge of a most elementary particle, that determines which type of non-gravitational field it interacts with, that is to say, information would need to be available to determine whether such a particle is visible or not, from the viewpoint of a given observer with the same sign of energy. Now, what I'm suggesting is that, not only are those the only fundamental parameters which can vary in a discrete way under such conditions and which actually allow to characterize the state of any matter particle on the most fundamental scale, and not only is it possible for the information that is required to determine the value of each of those parameters to be encoded on the surface of an elementary black hole, but in fact, this is the only information that *could* be encoded on an elementary surface.

I have already mentioned, in effect (in section 3.10), that each elementary unit of surface which is considered to correspond with one binary unit of information in the semi-classical theory of black-hole thermodynamics, actually contains four Planck units of surface (associated with a Planck unit of distance). Why exactly four Planck units of surface should be necessary to encode one fundamental unit of information has always remained unexplained. But in the context of the present semi-classical description of the degrees of freedom of a matter particle which is under the influence of an elementary black hole (containing one Planck unit of energy), the fact that we may need four units of area to determine the state of each elementary particle present on the quantum gravitational scale would no longer constitute a mystery, once we recognize that to each of those units of area there should correspond one discrete, elementary degree of freedom, because four microscopic parameters must be determined for each particle (one for the momentum direction, one for the handedness, one for the direction of propagation in time, and one for the sign of unified bidirectional charge), even though only the momentum direction of each particle contributes to determine the thermodynamic properties associated with the surface gravitational field of a macroscopic stable-state black hole, as required by the semi-classical theory of black hole thermodynamics.

Now, despite the fact that the area gap derived from loop quantum gravity is not explicitly and necessarily equal to a multiple of the Planck unit

of surface and is rather provided by a more complex equation that contains a free parameter that must be adjusted to match the results of the semi-classical theory, I believe that it is nevertheless appropriate to assume that the phenomenological unit of area provided by the semi-classical theory of black hole thermodynamics actually constitutes the true fundamental unit of surface. The problem with the theoretically derived area gap is that it is obtained on the basis of purely gravitational aspects of Planck-scale physics, unlike the area of the event horizon of a macroscopic black hole in the semi-classical theory, which also varies as a function of the electric charge (unlike the gravitational field of the object itself). Thus, it would only be appropriate to conclude that the quantum gravitational area gap really is not equal to four times the Planck area if there was no reason to believe that the value of this parameter derived from loop quantum gravity may need to be adjusted by taking into account the details of the currently incomplete unified theory of elementary particle interactions that must apply on this scale and which should allow the more natural value produced by the phenomenological approach to be recovered, under the assumption that four discrete, fundamental degrees of freedom, among which would be the sign of unified bidirectional charge, contribute to the measure of black hole entropy. Of course, this does not mean that current theoretical estimates for the value of the area gap are useless, because they do provide the correct order of magnitude for such a parameter and they do confirm that there must be such a limit to the continuity of space from a quantum gravitational viewpoint.

In any case, it must be clear that the conclusion stated above, to the effect that each Planck unit of area can be associated with one elementary degree of freedom for the matter that is constrained by the gravitational field of an elementary black hole, is not contradicted by the standard derivation of the measure of black-hole entropy, because, under such conditions, three out of each four units of missing information associated with a unit of surface that contains four Planck units of area are irrelevant to the definition of the thermodynamic properties of the gravitational field associated with the curvature of spacetime and can actually be ignored. Indeed, the handedness of all particles that reached the singularity of a stable-state black hole is fixed by the microscopic (quantum gravitational) degrees of freedom of the gravitational field associated with the torsion of spacetime, while their direction of propagation in time only influences the microscopic properties of the field of interaction associated with the unified non-gravitational charge and the sign of bidirectional charge merely determines which component of this

field encodes the information about the sign of charge, thereby leaving only the momentum direction of particles to be determined by the microscopic degrees of freedom of the surface gravitational field of a black hole associated with the curvature of spacetime. It is quite remarkable that such an exact quantitative result can be entirely derived from logical arguments made in the context of a semi-classical approximation. I believe that this conclusion, more than any other, illustrates the effectiveness of an unconventional approach, such as the one I came to adopt, for solving certain kinds of problems of particular importance in fundamental theoretical physics.

4.4 Large-scale structure

I remember, as a teenager, before I even learned about the existence of dark matter, having been deeply amazed and puzzled after reading in the newspaper that astrophysicists had determined that most of the visible matter in the universe, including our own galaxy, was located on the surface of giant voids of truly enormous proportions forming a bubble-like pattern in the matter distribution. I cannot say that I already expected, back then, that I would eventually be involved in developing a model that would help explain this troubling observation, but I did feel very strongly that this was something I needed to better understand. Anyhow, this stunning discovery and the mystery that initially surrounded it helped shape my early approach to the problem of gravitation in a way that turned out to be highly fruitful. What is truly remarkable is that the problem of voids has endured to this day, as we kept discovering empty regions of increasingly larger sizes that still defy traditional explanations, despite all the progress which was achieved in developing cosmological models that can more accurately reproduce those features.

I believe that the introduction of negative-energy matter will have a significant impact on theories of structure formation. Indeed, what emanates from the results discussed in the preceding section is that the formation of overdensities in the primordial distribution of positive-energy matter must have been accelerated by the presence of negative-energy matter underdensities, which developed as a result of the gravitational repulsion exerted by those positive-energy matter overdensities and whose effects on developing positive-energy structures would be indistinguishable from those of ordinary, positive-energy dark matter. What needs to be emphasized, in this regard,

is that, given that, in the early universe, the average matter density (not just that of negative-energy matter) was much larger than it currently is, then it follows that the underdensities which were present in the distribution of negative vacuum-dark-matter energy had a significant influence on positive-energy matter. As I previously mentioned, certain observations [29], performed after the first versions of the present document were published, appear to confirm that the formation of galaxies must have been accelerated, in the primordial universe, by effects of an unknown origin, because the observations in question have revealed the existence of galaxies which are too large at such an early epoch for their presence to be explained using more conventional models (which do not involve negative-energy matter). In fact, it had been known for some time that the first large elliptical galaxies appear to have formed too early after the Big Bang for their creation to be easily explainable using conventional models. But if we recognize that the presence of negative-energy matter underdensities must have played an important role on such a scale in the remote past, when the average density of negative vacuum-dark-matter energy was much larger, and its distribution much more homogeneous, then those mysteries can be explained quite straightforwardly.

Now, it is important to understand that the observationally confirmed absence of localized overdensities in the distribution of negative vacuum-dark-matter energy, on the scale of individual stars and galaxies means that there was no similar influence on negative-energy matter arising from the presence of underdensities in the early distribution positive-energy matter on smaller scales, because only the presence of negative-energy matter overdensities would allow the depth of those underdensities to grow large enough, on this particular scale, that they could significantly accelerate the formation of structures in the negative-energy matter distribution. Of course, one may ask why it is that there were no significant overdensities in the negative vacuum-dark-matter energy distribution on smaller scales during the matter-dominated era, even though (as I will suggest in section 4.5) the average density of negative vacuum-dark-matter energy was always very similar to that of positive vacuum-dark-matter energy? I wrestled with that question for a long time before I realized that the absence of negative-energy matter overdensities is due to the fact that only a negligible amount of baryonic negative-energy matter has survived the early annihilation of matter with antimatter.

What happens is that the presence of baryonic matter is necessary for

initiating the formation of structures in a distribution of vacuum dark matter on smaller scales, because only such matter can reduce its kinetic energy through the emission of radiation. It is well-known, in effect, that collapsed structures can only begin to form through gravitational instability when they are allowed to release energy in such a way, as otherwise their internal pressure remains too large to allow them to stabilize. But if one may assume that vacuum dark matter shares some of the properties of baryonic matter, due to the fact that zero-point vacuum fluctuations involve the same particles fluctuating in and out of existence in their virtual form, then one would be justified to assume that, under the same conditions of temperature, vacuum dark matter is merely as cold as it would be if it was composed of baryonic matter (both quarks and leptons)¹⁰. Under such conditions, it follows that, on a scale where ordinary baryonic matter is allowed to collapse into stable structures only through the emission of radiation, the density of vacuum dark matter itself can only begin to grow if ordinary baryonic matter overdensities with the same sign of energy are already present, because without the additional gravitational attraction produced by such an overdensity, vacuum dark matter would rather tend to disperse, if it is merely as cold as ordinary baryonic matter.

Thus, it would appear that it is the possibility for baryonic negative-energy matter and antimatter to have annihilated more completely than baryonic positive-energy matter and antimatter in the early universe, that explains that no significant, gravitationally repulsive overdensities of negative-energy matter were present in the early universe (and therefore also at later times) on the scale of stars and galaxies. Such an outcome would be expected to occur whenever a more limited violation of time-reversal symmetry (or indeed an absence of such violation) would affect the production of opposite-bidirectional-charge particles and antiparticles with negative action, than affected the production of opposite-bidirectional-charge particles and antiparticles with positive-action. What must be understood is that the presence of negative-energy matter overdensities is not unavoidable, be-

¹⁰It would be problematic to assume that it is not possible to say whether vacuum dark matter is either cold or hot or anything in between, because that would mean that its physical properties are undefined, despite the fact that it must behave unambiguously under any particular circumstances. But the fact that it appears that vacuum dark matter must be as cold as baryonic matter is certainly not undesirable given that it was shown that only cold dark matter can give rise to a bottom-up process of structure formation of the kind that is favored by astronomical observations.

cause, even if baryonic negative-energy matter must have been present in the universe initially, it can be almost completely absent during the matter-dominated era. Anyhow, one must conclude that immediately after most of the baryonic negative-energy matter and antimatter annihilated, during the radiation-dominated era, the average density of baryonic negative-energy matter did become completely negligible, because this is the only way one can explain that no overdensities later developed in the negative-energy matter distribution, on smaller scales, whose presence could have been revealed by weak gravitational lensing experiments.

That is not to say, however, that there can be no overdensities in the distribution of negative vacuum-dark-matter energy. In fact, it appears necessary to assume that vacuum-dark-matter overdensities with very large negative masses would be present inside the largest voids in the distribution of positive-energy galaxies, which do exert a localized gravitational pull on this vacuum dark matter. Given that the gravitational attraction exerted on negative vacuum-dark-matter energy by voids in the positive-energy matter distribution grows along with their size, it is possible to conclude, in effect, that most of the overdensities which are present today in the distribution of negative vacuum-dark-matter energy should be concentrated in the largest voids in the positive-energy matter distribution. In the absence of baryonic negative-energy matter overdensities, it would appear that only such large voids can exert a gravitational attraction large enough to allow negative vacuum-dark-matter energy overdensities to form and to grow to such proportions that they may actually influence the process of structure formation in the positive-energy matter distribution.

What's interesting here is that the presence of such very-large-scale overdensities in the negative-energy matter distribution is not only allowed by current observational data, it is actually required in order to explain the size of those voids. Indeed, in the absence of large negative-energy matter overdensities located within the largest voids in the positive-energy matter distribution, the unexpectedly large size of those voids would need to be explained through biasing, an approach which amounts to assume, without justification, that galaxies have a tendency to form preferably in those regions where the density of the cold dark matter is already larger, at the epoch of recombination. If I do not agree that biasing is the appropriate solution to the presence of unexpectedly large voids in the galaxy distribution, it is because it seems to me that biasing merely amounts to impose the required matter distribution without explaining it. In fact, it cannot even be assumed that

the solution offered by the biasing hypothesis is merely incomplete and that what one must do is explain why the largest overdensities in the primordial distribution of positive vacuum-dark-matter energy formed predominantly in those regions where baryonic positive-energy matter overdensities were located, as I have suggested would occur at later times on smaller scales, because any correlation between those two types of inhomogeneities that may have existed in the initial Big Bang state would have been lost, early on, due to the fact that any preexisting baryonic matter inhomogeneity would be wiped out on smaller scales, before new ones would be allowed to form during the matter-dominated era, influenced by the presence of the original inhomogeneities which had continued to develop in the distribution of vacuum dark matter, on a larger scale¹¹.

It should be clear though, that given the relatively small, current, average value of positive-energy matter density, only the largest of the voids in the positive-energy matter distribution should produce a localized gravitational attraction on negative vacuum-dark-matter energy that would be large enough to allow a measurable overdensity of such matter to exist inside those structures (the effect cannot have been more substantial in the early universe, because, despite the larger average density of positive vacuum-dark-matter energy, the amplitude of density fluctuations on such a large scale was then smaller). Thus, despite the absence of small-scale overdensities in the early distribution of negative vacuum-dark-matter energy, the presence of sufficiently large voids in the distribution of positive-energy galaxies should allow gravitational instability to arise in this distribution of negative vacuum-dark-matter energy, which, once triggered, would give rise to a self amplifying process that could produce more overdensity on such a scale. It is the fact that baryonic positive-energy matter overdensities can only gravitationally repel negative vacuum-dark-matter energy and make it spread, that explains that, in the absence of baryonic negative-energy matter overdensities, the only place where the density of this vacuum dark matter was allowed to

¹¹The fact that vacuum dark matter is only submitted to the gravitational interaction is what explains that dark matter inhomogeneities were allowed to grow larger, sooner than the inhomogeneities which were present in the baryonic matter distribution, despite the fact that those two types of fluctuations had the same scale invariant magnitude in the initial Big Bang state. For observational reasons, it is not possible to assume that vacuum dark matter overdensities began to grow at the same time as those in the baryonic matter distribution, but that the presence of vacuum dark matter merely accelerated the rate of structure formation that took place subsequently.

grow and to exert an influence on positive-energy matter is in the large-scale underdensities that formed in the positive-energy matter distribution. Nevertheless, it appears that the presence of large and growing overdensities of negative vacuum-dark-matter energy inside those voids in the positive-energy matter distribution would allow the voids themselves to grow larger than expected at an earlier time, due to the gravitational repulsion they would exert on the surrounding positive-energy galaxies.

There are still reasons to expect, therefore, that negative-energy matter must have contributed to the process of structure development on the largest scales, not just because this is not ruled out from an observational viewpoint, but also because it would actually help explain certain observations mentioned above, like the size of the largest voids in the galaxy distribution or the mass of the earliest galaxies. But given that the gravitational repulsion attributable to a large negative-energy matter overdensity would be similar to that which is attributable to the void in the expanding positive-energy matter distribution in which it would normally be located, then it follows that the presence of the overdensity would only enhance the equivalent gravitational repulsion of the void. I think that this is what explains that it is not immediately apparent that the gravitational forces involved are sometimes attributable in part to the presence of gravitationally-repulsive material. However, the presence of a negative-energy matter overdensity within such a large void in the expanding positive-energy matter distribution would produce some distinctive effects, as it implies that the structure can exert an unexpectedly large gravitational repulsion on the surrounding positive-energy matter. This is certainly a positive development, given that it has been known for some time that certain voids, apparent on a very large scale in the positive-energy matter distribution, do exert larger than expected gravitational repulsion on galaxies located in their periphery, a phenomenon which had remained unexplained until now.

When gravitationally-repulsive matter is present inside the largest voids in the visible matter distribution, it is also easier to reconcile our theory of structure formation with those observations which show that there is a much smaller number of galaxies in the Local Void than is predicted by computer simulations, because any galaxy that would form in the void would rapidly be expelled to the periphery by larger than expected repulsive forces. Also, given that the density of negative vacuum-dark-matter energy in the Local Void would not be as low as it would in our galactic neighborhood, it follows that the missing-mass effects attributable to negative-energy matter

underdensities would be more localized around those galaxies located nearer the void and this must have accelerated their formation. That may explain why a larger than expected number of very large galaxies in the Local Sheet are located on the periphery of the Local Void instead of in the more crowded areas, where most of the visible matter is concentrated. But it must be clear that in the absence of overdensities in the baryonic negative-energy matter distribution, negative vacuum-dark-matter energy overdensities large enough to trigger such phenomena can only develop on a very large scale, where the average density of positive-energy matter is quite low, but where the effects attributable to its absence are allowed to accumulate.

The fact that the only observed departures from what is currently predicted to occur on the scale of individual galaxies do not involve gravitational repulsion, but merely require the existence of additional gravitational attraction localized over the visible structures, would appear to confirm that those deviations are attributable to the presence of early underdensities in the negative vacuum-dark-matter energy distribution. It remains, though, that negative-energy matter is the source of additional gravitational instability which does not arise only from stronger gravitational attraction, but also from the gravitational repulsion exerted on positive-energy matter by the overdensities of negative-energy matter that are present on a very large scale. This means that, starting from the same relatively smooth initial distribution of positive-energy matter that is revealed by the low amplitude of cosmic microwave background temperature fluctuations, we can actually expect inhomogeneities to develop more rapidly in this matter distribution, due to the presence of smaller-scale underdensities and larger-scale overdensities in the distribution of negative vacuum-dark-matter energy. But given that a certain portion of the fluctuations which are observed in the temperature of CMB radiation can be attributed to the presence of inhomogeneities in the primordial distribution of negative vacuum-dark-matter energy, then one must conclude that the magnitude of the inhomogeneities which were present, at the epoch of last scattering, in the positive-energy matter distribution is smaller than what seems to be indicated by observations of CMB temperature fluctuations, which means that their magnitude has not necessarily developed to become larger than one would otherwise expect, particularly on smaller scales (I will return to this question in section 4.9).

Now, even though it appears that most of the baryonic negative-energy matter has vanished at a very early time, shortly after the quark-hadron transition, it cannot be the case that the only negative-energy matter that

remains in the universe today is vacuum dark matter, because negative-energy neutrinos, at least, must have survived the early period of matter-antimatter annihilation and be present with an average density as large as that of positive-energy neutrinos, along with negative-energy radiation. Indeed, even if neutrinos and antineutrinos existed in the exact same numbers initially, they cannot be expected to annihilate completely with one another, given their low mass and the weakness of their interactions. But this negative-energy matter would have been hot when it decoupled from the rest of matter and radiation and therefore it must have remained homogeneously distributed and cannot have contributed much to the process of structure formation, particularly on smaller scales. In fact, given that it is known that positive-energy neutrinos themselves cannot exert much influence on the formation of structures in the positive-energy matter distribution, due to their low mass, then it is certainly not necessary to take into account any potential effect that could be attributed to the presence of underdensities in the distribution of negative-energy neutrinos, even if it seems that the fact that this matter is still homogeneously distributed on the scale of galaxies and clusters could allow the missing-mass effect attributed to those underdensities to be localized around gravitationally repelling positive-energy structures.

What must be retained from all this is that the additional influence which may be exerted, under certain conditions, by negative vacuum-dark-matter energy on the process of structure formation in the positive-energy matter distribution may be significant enough to have given rise to structures which are more developed than those which are predicted to arise by the conventional cold-dark-matter model at various cosmological epochs, either at an early time on smaller scales, or at later times on larger scales. This is not a problem, but rather an advantage of the proposed approach, because it is no secret that the most recent observations have revealed the existence, on both a large and a very large scale, of structures whose existence has become increasingly more difficult to reconcile with conventional models of structure formation. It is obvious to me that such observations, and the bubble-like pattern of the matter distribution in general, can be much more easily explained if we allow for the existence of a parallel distribution of invisible, gravitationally-repulsive, negative-energy matter submitted to mutual gravitational attraction among particles of the same kind.

Before concluding this section, I would like to mention the existence of another astronomical phenomenon which can be expected to occur as a con-

sequence of the presence of very-large-scale overdensities in the distribution of negative vacuum-dark-matter energy. It involves an effect which might be called *repulsive gravitational lensing* and which is merely the counterpart to ordinary gravitational lensing that would be produced when the visible light from a distant source would be gravitationally repelled while it travels through a negative-energy matter overdensity, on its way to our telescopes. While ordinary gravitational lensing produces arcs of light, repulsive gravitational lensing would produce blobs of light, because they would distort the image of the background structures in such a way that the objects observed would appear to be more densely packed in space, behind the invisible negative-energy matter overdensity located in the foreground. In fact, such divergent gravitational lensing phenomena could also be caused by the presence of a positive-energy matter underdensity located between a distant light source and the observer who measures its position in the sky, just like ordinary gravitational lensing could also potentially be enhanced by the presence of underdensities in the negative-energy matter distribution, superposed on the positive-energy objects in the foreground. Such an effect, therefore, would not be easily distinguishable from that which is produced by the presence of voids in the positive-energy matter distribution, as negative vacuum-dark-matter energy overdensities can be expected to exist only inside such voids and only in smoothly distributed form, due to the absence of baryonic negative-energy matter overdensities. This may explain why weak gravitational lensing experiments have not yet revealed the existence of a phenomenon of this kind that could only be attributed to the presence of a gravitationally repulsive matter overdensity.

In face of the mounting difficulties we have encountered in recent years in trying to make sense of a growing amount of unexpected empirical results, I think that the time has come to recognize that simply adjusting the free parameters of the cold-dark-matter model is no longer an adequate approach for addressing the challenges raised by the observed large-scale features of our universe. But even if the words ‘dark matter’ are contained in the name of the currently favored cosmological model, it does not mean that rejecting this model requires completely abandoning the idea that invisible forms of positive energy may play a role in the development of large-scale structures, because it remains that a certain phenomenon attributable to local variations of vacuum energy density can be expected to have consequences similar to those which were once attributed to conventional cold dark matter. Thus,

I believe that what is required to make the current models more acceptable is merely an additional ingredient that would strengthen the gravitational forces responsible for sculpting the large-scale matter distribution in ways which would allow to appropriately describe certain phenomena that would otherwise remain unexplained.

It is merely the fact that a void in the positive-energy matter distribution is expected to exert a gravitational repulsion on the expanding positive-energy matter that surrounds it, on a cosmological scale, that prevents us from drawing the obvious conclusion that unseen matter must be present in the largest voids that may slightly, but very systematically enhance their gravitational repulsion and allow those empty spherical structures to grow to unexpectedly large sizes. The early proposals that the largest voids might have formed as a consequence of explosive processes that would have taken place in the early universe were thus based on the right intuition, but they failed because they did not involve gravitation as the repulsive force. It would therefore be the traditional reluctance to consider the possibility that gravitationally-repulsive matter may exist, as well as the ignorance of the fact that such matter must necessarily be dark and be gravitationally attracted to itself, that would explain the difficulties we experience in trying to make sense of the most recent data regarding the processes that take place in our universe on a very large scale.

4.5 The flatness problem and matter creation

In the introductory section of this chapter I mentioned that there are two broad aspects to what I call the inflation problem, which are the flatness problem and the horizon problem. Here I would like to discuss the first category of issues. Despite the commonly held belief that this problem has been solved by inflation theory, I think that it is still important to understand the difficulties it raises for cosmology, given that the validity of inflation has not yet been definitely confirmed and even if there occurred an initial phase of accelerated expansion, it may not necessarily produce the desired outcome. As I previously mentioned, the flatness problem arises from the fact that the present density of positive matter and vacuum energy appears to be fixed to its critical value, while we have no idea what the constraint is that would require such an extremely precise adjustment of parameters as would have to occur in the early stages of the Big Bang in order to produce the observed

outcome. The problem is that if the faintest of deviation away from a critical rate of expansion had taken place at such an epoch, it would have given rise to a much larger deviation away from flatness at later times, while what we observe is a universe with a positive density of energy that is still critical to a very good degree of precision. The truth, therefore, is that according to current knowledge, the Big Bang model, while mathematically consistent, is nevertheless incomplete, given that the initial conditions, it would seem, cannot be determined by the theory.

Of course, this does not mean that we can't determine a unique rate of expansion at any time in the past by evolving the current state backward in time, which would actually allow us to predict that the density of positive matter, radiation, and vacuum energy $\rho(t)$ was closer to its critical value $\rho_c(t)$ (associated with a density parameter $\Omega = \rho/\rho_c$ equal to 1) in the past. Only, we cannot explain why the current density of positive energy ρ_0 itself is fixed to its critical value $\rho_{c,0}$ to such a high degree of precision. Thus, while relativity theory enables a positive-energy observer to predict what the rate of expansion of the universe was at different times in the past, given the current densities of positive matter, radiation, and vacuum energy, according to the traditional approach this is only true in as much as the rate of expansion at the present time is empirically determined through a measurement of the Hubble constant H_0 , but the model remains well-defined for any value of ρ_0 and H_0 . Yet, I believe that there is much less freedom than is usually assumed in fixing the initial variation of the specific rates of expansion that give rise to the present specific densities of positive- and negative-energy matter. What I will now explain is that despite the conventional assumption to the effect that this initial condition is left unconstrained in the standard Big Bang model (without inflation), there does exist an unavoidable requirement for the current energy density of matter and vacuum to be very precisely equal to the critical value associated with a flat space, from the viewpoint of both positive- and negative-energy observers.

One thing must be clear before we attempt to explain the current flatness of space on the cosmological scale and this is that there is an upper limit to the positive and negative contributions to the density of matter energy. This means that space cannot continue to contract (in the past direction of time) beyond the point at which a maximum amount of matter and radiation energy of positive or negative energy sign is contained in every elementary unit of space. It would be incorrect to assume that the initial value of the density parameter Ω cannot be determined, due to the 'fact' that the initial density

of matter is infinite in the very first instant of the Big Bang. Indeed, from a quantum gravitational viewpoint, there is no time zero at which the density of matter is infinite, only a minimum significant time at which positive and negative energy densities have a maximum, but finite magnitude. Given that, in the context of my interpretation of matter as being equivalent to missing vacuum energy, a maximum value of matter energy density is determined by the natural vacuum-stress-energy tensors associated with the upper limits of the positive and negative contributions to vacuum energy density, then this must be assumed to be the maximum magnitude of the positive and negative contributions to the density of matter energy in the state that emerges from the initial singularity.

What needs to be explained, therefore, is merely why it is that the rate of expansion of space did not begin to differ from its critical value immediately after the universe emerged from this state of maximum positive and negative energy densities that is uniquely determined by the natural scale of quantum gravitational phenomena. The initial positive and negative densities of non-gravitational energy are not arbitrary, but the problem is that there is too much freedom in fixing the early variation of the rate of expansion which determines the average density of matter at all later times. From the conventional viewpoint, it would appear that the early variation of the rate of expansion that gives rise to a flat space at the present time is merely one alternative among an enormous range of possibilities. What I will explain, however, is that, while the current value of gravitational potential energy for the universe as a whole (which is fixed by the present average densities of positive matter and vacuum energy) and the currently observed kinetic energy of expansion (which is determined by H_0) appear to constitute free parameters of the standard model of cosmology, they are not really independent variables in the context where energy must be null for the universe as a whole. In fact, under such conditions, the early variation of the rate of expansion measured by a positive-energy observer must be adjusted not merely to a level of precision that would allow space to keep expanding until the present epoch, but to such an extent that space can be expected to remain perfectly flat on the largest scale for an arbitrarily long time. I will show that this constraint can only be fulfilled when a maximum density of negative-energy matter is assumed to have been originally present in the universe alongside that of ordinary, positive-energy matter.

In the context of the model I have proposed to integrate negative-energy matter to gravitation theory, it may seem like the presence of such negative-

energy matter would change nothing to the conclusion that flat space is an unlikely possibility for the present state of the universe, because a uniform distribution of negative-energy matter exerts no influence on the gravitational dynamics of positive-energy matter on the largest scale, for reasons I have explained in section 2.6. The present specific density of negative-energy matter would in fact be independently subjected to the same excess of freedom as affects that of positive-energy matter, given that the variation of the specific rate of expansion of negative-energy matter is determined only by the density of matter with the same energy sign, and it would appear that this expansion rate could vary as freely as the specific rate of expansion of positive-energy matter initially. In any case, if space was negatively curved from the viewpoint of a negative-energy observer, this would not merely be a consequence of the fact that the energy of matter that determines the expansion rate measured by such an observer is indeed negative, as if negative-energy matter could accelerate its specific rate of expansion through the gravitational repulsion it would exert on itself, like we would expect from a traditional viewpoint, because, as I explained in section 2.4, negative-energy or negative-mass matter does not exert a gravitational repulsion on matter with the same energy sign. Thus, in principle, space could just as well be positively curved and closed from the viewpoint of a negative-energy observer, because the property of gravitational attraction or repulsion is not an absolute feature of matter with a given energy sign. Yet, despite this state of affairs, it turns out that the presence of negative-energy matter is, in fact, required (as I mentioned above) to explain why it is that we are allowed to expect that space should be perfectly flat from the viewpoint of a positive-energy observer.

Although the alternative solution I will propose to the flatness problem is quite simple, it was actually one of the results which I had the most difficulty deriving among those that figure in this report. Part of the difficulty arose from the fact that there are conflicting accounts of what constitute the many contributions to the energy budget of the universe and how their magnitudes may vary as a function of the values assumed by various physical parameters. Thus, while I always had the intuition that, in the context where the presence of negative-energy matter cannot be ignored, a natural solution to the flatness problem might become possible once we recognize the necessity to appropriately apply the principle of energy conservation to the Big Bang, it was not clear which contributions could balance one another out exactly in order to produce a universe out of nothing. But when I finally figured

out what the various contributions to the energy budget of the universe are in the presence of negative-energy matter, and which must be considered independent from which others, and which would need to have the same magnitudes in the initial Big Bang state, then it became clear that, under such conditions, one must actually observe space to expand at precisely the critical rate when we require the energy to be null for the universe as a whole. Before I explain why it is exactly that applying this theoretically motivated constraint may have such far-reaching implications, however, I would like to describe what the motives are that justify assuming that the energy of the universe must, in effect, be null.

I already discussed the importance and the unavoidable character of the constraint imposed by the requirement of relational definition of physical attributes in the preceding two chapters of this report. Basically, what must be understood, concerning the problem at hand, is that the total energy of the universe constitutes one such property which definitely cannot violate the rule that it be characterized in a purely relational way. What is implied by such a requirement is that, even if the Big Bang was not considered to constitute a creation event at which any conserved physical quantity must be created out of nothing, from the viewpoint of an observer of any energy sign, the universe would still need to have a vanishing, total, average energy density. Indeed, one might argue that the requirement of invariance in time of conserved physical quantities does not apply to such a singular event as the Big Bang at which time itself may come into existence, or alternatively that the Big Bang does not even constitute an absolute beginning to time, given that evolution could perhaps be continued to times before the initial singularity if a quantum bounce occurs, as implied by the most promising quantum gravitation theories. But when we recognize the unavoidable nature of the constraint of relational definition of the physical attribute of energy, it emerges that the universe as a whole cannot have a non-zero energy, even if the Big Bang does not constitute a creation event at which any conserved quantity must be created in equal positive and negative amounts.

This conclusion simply follows from the fact that if it was possible to measure a non-zero value of matter energy for the universe as a whole, then this value would have to be either positive or negative and this would allow the particular direction of time relative to which this positive or negative energy would propagate to be singled out as an absolutely defined direction, unless this energy is compensated by an opposite energy of the gravitational field of the universe. A similar condition would also apply to the momentum

of the universe, given that any non-zero momentum of matter arising from a collective motion of positive-energy matter relative to negative-energy matter, would allow to single out a particular direction in space as being that along which this positive or negative momentum is directed, unless it is compensated by an opposite momentum of the gravitational field of the universe, such as would normally need to develop in order to produce such an outcome, whenever there is an absence of collective motion initially (for reasons I have discussed in section 2.11). Here, the fact that there exist both positive- and negative-energy particles propagating forward in time is no different from the fact that there may exist particles with both positive and negative momenta propagating in one and the same direction of space. What this means is that the total energy of the universe, just like its total momentum, must remain null under all circumstances if one is to avoid giving preferred status to one particular direction of space or time that would be significant on the scale of the universe as a whole, as if this direction could be related to some reference point outside that universe, in violation of the requirement of relational definition of physical attributes.

The validity of the above argument is reinforced by the fact that, even in the context where positive and negative energy observers are allowed to experience different rates of expansion of space on the global scale, if the total density of matter energy is null in the initial Big Bang state, it always remains null from the viewpoint of any observer as expansion takes place. Indeed, based on the developments introduced in section 4.2, it would appear that when the average, specific density of negative-energy matter is growing relative to that of positive-energy matter, as a consequence of the emergence of a difference between their specific rates of expansion (the rates of expansion experienced by negative- and positive-energy observers respectively), the ratio of the average densities of positive- and negative-energy matter determined by a positive-energy observer must remain invariant, because the density of negative-energy matter measured by such an observer is modified by the same metric conversion factor which fixes the density of vacuum energy, while the density of vacuum energy must grow in proportion to the magnitude of the divergence between the scale factors experienced by opposite-energy observers. As a result, any variation of the average, specific density of negative-energy matter relative to that of positive-energy matter remains unobservable for a positive-energy observer (even though matter energy can be exchanged with radiation energy as a result of matter-antimatter annihilation). For those reasons, I believe that the commonly held opinion to

the effect that it may not be absolutely necessary to require the universe to have a null value of energy, even when matter is not created out of nothing at the Big Bang and the principle of conservation of energy is not explicitly required to apply, cannot be justified. The fact that, by taking a different stance, I will achieve significant progress in describing the early stages of the universe's expansion will serve, I hope, to vindicate the legitimacy of my viewpoint.

Now, what most people already recognize concerning the energy content of the universe is that, for a flat universe with a zero cosmological constant, the negative gravitational potential energy of positive-energy matter and radiation is balanced by the positive kinetic energy of expansion of this matter. When that is not the case, then an additional amount of energy is present that is attributable to the gravitational field itself (or the curvature of space) and this energy tends to become dominant very rapidly (regardless of whether it is positive or negative) as space expands, because, while the gravitational potential energy of matter decreases in inverse proportion to the volume, the energy associated with the curvature of space decreases as the inverse of the surface enclosing that volume. What may be difficult to understand is the fact that the kinetic energy of expansion is actually a property of the expanding space, which means that it must be considered an energy of the gravitational field itself and not really an energy of matter, despite the fact that the sign of this energy varies as a function of the sign of energy of the observer which is assessing its value. Indeed, the initial value equation for a homogeneous and isotropic universe, which is derived from the general-relativistic gravitational field equations under the condition that energy is conserved for the universe as a whole, is usually written as

$$E = K + V(a) = \left(\frac{1}{a} \frac{da}{dt}\right)^2 + \left(\frac{-8\pi\rho}{3} - \frac{\Lambda}{3} + \frac{k}{a^2}\right) = 0 \quad (4.3)$$

where E is the gravitational energy of the universe, K is the kinetic energy of expansion, and $V(a)$ is the Friedmann potential as a function of the scale factor $a(t)$ in the presence of a cosmological constant Λ for a universe with an average matter density ρ . Here the spatial curvature parameter, which I redefine as $-k/a^2$ and which is always precisely equal to zero for a flat universe, appears as just one particular (reversed) contribution to the Friedmann potential, but, when it is possible to assume that the magnitude of the cosmological constant was negligible initially, this equation can be rewritten

as

$$E_g = K + U(a) = \left(\frac{1}{a} \frac{da}{dt}\right)^2 - \left(\frac{8\pi\rho}{3}\right) = \frac{-k}{a^2} \quad (4.4)$$

which clearly shows that the spatial curvature parameter is the outcome of the imperfect cancellation of the gravitational potential energy of matter by the kinetic energy of expansion.

Thus, whenever the gravitational potential energy of matter $U(a)$ is not matched by a kinetic energy of expansion K that's exactly its opposite, the energy E_g associated with the gravitational field or the curvature of space itself, which is given by $-k/a^2$, is not zero and contributes to alter the expansion rate. If k is positive this excess of gravitational energy is negative, which means that the negative gravitational potential energy of matter contributes predominantly to determine the gravitational field, as must be the case when the source of this gravitational field has positive energy, while when k is negative there is a positive excess of gravitational energy, which means that the positive kinetic energy of expansion (which is also an energy of the gravitational field) contributes predominantly to determine the gravitational field of the universe. The gravitational energy E_g associated with the present value of the curvature parameter $-k/a^2$ must therefore be considered to consist of a residual measure of energy, which could in principle assume any positive, negative, or null value depending on the current value of the scale factor and on whether k is equal to -1 , $+1$, or 0 . There is no *a priori* reason, however, to assume that the measure of gravitational energy associated with the curvature of space on the cosmological scale should be the same for positive- and negative-energy observers at the same epoch, because the kinetic energy of expansion varies as a function of the rate of expansion, which is an observer-dependent quantity in the context where, as I explained in section 2.6, only the average density of positive-energy matter contributes to determine the gravitational field that influences the expansion rate measured by a positive-energy observer, while, in principle, a negative-energy observer could measure different magnitudes for both the average density of negative-energy matter and the rate of expansion it contributes to determine (even before the early annihilation of matter with antimatter), for reasons I discussed in section 4.2.

It must be clear that, even though it is usually assumed that, in a general-relativistic context, the initial value equation expresses the requirement of gravitational energy conservation for the universe as a whole, what the origi-

nal form of the equation really means is that when an additional term, which is provided by the negative of the spatial curvature parameter $-k/a^2$, is added to the equation that would otherwise express the nullity of gravitational energy, then the gravitational energy of the universe can be conserved even in those cases where it would not really be null initially, but it does not really amount to require that the universe comes into existence with zero gravitational energy. What equation (4.4) means is that, once it is assumed that the cosmological constant Λ is negligible initially, then it is only when the free parameter $-k/a^2$ associated with the curvature of space is zero that the positive kinetic energy of expansion K can balance the negative gravitational potential energy $U(a)$ attributable to the presence of positive-energy matter. The true measure of gravitational energy for the universe as whole, therefore, is really the energy E_g which is associated with the curvature of space (which would justify that we refer to this energy as the actual gravitational energy of the universe) and it is only when this energy is null that the gravitational field does not contribute energy on the cosmological scale. But it is usually assumed that this curvature parameter can also be positive or negative and the universe be positively or negatively curved, so that the degree of curvature at any given time would depend on the initial value of the kinetic energy of expansion when the density of matter and radiation was maximum. It must be acknowledged, however, that from the viewpoint of positive-energy observers at least, space does have a flat geometry, to a relatively good degree of precision, and this means that there must be a reason why the curvature parameter has a null value.

I believe that what allows the value of gravitational energy E_g associated with the spatial curvature parameter to be null for an expanding zero-energy universe is the fact that the gravitational potential energy of matter experienced by a positive-energy observer can be arbitrarily large, even when negative-energy matter is present and the total energy of matter itself is null. Indeed, when a large density of negative-energy matter is present initially, a flat universe can actually have zero energy, despite the fact that from a traditional viewpoint it would appear that if the energy contained in the gravitational field we experience was null (as would occur if the negative gravitational potential energy of matter was compensated by the kinetic energy of expansion) the energy of the universe would still be positive (because the energy of matter would not cancel out). It is only from a traditional perspective that it would appear impossible to require our flat universe to have zero energy. In the absence of a large density of negative-energy mat-

ter, the primordial universe would actually need to have a positive curvature in order to have zero energy, because only then could the negative energy contained in the gravitational field compensate the large positive energy of matter (while the gravitational field of a negatively-curved universe would contribute even more positive energy, as the positive kinetic energy of expansion would overcompensate the negative gravitational potential energy of positive-energy matter to provide a positive gravitational energy of curvature E_g). In fact, it seems that it is only for a closed universe that does not expand at all, that the positive energy of matter could be entirely compensated by the residual energy of the gravitational field in the initial state of maximum matter energy density, because according to certain accounts, when the energy density is that high, the gravitational potential energy is actually equal in magnitude to the energy of matter. Again, however, the problem is that the universe is not highly curved, but in all likelihood almost perfectly flat.

At this point, it is important to mention that the idea that the energy of the universe should perhaps be required to be null is not new. Thus, it was once suggested [34] that the universe could fluctuate into existence if the positive energy of matter could be compensated by its negative gravitational potential energy, at least in the very high density of a primordial state. The problem was that it appeared that such a highly curved universe could never be produced as a fluctuation out of nothing, because, if it actually has zero energy, it would only be allowed to expand for a very short period of time before immediately recollapsing back to the vacuum. Creation out of nothing was eventually salvaged from this severe failure by assuming that once in a while inflation may occur when a universe is fluctuating out of the vacuum, which would enable its expansion rate to start growing exponentially, thereby giving rise to a flat space which would keep expanding indefinitely.

I will not immediately discuss any motives we may have to resist appealing to inflation in order to obtain an expanding, zero-energy universe or indeed to solve any other problem in cosmology, but given that, very early on, I chose to explain known facts with principles which are themselves known to be valid with absolute certainty (even if certain consequences of applying those fundamental principles may not yet be recognized as unavoidable), then I was led to explore the possibility that there may exist a more natural way by which the requirement of null energy could be satisfied in the presence of negative-energy matter. In order to proceed in this direction, however, one must first acknowledge that if the negative gravitational potential en-

ergy of positive-energy matter exactly balances the positive kinetic energy of expansion for a flat universe, then this gravitational potential energy cannot also balance the positive energy of matter itself, as earlier proposals required assuming. This does not mean that the magnitude of the gravitational potential energy experienced by a positive-energy observer cannot be equal to the magnitude of positive matter energy initially, only that this is not an appropriate and sufficient condition for obtaining a zero-energy universe. In fact, as I mentioned above, it does appear, according to certain accounts, that in the initial Big Bang singularity (or indeed any other singularity) the positive energy of matter is equal in magnitude to its negative gravitational potential energy and this is precisely the reason why it was so difficult for me to realize that it is not appropriate to merely require the gravitational potential energy to compensate the energy of matter in order to obtain a universe with zero energy.

What I have realized is that, in a zero-energy universe, any residual measure of gravitational field energy associated with the initial value of the spatial curvature parameter $-k/a^2$ determined using the metric properties of space experienced by positive-energy observers must necessarily balance the residual energy of matter obtained by adding the opposite contributions of positive- and negative-energy matter. Now, if the curvature parameter is null, like gravitational energy itself, in the case of a flat universe (for which the kinetic energy of expansion experienced by an observer with a given energy sign precisely balances the gravitational potential energy of matter with the same sign of energy), then it can only mean that, in such a case, the energy of matter must itself add up to zero. Normally that would not be possible, because only an empty universe would have a null, average density of matter energy. But in the presence of negative-energy matter, a high-density universe can actually have a null matter energy, as long as the average densities of positive- and negative-energy matter have exactly the same magnitude initially.

Now, the idea that negative matter energy could compensate positive matter energy, from the viewpoint of a positive-energy observer, may perhaps appear problematic in the context where I have explained (in section 2.11) that those two energies are conserved independently from one another. One must keep in mind, however, that negative-energy matter inhomogeneities, at least, do have an effect on positive-energy matter and if it may seem that the presence of a homogeneous distribution of negative-energy matter is without consequences for positive-energy observers, given that it exerts

no influence on the expansion rate they experience, it is merely because a uniform distribution of negative-energy matter is equivalent to a void of universal proportions in the distribution of positive vacuum energy. But this void does exist, from the viewpoint of a positive-energy observer, and it does constitute missing positive energy, even if it has no effect on positive-energy matter on the global scale. Therefore, if this missing energy is to constitute an objective fact (which can be communicated between opposite-energy observers) then it must contribute to the matter energy budget, even from the viewpoint of a positive-energy observer.

We may, in fact, consider that the way by which negative-energy matter does contribute to determine the gravitational field experienced by a positive-energy observer on a global scale is precisely by reducing the energy of matter, which allows the magnitude of the energy of the gravitational field to itself be reduced when the total energy is required to be null for the universe as a whole. Anyhow, either the negative energy of matter remains totally uncompensated by the energy of the gravitational field associated with a positive-energy observer, in which case this gravitational field energy would alone need to compensate the positive energy of matter, which would imply that there can be no contribution by a positive kinetic energy of expansion (so that the universe should not expand at all), or both the positive and the negative portions of the energy of matter must be compensated by the same gravitational field energy, in which case expansion is allowed to occur, but we must explain why the total, average density of matter energy was initially so close to zero that the energy of the gravitational field (associated with the global curvature of space) was itself required to be perfectly null. Clearly, the second option is the only one that *could* be viable, and therefore I will concentrate on explaining why the total, average energy of matter which balances the energy of the gravitational field for the universe as a whole cannot be as arbitrarily large as one might expect.

It must be clear, first of all, that even if we were to assume that there were as many positive-action particles as there were negative-action particles in the initial state of maximum matter energy densities (as may be required in order that the initial matter distributions be as homogeneous as they are observed to be) in principle it would still be possible for the magnitude of the average density of positive matter energy to be larger or smaller than that of the average density of negative matter energy, even in a universe with zero energy, when thermal energy would differ for those two matter distributions. In the absence of an appropriate constraint this would, in effect, be allowed as long

as the difference between the magnitudes of the positive and negative energies of matter is compensated, from the viewpoint of a given observer, by the energy of the gravitational field associated with the curvature of space, which is determined by the rate of expansion measured by that observer (because while negative matter energy can compensate positive matter energy, only the gravitational field experienced by a positive-energy observer can contribute to cancel out any non-zero, average energy density of matter determined by such an observer). Under such conditions, the magnitudes of the positive and negative contributions to the energy of the universe could be equal initially, even if the average densities of positive and negative matter energy were not themselves equal, and therefore the total energy could in principle be null regardless of the amount of energy contained in the gravitational field. It may, therefore, seem like a condition of null energy for the universe as a whole does not provide sufficiently strong a constraint to necessarily give rise to a flat universe. But, in fact, I came to realize that this condition is much more constraining than one may expect for gravitational energy and the rate of expansion and that it actually allows to predict that the geometry of our universe can only be observed to be flat on the largest scale.

It is important to point out, however, that the nullity of the energy of matter cannot be fixed as an independent consistency requirement, because that would require assuming that there cannot even be local fluctuations away from this zero energy for matter, while this is required in order to explain the observed inhomogeneities present in the initial distribution of matter energy on a scale larger than the cosmic horizon. But in the absence of such a constraint, local fluctuations above or below the average zero value of matter energy density could, in effect, be present in the initial Big Bang state, even if the average densities of positive and negative matter energies were required to cancel out on the global scale, so as to allow the zero-energy universe to have a flat geometry, as long as there is, in effect, as much overdensity as there is underdensity in the positive- and negative-energy matter distributions on a sufficiently large scale. Such fluctuations in matter energy would simply need to be compensated by *local* variations in the kinetic energy of expansion, above or below the value associated with a critical expansion rate. Thus, even when the density of positive- and negative-energy matter particles is maximum, as must have been the case in the first instants of the Big Bang, fluctuations in the energy density of matter would be possible for a zero-energy universe, given that local variations of gravitational field energy could compensate local variations in the thermal energy of matter

particles and maintain the total energy of matter and gravitational field to a null value. If there was less positive than negative matter energy in a certain location initially, then there would simply need to be more positive gravitational energy and therefore more positive kinetic energy of expansion from the viewpoint of positive-energy observers and less negative kinetic energy of expansion from the viewpoint of negative-energy observers.

Thus, inhomogeneities could be present in the initial distribution of matter energy, even if the density of positive-energy particles (the number of positive-action particles in a volume of space) was required to everywhere equal that of negative-energy particles, so as to allow a homogeneous initial matter distribution, because, *locally* at least, the nullity of energy can arise from a compensation between the energy of matter and the energy of the gravitational field. Indeed, a local variation in the energy of the gravitational field (attributable to a local variation of the kinetic energy of expansion above or below the value associated with a critical expansion rate) can be made to compensate any local difference between the magnitude of the density of positive matter energy and that of negative matter energy, just like the global measure of gravitational field energy which is attributable to the difference between the observer-dependent gravitational potential energy of matter and the observer-dependent kinetic energy of expansion could in principle compensate any difference between the magnitude of the *average* cosmic densities of positive and negative matter energy. However, in section 4.9 I will explain that a certain unavoidable constraint actually limits the amplitude of those fluctuations in the initial Big Bang state and therefore it cannot be expected that there would occur large deviations from zero gravitational energy locally if this condition is also obeyed globally.

But even if local fluctuations in the density of matter energy are clearly unavoidable, it remains to explain why it is that such a compensation of matter energy by gravitational energy is not allowed to take place on a global scale, as required if space is to be flat for the universe as a whole. Indeed, as I mentioned above, if a residual gravitational energy associated with the spatial curvature parameter $-k/a^2$ could also compensate a difference in the magnitude of the initial, average densities of positive and negative matter energy on a global scale, then it should be possible for the magnitude of the kinetic energies of expansion experienced by positive- and negative-energy observers to be larger or smaller than the magnitude of the gravitational potential energies of their associated matter. Under such conditions, the rates of expansion would no longer need to be critical, even in a zero-energy universe. It is cer-

tainly true that a homogeneous distribution of negative-energy matter exerts no influence on the specific expansion rate of positive-energy matter which determines the kinetic energy of expansion measured by a positive-energy observer, but this is significant merely in the sense that only the energy of the gravitational field perceived by a positive-energy observer can contribute to the energy budget that must add up to zero on a global scale, from the viewpoint of such an observer. For reasons I previously mentioned, it is still necessary to assume that both the positive and the negative energy of matter contribute to the total value of energy measured by a positive-energy observer.

What must be clear also is that, in the context where the energy of the universe is required to be null, if space was positively curved and closed from the viewpoint of a positive-energy observer, it would need to be negatively curved and open from the viewpoint of a negative-energy observer. Indeed, the gravitational field of a universe that would be positively curved, from the viewpoint of a positive-energy observer, would have a negative energy and could, therefore, only compensate an excess of positive matter energy (through a reduction of the positive kinetic energy of expansion). But while it is true that, even from the viewpoint of a negative-energy observer, an excess of positive matter energy would require the contribution of a gravitational field with negative energy, such a gravitational field would be associated not with a smaller positive kinetic energy of expansion, but with a larger negative kinetic energy of expansion and a higher than critical expansion rate, which would actually give rise to an open universe. If the total energy of matter was instead negative initially (before the early annihilation of matter and antimatter took place), as would occur if negative-energy matter particles contributed more energy than positive-energy matter particles on the average, then the opposite would be true and the universe would need to be closed from the viewpoint of a negative-energy observer and open from that of a positive-energy observer. Now, while those two mutually exclusive configurations may appear to merely consist of two additional possibilities, no different from the case where the average density of matter energy happens to be null initially, just like the energy of the gravitational field, there is actually a very important distinction between the case of a flat universe and that of the curved space configurations. This essential difference has to do with the fact that, in the case of a flat space, the universe would be open from both the viewpoint of a positive-energy observer and that of a negative-energy observer, while in all the other possible cases it seems that

the universe would need to be open for an observer with a given energy sign and closed for an observer with the opposite energy sign.

I believe that if the average density of positive matter energy must exactly compensate the average density of negative matter energy in the initial Big Bang state, even though local fluctuations away from the zero energy of matter are allowed to be present to a certain extent (as long as they are compensated by opposite local fluctuations in gravitational energy), it is precisely because, in the absence of any other contribution to the energy budget, if matter energy was not null, then space could not be flat and open from the viewpoint of all observers. If an excess of positive or negative gravitational energy was allowed to compensate an excess of negative or positive matter energy (respectively) on a global scale, then this excess gravitational energy would give rise to a universe which would be open for one observer and closed for an observer with opposite energy sign. But given that the difference between the volume of a closed universe and that of an open universe would in principle be infinite, it follows that such a configuration would be characterized by an arbitrarily large positive or negative density of vacuum energy. Indeed, from the viewpoint of the developments discussed in section 4.2, it would follow that if gravitational energy was negative and the universe was closed from the viewpoint of a positive-energy observer and open from the viewpoint of a negative-energy observer, as would appear to be required if it is to compensate a positive total density of matter energy, the density of vacuum energy should be positive with a maximum amplitude, while if the opposite was true and gravitational energy was instead positive, as would appear to be required if it is to compensate a negative total density of matter energy, then the density of vacuum energy should be set to its maximum negative value right at the Big Bang¹².

The problem is that the maximum positive or negative value of the cos-

¹²In such a context, it should be clear that it cannot be assumed that the energy of matter in a universe with a non-zero curvature of space is compensated by an opposite energy that would be contained in the vacuum as a result of this curvature. This is not merely a consequence of the fact that under such conditions vacuum energy would actually have the same sign as that of the overabundant matter, it is also unavoidable in the context where the portion of vacuum energy associated with the cosmological constant gives rise to its own contribution to gravitational potential energy, while this contribution is known to always compensate that which is contained in the uniform portion of vacuum energy itself, so that the vacuum and its gravitational field together contribute nothing to the energy of the universe from either the viewpoint of positive-energy observers or that of negative-energy observers.

mological constant which would be associated with *any* such configuration would appear to forbid the emergence of an observer, because even if the gravitational force exerted by the cosmological constant on the specific expansion rates generally contributes to reduce the magnitude of the average density of vacuum energy, if this magnitude had been maximum right when the expansion process began, then it could never have been reduced to a level compatible with the emergence of an observer before the average matter density itself would have become too low to allow for the development of structures (which we may assume to be essential for the existence of an observer). Only a universe with precisely balanced initial contributions to the energy of matter and therefore, also, to the energy of the gravitational field, is allowed to be experienced as a long lasting process by a physical observer that is part of that universe, when it is appropriately required that the universe itself has null energy. It is only when space is flat on a global scale and there is no energy in the gravitational field that the magnitude of the density of vacuum energy can be different from its maximum theoretical value initially. But given that this is required if an observer is to be present at some point in the universe to measure any value of gravitational energy, then one must conclude that the kinetic energy of expansion determined by a positive-energy observer would always precisely compensate the gravitational potential energy attributable to positive-energy matter and the same would be true (independently) for negative energy matter and its gravitational field, from the viewpoint of a negative energy observer. What I'm suggesting is that this is allowed to occur in the case of a zero-energy universe when the energy of matter itself is null in the very first instants of the Big Bang (before the annihilation of most baryons with their antibaryon counterparts), as becomes possible in the presence of negative energy matter.

When this is properly understood, it becomes clear that the 'extra' principle which would allow to fix the expansion rate of our universe to its critical value is nothing else but the requirement of relational definition of physical attributes, which requires the sum of all energies to be null, for the universe as a whole, from the viewpoint of both positive- and negative-energy observers. In the context of the generalized gravitation theory introduced in the second chapter of this report, and given the interpretation that was proposed in section 4.2 for the vacuum-energy term, this constraint actually allows to determine which solution of the gravitational field equations is the appropriate one for a description of the expanding universe. It is, therefore, by applying this very basic principle, in the context where it is recognized that negative-

energy matter must also contribute to the universe's initial energy budget, that it becomes possible to explain not only why there is expansion, but why it is that the rate of this expansion is still critical, even long after the Big Bang. Space is flat and the rate of expansion remains critical, because the universe must be open from both the viewpoint of positive energy observers and that of negative energy observers and the precision with which the initial rate of expansion was adjusted to its critical value is merely a reflection of the exactness of this requirement¹³.

It must be clear, however, that in the context where the initial, average density of negative-energy matter can be reduced to a greater extent than that of positive-energy matter following the annihilation of baryonic matter and antimatter that takes place early on, during the Big Bang, it is possible for the average value of vacuum energy density, or the cosmological constant, to grow from its initial zero value toward a larger, positive value during the matter-dominated era, because under such circumstances the rate of expansion experienced by positive-energy observers is reduced more rapidly than that which is experienced by negative-energy observers, due to the larger gravitational pull exerted by positive-energy matter, which allows the scale factors experienced by opposite-energy observers to diverge. But it is not to be expected that this divergence could develop to an arbitrarily large magnitude, because the weak anthropic principle also forbids the cosmological constant from becoming so large, as a result of this divergence, that it would no longer be compatible with the presence of a (positive-energy) observer at the present time.

What must be retained from all this is that, if it was not for the fact that the presence of a homogeneous distribution of negative-energy matter exerts no influence on the expansion rate of positive-energy matter (as explicitly stated in the formulation of principle 6 from section 2.14 and for reasons I

¹³It must be noted that the same constraint allows one to expect that there is no difference between the average states of motion of positive- and negative-energy matter on the largest scale that could have given rise to a non-zero momentum for the universe as a whole, because such a momentum for matter would need to be compensated by an opposite momentum of the gravitational field, as I previously mentioned, and if the gravitational field had non-zero momentum on a global scale, it would also need to have non-zero energy and this is not possible for a universe with flat geometry. It is therefore possible to predict that there must be an exact correspondence between the global inertial reference systems experienced by positive- and negative-energy observers in Einstein's elevator experiment far from any local matter inhomogeneities, so that under such conditions positive- and negative-energy bodies should have the exact same acceleration.

have explained in section 2.6), then, even if the total energy of matter was null initially, it would not be possible to conclude that the observed expansion rate must be the critical rate associated with the density of positive matter energy, because under such conditions the gravitational potential energy of matter that would need to be balanced by the kinetic energy of expansion would actually be zero (because the total density of matter energy that would determine the strength of the gravitational field would itself be null initially), which means that the kinetic energy of expansion would also need to be zero and the universe should not expand at all. But if the requirement of energy conservation did not apply to the gravitational field and the universe did expand, as we would normally assume, then the expansion rate would not be submitted to any deceleration and the universe would explode like a negatively-curved universe with a null matter density. The independence of the expansion rates of positive- and negative-energy matter from the presence of matter with an opposite energy sign, which follows from my description of negative-energy matter as consisting of voids in the positive-energy portion of the vacuum, is therefore an essential ingredient of the alternative solution to the problem of flatness that is proposed here. This condition is especially constraining in the context where the initial matter distribution must be highly homogeneous on a macroscopic scale (for reasons I will explain in section 4.9), so that there cannot even exist significant local perturbations of the rate of expansion of matter with a given energy sign by matter with an opposite energy sign on a large scale.

It is only after I realized that, from the viewpoint of a positive-energy observer, the presence of negative-energy matter does not contribute to determine the gravitational potential energy of the universe (which in the case of a universe with an overall flat geometry is compensated by the kinetic energy of expansion), that I was able to understand that, despite what is usually assumed, it is, in effect, not only the current *variation* of the specific rate of expansion of positive-energy matter which is determined in part by its energy density, but actually also the current specific rate of expansion itself. It took me a certain time to recognize that the variation of the rate of expansion must indeed be considered to depend on the density of (positive) matter energy, as most people may consider obvious (but unlike one would perhaps expect for a universe with null matter energy), yet my questioning has allowed me to realize that the relation which exists between the rate of expansion and the density of matter energy is actually much more constraining than is usually assumed.

As a result, I'm allowed to conclude that, even in the absence of inflation, it is not necessary to assume that the present density of positive matter and vacuum energy is critical purely for aesthetic reasons, because, in fact, it is possible to explain why the universe is so perfectly balanced, when one recognizes the necessity to properly apply the requirement of relational definition of physical attributes to the energy of the universe as a whole, which requires it to remain null even when it is hypothesized that bidirectional time may be continued past the initial Big Bang singularity, following a quantum bounce, so that matter does not need to be created out of nothing. Furthermore, it appears that it is the fact that an observer can only measure a value of vacuum energy density that is compatible with the conditions of her own existence, that explains that it is not merely the total energy content of the universe that is observed to be precisely null, but to a very good degree of precision, also, the total energy of matter in the initial Big Bang state. In any case, it is now possible to understand that the flatness of space is not a mere possibility that emerged as a byproduct of an uncertain process of inflationary expansion, but rather constitutes a basic consistency requirement that must be satisfied by any viable cosmological model.

When I will discuss the horizon problem, in section 4.9, I will explain what justifies assuming that the distribution of positive and negative matter energies in the initial Big Bang state was sufficiently homogeneous that no macroscopic event horizons (understood as any event horizon larger than that which is associated with an elementary black hole) would be present on any scale. But it can already be appreciated that, in the context where the initial distribution of matter energy is uniform to a very high degree and the *local* rates of expansion of positive- and negative-energy matter only vary in such a way as to allow the kinetic energy of expansion to compensate any difference between the amplitudes of their opposite energy densities, as I'm here assuming, then the expansion of space must remain almost perfectly isotropic on the largest scale; which certainly constitutes an appropriate conclusion from an observational viewpoint. The fact that, from a traditional perspective, such an outcome would only be allowed to happen as the consequence of an early phase of inflationary expansion, therefore, no longer constitutes a decisive argument in favor of inflation theory, because from my perspective an initial period of accelerated expansion is no longer necessary to produce such an outcome.

In the context where the sum of all energies which can be measured by a

given observer must be null, for the universe as a whole, it also emerges that the often met remark, to the effect that the observed equilibrium between open and closed universe is improbable, as it requires a delicate balance between the kinetic energy of expansion and the gravitational potential energy of matter, is irrelevant, because, on the basis of the hypothesis that an observer must be allowed to exist in the universe to determine this expansion rate, such an observation, far from being improbable, is actually unavoidable. The tentative solution to the problem of flatness provided by inflation, therefore, appears to simply be unnecessary, because even when the initial density of positive-energy matter is very high, the energy of the gravitational field is required to be null in a zero-energy universe, which means that the universe must necessarily have a critical density of positive matter energy and enough kinetic energy to keep expanding forever (at an ever slower rate), even if a non-zero cosmological constant develops later on, as a result of a variation of the ratio of positive to negative matter energy densities attributable to the annihilation of matter with antimatter, because even a negative average value of vacuum energy density would have a tendency to be reduced under its own influence, instead of allowing space to recollapse from the viewpoint of a positive-energy observer.

It would therefore appear that the idea that the initial push of inflation is necessary to explain that there is any expansion at all is incorrect, because, if a state of maximum energy density must be present in the universe at some point (if such a condition needs to be satisfied independently from whether there is expansion or not, as I will suggest in section 4.9), then expansion at a proportionately high rate does actually become an absolute necessity, if energy is to be null at all times for the gravitational field, independently. But even if it is the presence of an observer that requires this latter condition to be satisfied, this does not mean that it is necessary to appeal to the anthropic principle in order to explain the fact that the universe has not yet recollapsed, because what is required by the presence of an observer is not merely that the universe is still expanding at the appropriate rate for life to exist, but really that it has a perfectly null cosmological constant initially (which can only happen when its rate of expansion is critical). It is true, though, that if the universe did not expand at a rate that would have been too large or too small to allow for the emergence of an observer, it is not only because the cosmological constant was null initially and the rate of expansion critical, but also because the average density of vacuum energy did not later grow to a much larger value that would have accelerated or decelerated the critical

rate of expansion to such an extent that it would have become incompatible with the presence of an observer.

But in the context where bidirectional time may extend past the initial singularity, one would also need to recognize that it is not necessary for matter to be created out of nothing, because matter and radiation could have been present before the Big Bang that would have been submitted to a quantum bounce as a result of the contraction of space that would have taken place in the future direction of time. In such a case, there would be no meaning to ask how it is that matter was created, because matter simply exists and is not produced by the Big Bang. In fact, in chapter 5 I will explain that there may be good reasons to believe that this persistence in time is actually an essential requirement and that it may need to be extended in a certain, more unexpected way, in order to avoid the hypothesis that matter must have been created out of truly nothing. From that viewpoint, it would appear that it is only when we ignore the limitations imposed on the growth of positive and negative energy densities in a quantum gravitational context, that a problem may arise with the fact that matter appears to be present in the very first instants of the Big Bang, even if it cannot be created out of nothing. Ironically, it is precisely because we assume the existence of an early phase of inflationary expansion that must leave the universe totally empty, that we need to justify the presence of matter in our universe, by assuming that it was created at a later time by a process of reheating. But given that inflation may no longer be required to explain flatness itself, then it is certainly not inadequate to conclude that there may, after all, be no substance to the problem of matter creation. The idea that only inflation allows to explain the relatively large ‘initial’ density of positive-energy matter would then be incorrect, because, in fact, the hypothesis that there occurred an early phase of inflationary expansion is precisely what makes it more difficult to explain the existence of a hot Big Bang.

Now, if there actually exists a history unfolding past the state of maximum matter density, then it is necessary to assume that the same density of matter energy must have existed in the moments immediately *preceding* the singularity, as existed in the moments immediately following it. This means that the expansion rate, following the quantum bounce in the *past* direction of time, must be as large as it was in the moments immediately following the singularity in the future direction of time (from the viewpoint of both positive- and negative-energy observers), given that this expansion rate must be critical if energy is to be null on the other side in time of the maximum-

density state as well. Therefore, it seems that the conditions necessary for the existence of an observer may also exist in that portion of history. In fact, the initial expansion rates can only be equal on both sides in time of the singularity under the condition that the distribution of matter energy in the ‘final’ state which would be reached while space collapses in the *future* direction of time, in that unknown portion of history taking place before the Big Bang, is as homogeneous as the distribution of matter that provides the ‘initial’ boundary conditions for the current one, because otherwise macroscopic inhomogeneities could potentially survive the quantum bounce that would influence the local rates of expansion. In section 4.9 I will explain why this hypothesis is appropriate and it will then be clear that it does not even need to be confirmed by making use of the detailed mathematical framework of a quantum theory of gravitation.

It is remarkable that despite our ignorance of the exact nature of the laws which apply at the Planck time, it is nevertheless possible to predict what the exact variation of the rate of expansion of the universe was when the average densities of positive and negative matter energy were maximum. But it is also possible to predict that regardless of what happens to the ratio of the average densities of positive- and negative-energy matter (as a result of matter-antimatter annihilation in particular), both the average, specific density of negative matter energy plus vacuum energy and the average, specific density of positive matter energy plus vacuum energy must remain critical if they originally were, given that a flat geometry is the one configuration whose radius of curvature does not change with time. Indeed, when the cosmological constant grows from its initial zero value into a larger positive value, there is more positive vacuum energy to accelerate the specific rate of expansion of positive-energy matter at later times, but this additional positive energy also contribute to the total density of energy that determines the curvature of space experienced by a positive-energy observer, which means that this density remains critical if it initially was and the same is true for a negative-energy observer. This is allowed as a consequence of the fact that the uniform portion of vacuum energy is conserved independently from the energy of matter and can actually be created, even when it does not exist initially, because it is compensated by an associated variation of gravitational potential energy which, under such conditions, can actually grow (reach larger negative or positive values) while space is expanding, exactly as would occur during a hypothetical phase of inflationary expansion.

At this point, one may recall the conclusion I arrived at in section 4.3

to the effect that the amount of positive-energy dark matter (like that of negative-energy dark matter) cannot be assumed to rise on the global scale, despite the fact that the density of vacuum-dark-matter energy should grow locally, along with the inhomogeneity of the visible matter distribution. In light of the developments introduced in this section, it would appear that this conclusion is fully justified, because if the amount of positive-energy dark matter was allowed to grow in such a way, then, in the context where the universe was flat and the matter distribution highly homogeneous in the initial Big Bang state (or in the state immediately following inflation), if the growth of inhomogeneity that follows was to contribute to increase the total amount of positive vacuum-dark-matter energy, this would slow down the rate of expansion determined by a positive-energy observer, because only the positive density of vacuum-dark-matter energy would exert an influence on this expansion rate. As a consequence, the universe would acquire a positive curvature on the cosmic scale (while a similar phenomenon would be experienced by negative-energy observers). But this means that the energy of the gravitational field of the universe would become negative, while the energy of matter and radiation could still be null in principle (as there may occur a similar growth of negative vacuum-dark-matter energy density). As a result the condition of null energy would be violated for the universe as a whole. We may, therefore, conclude that it is, in effect, an absolute requirement for vacuum dark matter to already be present in homogeneously distributed form (still distinct from the uniform portion of vacuum energy associated with the cosmological constant, which obeys a different equation of state), before it begins to grow locally, as a result of the formation of inhomogeneities in the matter distribution.

To summarize, we are in a situation where both the magnitude of the average, initial density of positive matter energy determined by a positive-energy observer and the magnitude of the average, initial density of negative matter energy determined by a negative-energy observer are fixed to the maximum theoretical value that is determined by the natural scale of quantum gravitational phenomena. However, it is not just possible, but actually unavoidable, for those two opposite energy densities to be equal at the exact same time in the initial Big Bang state, from the viewpoint of any observer (made of matter with either a positive or a negative energy sign), in a zero-energy universe, given that the energy of the gravitational field must then itself be null and space be flat from the viewpoint of both positive- and negative-energy

observers, in the context where opposite-energy observers would necessarily measure opposite values for the non-vanishing gravitational energy of the universe if matter-energy was not null initially, while this would result in opposite space curvatures, which would give rise to a maximum average magnitude of vacuum energy density that would be incompatible with the emergence of an observer at later times. But this condition of null gravitational energy can only be satisfied when the kinetic energy of expansion measured by an observer with a given energy sign precisely balances the opposite gravitational potential energy attributable to matter with the same energy sign and this is what explains that space still expands at a critical rate long after the Big Bang.

The problem that there was with the traditional approach is that, if we required energy to be null for the universe as whole, we could not balance the very large positive energy of matter (characteristic of quantum gravitational phenomena) that existed initially in our flat universe, so that it always appeared inappropriate to try to justify the flatness of space as being the consequence of a condition of null energy that would apply to the Big Bang, despite the fact that gravitational energy itself really is null for a flat universe (given that the kinetic energy of expansion is the exact opposite of the gravitational potential energy of matter). This is the reason why we failed to understand that applying a condition of zero energy to the universe could actually provide the basis for an explanation of the flatness of space that does not require assuming that the null energy of the gravitational field determined by a positive-energy observer is a mere coincidence or an outcome of inflation.

4.6 The problem of time asymmetry

It is remarkable that, at this point into my discussion, I have already been able to provide independent solutions to two of the worst fine-tuning problems of cosmology, guided merely by an unwavering confidence in the validity of well-known physical principles. It is significant, also, that both the solution to the cosmological-constant problem and that which was proposed to the flatness problem involved considering the balancing effects of negative-energy matter in order to provide additional constraints on the values of physical parameters. But before I can address other aspects of the inflation problem, it will be necessary to delve a little deeper into what really constitute the

many facets of the problem of time asymmetry from a classical viewpoint. This will allow me to properly identify the nature of the deep contradiction that still dwells at the heart of theoretical physics, as a result of the apparent incompatibility between the time-symmetric laws of classical mechanics and particle physics and the unidirectional laws of thermodynamics and statistical mechanics.

Before engaging in a discussion of the problem of time asymmetry what one must first decide is whether irreversibility is real, or whether it is a mere consequence of the way we describe the state of a system. It has been argued, in effect, that it is only as a consequence of adopting a particular coarse-graining and due to the choice that is made regarding what details of the microscopic state of a system are to be ignored, that irreversibility occurs. If that was the case, then the continuous increase of entropy, which under certain conditions appears to characterize the evolution of physical systems with a large number of microscopic degrees of freedom, would be a purely subjective notion, significant merely in the context where there are practical limitations on our ability to perceive the evolution of a physical system down to its most intricate details. Under such conditions, even if entropy (as a measure of the number of possible, distinct, microscopic states of a system that are compatible with an appropriate choice of observable macroscopic parameters) was to vary, the changes which are taking place would have no fundamental significance and the observation of certain regularities regarding entropy growth would not require explanation, given that the quantity involved would merely be a subjective notion. But despite the fact that this idea is still quite popular among those who have not seriously studied the question of the origin of time asymmetry, it is no longer viewed by most specialists as an appropriate solution to the problem of the origin of irreversibility, but rather as an attempt at easily disposing of the problem without really explaining anything.

It was pointed out by Roger Penrose that the growth of entropy involved in most irreversible thermodynamic processes is so large that it is only marginally dependent on the choice of coarse-graining. Thus, it appears that the degree of *appropriateness* of any particular coarse-graining itself varies dramatically in the course of certain processes which are occurring all the time in our universe. The truth is that, even if we were to follow the detailed evolution of all the microscopic physical parameters of a large system in a non-equilibrium state, certain aspects of this evolution could still be characterized as unidirectional. What this means is that we are not just

shuffling an initially well-ordered deck of cards (to use a simple analogy) which would merely be losing a subjective amount of structure. When we are considering an ordinary deck of cards, all configurations are equivalent, despite the particular significance we attach to the ‘ordered’ configuration. But in our universe the changes which are taking place when entropy is observed to be growing can be characterized in a more objective way, due to the nature of that portion of entropy that is attributable to the gravitational field. Indeed, as I have explained in section 3.10, the measure of entropy associated with black-hole event horizons does not grow merely as a consequence of adopting a certain *arbitrary* definition regarding what parameters should characterize the macroscopic state of such a system and what information remains unavailable, and therefore it gives rise to a less subjective notion of irreversibility. Another distinction of the evolution which is actually taking place on a macroscopic scale in our universe is that the probability to return to a former state of lower entropy never stops diminishing, because the entropy is in principle allowed to grow without limit.

It must be clear, though, that this is not just a consequence of the expansion of space. It was once suggested, in effect, that the growth of entropy associated with all irreversible processes could be a consequence of universal expansion, given that the thermodynamic arrow of time is oriented in the same direction as what is sometimes called the cosmological arrow of time, which is merely the direction of time in which space is expanding. But it was later pointed out that this assumption is inappropriate, because in such a context one would need to assume that the arrow of time should immediately reverse under conditions where space would begin contracting, while there is no independent motive to justify that conclusion. Indeed, the expansion of space is a global phenomenon, while an expanding gas in a container is a local phenomenon which we have no reason to expect would be so drastically affected by what happens to the relative motion of distant galaxies as to start behaving anti-thermodynamically and retract into a smaller volume the moment space would begin contracting on a global scale. This conclusion is certainly appropriate, given that if we were to assume that space contraction alone is sufficient to give rise to a reversal of the arrow of time then we should probably also have to assume that the thermodynamic arrow of time reverses in the presence of a strong enough, attractive, local gravitational field, while of course there is no evidence at all that this is happening.

It is usually understood, however, that while we are allowed to consider entropy as missing information, an objective characterization of temporal ir-

reversibility does not require assuming that information is actually vanishing from reality when entropy is rising. It is certainly true that, when the exact evolution of a system that is not in a state of thermal equilibrium cannot be followed down to its most intricate details, we may lose sight of information concerning its exact microscopic state and therefore more information than is available afterward may be needed to describe it. But if ignorance is growing, it is only because the macroscopic parameters we use to describe the state of a system are leaving aside an increasingly larger portion of the information that would be required to accurately describe its exact microscopic state. Thus, even if certain physical parameters which allow to objectively assess the growth of entropy evolve irreversibly, the amount of structure present on a microscopic scale remains unchanged as those transformations are taking place. It is simply the fact that, regardless of how well chosen they are, *macroscopic* parameters are increasingly less efficient at providing a full description of the structure contained in the exact microscopic state of our universe, that makes it look like information is being lost when the number of microscopic states which can potentially be occupied is growing with time.

In other words, it is merely the difficulty to keep track of all the changes taking place in the most detailed description of the state of a system that is growing with time in an irreversible way, but no information, or no microscopic structure is really vanishing in the process. When one recognizes that there does exist a minimally coarse-grained definition of the state of a system associated with what would be a maximum level of knowledge of its microscopic configuration (regardless of whether this knowledge can actually be obtained or not by a given observer at a specific moment), then one has no choice but to also recognize that it provides a measure of information that is unchanging. In the next section, I will show how certain usually unrecognized variations in the amount of information required to describe the exact microscopic state of the gravitational field are crucially involved in allowing information to be conserved, even when black holes are involved and the growth of entropy constitutes a more objective change. But it is already possible to acknowledge that the conclusion that entropy growth does not require the minimally coarse-grained measure of information to vary is appropriate from a theoretical viewpoint, because the conservation of information is a requirement of quantum unitarity (or of Liouville's theorem in a classical context), as I have mentioned in section 3.10.

Now, if entropy is indeed increasing in the future, from the viewpoint of an appropriately defined choice of coarse-graining, then it means that entropy

was definitely smaller in the past. What is deduced from observations, in fact, is that entropy continuously decreases in the past, in every place we look and as far back in time as we can probe. This is a condition that is far more constraining than simply assuming that the universe is not in a state of thermal equilibrium at the present time, which would certainly also allow entropy to grow larger in the future. What we might be justified to expect, in effect, is that entropy should rise in the past, just as it does in the future, given that it is not already maximum at the present moment. This would appear to be implied by the fact that there is a higher probability that such states be reached as evolution takes place randomly, because there is a much, much larger number of allowed microscopic states compatible with a condition of higher entropy than there are microscopic states compatible with a condition of lower entropy. Only for an isolated system, with a finite number of microscopic degrees of freedom, would there be a chance that evolution could momentarily take place toward a lower entropy state as a mere statistical possibility. Such fluctuations would not constitute violations of the second law of thermodynamics, given that this law is probabilistic in nature. Thus, we may consider that the evolution we observe to be taking place in general in the future direction of time is in line with expectations arising from both classical and statistical mechanics.

The real problem is with the past. Due to the time-symmetric nature of fundamental physical laws it would appear, in effect, that when a macroscopic physical system with many independent microscopic degrees of freedom evolves in the past direction of time, starting from a present non-equilibrium state of relatively low entropy, its entropy should grow (regardless of the details of its microscopic configuration) for the exact same reason that we expect its entropy to grow in the future, when evolution occurs in a random way. But in our universe, entropy was clearly not larger in the past than it now is and the truth is that there is no evidence from astronomical observations that any large-scale, entropy-decreasing phenomena has ever taken place and no written account of any person having ever observed any significant departure from constant, or continuously increasing entropy at any occasion in our entire history. Thus, while we can determine the probability of the statistically significant properties of future configurations from a knowledge of the current state of a system, the probability of past configurations cannot in general be appropriately estimated based on that same knowledge. In fact, even if entropy was continuously increasing in the past, from its present non-maximum value, we may still have a problem, because

from the forward-time viewpoint the evolution that would have taken place in the past would have occurred with diminishing entropy in the future and this aspect would also be unexplained, unless we are dealing with a *momentary* fluctuation. Thus, it seems that what must be explained is not merely why it is that entropy does not increase in the past, but why it is not already maximal and unchanging in both the past and the future.

It was suggested that the conclusion that entropy should increase in the past may not be valid, because even a macroscopic system with a large number of independent microscopic degrees of freedom could perhaps be so carefully prepared that it would be allowed to retrace an unnatural entropy decreasing evolution as it evolves backward in time. Thus, it was argued that it is the details of the present microscopic state of the universe that explains that it evolves toward apparently less probable states in the past. But unsurprisingly, this argument dates back to a time when quantum chance and classical instability had not yet been discovered. In the present theoretical context, however, such an argument simply no longer makes sense, despite the fact that it is often still used to try to justify the kind of evolution that is taking place in the past direction of time. The hypothesis that a reversal of the motion of every particle in an irreversibly evolving system would bring it back to its preceding lower entropy state would actually be true only for a very limited period of time, as short in fact as the system is large and its entropy growth in the future significant¹⁴.

It is certainly right that a true reversal of time would actually have to involve more than a simple reversal of the motion and rotation of all components of a system, as I explained in chapter 3, but even if such a time-reversal operation was applied to the whole universe, there is absolutely no reason to believe that, in the absence of any constraint, the past evolution would be likely to evolve toward lower entropy states, because the only violation of symmetry that might occur as a result of such a time reversal would not

¹⁴The experiments which are sometimes mentioned as having confirmed that a reversal of the motion of all particles in the final state of a macroscopic system are observed to induce anti-thermodynamic evolution are misleading, because the processes involved take place under carefully controlled conditions, where random perturbations are absent over the totality of the short period during which the phenomena occur and therefore they merely confuse us into believing that the mystery of the continuous diminution of entropy that is taking place in the past direction of time is explainable as being the mere consequence of an improbable configuration of the present state, while this is clearly impossible under more general conditions.

be such as to allow anti-thermodynamic evolution. In any case, even if we were to assume that a system could be so carefully prepared that despite the known sensibility to initial conditions which exists even in a classical deterministic context and despite the inherently random nature of quantum processes, the system would nevertheless follow an evolution so unlikely that its entropy would be continuously decreasing all the way back to the first instants of the Big Bang with absolute precision, we would still be left with having to explain why it is that the present state of the universe happens to be of such an unlikely nature that it allows this kind of awkward evolution to take place. Clearly, this attempt at explaining the occurrence of the lower entropy states into which the whole universe evolves in this direction of time we call the past cannot be considered satisfactory.

What is also problematic with the assumption that the entropy-reducing evolution which we observe to take place in the past direction of time could be the mere outcome of a precise adjustment of the present microscopic state of the universe is that, even if we take this as an explanation for the diminishing entropy, we still cannot explain why such an adjustment does not occur for the future instead of the past, because even if that was the case it would simply seem like the past is replaced with the future and the future with the past and we would still not be able to explain why there is, in effect, an asymmetry. What we should actually expect to observe, if it was a precise adjustment of initial conditions that explained the occurrence of time-asymmetric behavior, is a situation where entropy would be continuously decreasing in various regions of the universe whose initial microscopic states would have been carefully prepared so as to produce anti-thermodynamic behavior, but not all in the same direction of time, that is, not all in the past direction for all locations. There is absolutely no reason to expect that such carefully prepared systems would all be set so as to evolve with diminishing entropy in only one particular direction of time, because time itself does not impose such a requirement. But we do not observe multiple, oppositely directed arrows of time in our universe and this is precisely what would have to be explained for such an approach to be made valid. We cannot assume that the reason why entropy-decreasing evolution is not occurring toward the future, from time to time, in certain locations, is that the precise initial conditions required to produce it are too unlikely, while we would also be assuming that the precise ‘final’ conditions required to produce a decrease of entropy in the past are, for their part, allowed to occur, even if they are no less improbable. The rules of probability applied to initial conditions

would lead us to predict that entropy should increase in the past, just as it increases in the future, and therefore they cannot alone explain the existence of an arrow of time, even if they do at least explain why it is that entropy does not decrease in the future.

Now, even if we were to recognize that the situation in which multiple coexisting subsystems would be set to evolve with decreasing entropy in opposite directions of time would probably be highly unstable, as the precise configuration required to produce a decrease of entropy in a given region would be subject to interference by what happens in another region where entropy would be decreasing in the opposite direction of time, there is no reason to believe that such a mixture of oppositely evolving subsystems should, through some kind of interference, give rise to a universe with a single well-defined direction of its arrow of time, as required by observations, that is, by our memory of past events. What must be clear is that, if we do not expect to frequently observe such carefully prepared subsystems evolving with diminishing entropy in the future, then we should not expect to observe the entire universe itself to evolve in such an unnatural way in the past, but this is precisely what is happening all the time, and if that is indeed the case then there must be another explanation to it.

It is only as a consequence of the fact that, for practical reasons, our thought processes are always functioning in the direction of time in which entropy is rising (thereby giving rise to a psychological arrow of time) that we usually fail to recognize that the kind of evolution that takes place in the past direction of time is amazingly abnormal from a purely probabilistic viewpoint. Thus, while it is certainly true that the present state of the universe is relatively unlikely configured, for example in the sense that, if time was reversed a local tendency for matter particles to disperse would momentarily turn into one for particles to convene, while a tendency for wave fronts to spread would turn into one for wave fronts to converge, this is explainable as merely being a consequence of the fact that the original state in the past that gave rise to the present ‘final’ state was, itself, in a highly unlikely configuration, even from a purely macroscopic viewpoint. It’s not the final states which are inexplicably organized, but really the initial state (in the distant past) that gave rise to them.

One of the oldest attempts at solving the problem of the origin of the arrow of time, which must also be considered inadequate, was originally proposed by Ludwig Boltzmann, the originator of the kinetic theory of gases. It was based on the recognition that there always occur fluctuations to lower

entropy states for randomly-evolving, isolated systems which are in a state of thermal equilibrium. Over a very long time-scale, it should sometimes happen that those fluctuations would be so significant as to bring even a system in thermal equilibrium into a state with an entropy so low that any subsequent evolution would likely be characterized by a continuous increase of entropy. Thus, it was proposed that the universe, as the ultimate isolated system, really starts in a maximum entropy state, which would presumably be a likely state to be randomly chosen as our initial conditions, and then remains in such a state during most of its existence, but that once in a while, as it evolves in either the past or the future, it simply fluctuates to a much lower entropy state from which it would naturally be expected to evolve with continuously increasing entropy back to its more likely, maximum entropy state in the same arbitrarily determined direction of time, which we would then call the future, regardless of its actual (relative) orientation. The fact that such an evolution would perhaps appear to be similar to that which we presently observe to occur at the level of the universe as a whole then suggests that this is what explains the continuous growth of entropy in one single direction of time that characterizes the evolution of all systems which have not yet reached back a state of thermal equilibrium.

It should be clear, however, that in such a context, the only reason we would have to expect to observe the universe in a phase of continuously growing entropy, instead of finding it in one of the much, much more common phases of unchanging maximum entropy would be that this entropy growth is necessary for the presence of an observer which can witness such an evolution. Indeed, the fact that we are allowed to experience a memory of past events and to have a persistent conscious existence is dependent on the condition that there exists a well-defined thermodynamic arrow of time. The problem, however, is that if such a requirement was to be satisfied merely as a consequence of the occurrence of a fluctuation in an otherwise unchanging maximum entropy state, then we should not expect to observe entropy to be so low in all parts of the universe and as far back in time as the epoch of the Big Bang. A much more localized and ephemeral fluctuation, that would provide the observer with no records of a long-lasting history, would do just as well for allowing such a condition at the present time and given that such a fluctuation would be more likely to occur than a long-lived fluctuation involving the entire universe, then based on this kind of argument what we should experience is a short-lived fluctuation.

The question, therefore, remains: Why is the universe evolving irre-

versibly in one single direction of time in all locations and throughout its entire lifetime? One cannot hope to satisfy the requirement imposed by the time-symmetric nature of fundamental physical laws by simply postulating that the universe actually evolves without any constraint, either in the past or the future, because that would leave the very property of irreversibility unexplained. As Boltzmann himself appears to have realized, the entropy-fluctuating universe scenario is ineffective for explaining this very constraining aspect of reality and therefore cannot count as a valid solution to the problem of time asymmetry.

Now, the fact that I'm suggesting that the random nature of elementary physical processes and the sensibility to initial conditions is what allows to reject the possibility that it could be a precise adjustment of the present conditions that would completely explain the diminution of entropy that is observed to take place in the past direction of time does not mean that I'm agreeing with the opinion that irreversibility is occurring at a fundamental and irreducible level in our description of physical processes, as was once proposed by some of those who pioneered the study of chaotic systems. I do not believe that we must equate unpredictability and randomness with irreducible time asymmetry, even if, in its most general form, statistical mechanics, as a probabilistic theory, is dealing with systems in non-equilibrium states whose evolution is inherently irreversible. The fact that quantum field theory can be considered to be a more fundamental instance of statistical mechanics, while it is definitely a time-symmetric theory, clearly indicates that my position is justified¹⁵. It would certainly not be appropriate to abdicate the requirement of symmetry under a reversal of the direction of time simply to provide an explanation for the observed unidirectionality of thermodynamic processes, in the context where our most valuable physical theories are all time-symmetric at the most elementary level of description.

The difficulties that we are experiencing in trying to identify the constraint that allows to derive irreversible evolution from time-symmetric physical laws should not be allowed to become a justification for abandoning some of the requirements we have very good reasons to believe must constitute part of a fully satisfactory solution. We would not be wise to reject a theoretical framework that works so well, even if it may seem that it cannot explain

¹⁵I will explain in the latter portion of chapter 5 under which conditions irreversibility can be expected to enter the quantum-mechanical description of elementary particle processes and what consequences this must have on observed aspects of reality.

every aspect of reality, simply to follow an alternative approach which also cannot be made to describe all significant aspects of reality. The challenge consists in actually explaining irreversibility, not in decreeing that it is the foundation of reality, when this would require abandoning most of everything else we have learned. I believe that the fact that we have not yet been able to achieve this objective is not an indication that our most fundamental theories are wrong, but merely a proof that we still do not fully understand all the consequences of the physical principles upon which they were built.

It is important to note, in this regard, that it has also been proposed that it is perhaps a fundamental irreversibility of the quantum measurement process that allows to explain the asymmetry of the evolution in time of observable physical phenomena that does not appear to characterize the evolution that takes place in between measurements. But while I do not want to immediately enter into a discussion of how irreversibility intertwine with quantum theory, I must point out that it would be circular reasoning to assume that it is the measurement process that gives rise to thermodynamic irreversibility, while it is already recognized that it is the irreversibility of the processes taking place in the environment with which a quantum system interacts that is involved in giving rise to the decoherence effect that characterizes all quantum measurements. But even if we were to follow such a route, it is not clear what would explain that this same unidirectionality does not instead operate toward the past rather than the future. After all, there is no sign of an intrinsic asymmetry regarding the direction of time in the equations of quantum theory. Why would quantum evolution always pick the same one particular direction of time instead of another during those processes that can be qualified as measurements? Once again, even if for pure convenience it was assumed to be the case that quantum theory, or a hypothetical process of actualization of potentialities, was to show preference for one direction of time instead of another, we still would be left with as great a mystery to explain, because time itself does not provide the means for such a distinctive feature to arise.

I do agree that irreversibility (just like time itself) is real and constitutes an objective aspect of physical reality and is not just a consequence of some arbitrary choice regarding the level of coarse-graining, but what I will try to demonstrate is that the suggestion that it is no longer appropriate to conceive of reality in terms of elementary particles obeying time-reversible physical laws is not justified, even when we are dealing with complex systems which exhibit strong non-linearity or highly irreversible evolution. As we will

progress, it will become clear that the idea that there should be no laws more fundamental than those which currently apply only under those particular conditions is excessive in proportion to the very specific nature of that most extraordinary property of physical reality we are trying to explain.

4.7 Gravitational entropy

Now that I have properly defined and circumscribed the problem of time asymmetry, I would like to discuss one further problem that arises when one acknowledges the objective nature of the growth of entropy that follows the formation of a black hole or indeed of any overdensity in the matter distribution and which is attributable, in part, to a local increase in the amount of information required to describe the exact microscopic state of the gravitational field. As a consequence of the progress I have achieved in better understanding the origin of the growth of gravitational entropy that takes place under such circumstances, I will be able to provide a decisive solution, not only to the problem of the unaccounted growth of missing information associated with the formation of event horizons, but also to the problem of the violation of the conservation of information which appears to take place in the context where the expansion of space is continuously creating new, elementary, quantum gravitational units of space in the vacuum.

What is already known, concerning gravitational entropy, is that it grows when the mass of an astronomical object and the strength of its gravitational field are rising. Thus, when gravitational attraction is involved, the natural tendency for matter to spontaneously disperse into a larger volume of space is overcome and the decreasing entropy of matter that follows is compensated by the even larger increase of entropy presumably attributable to the growth of missing gravitational field information. In fact, we currently have no exact definition for the entropy attributable to the gravitational field in a general context and it is merely a knowledge of the exact formula for black-hole entropy that allows one to estimate the magnitude of this entropy in the absence of event horizons. In any case, the prevailing character of gravitational entropy means that, when a large enough amount of matter is present in a given volume of space, particles with the same sign of energy are allowed to become more densely packed, because such an evolution is favored from a thermodynamic viewpoint, in the context where there are more possible microscopic configurations of the gravitational field compatible with a state of

higher density. Only the expansion of space could perhaps allow this natural tendency to be surmounted, if the growth of matter entropy attributable to cosmic expansion was to become rapid enough that it would overcompensate the growth of gravitational entropy that occurs as a result of the formation of inhomogeneities.

In section 3.10 I have explained that the presence of event horizons provides us with a unique set of macroscopic physical parameters, which allow a natural definition of coarse-graining and therefore an objective measure of entropy growth. Indeed, what emerges from the semi-classical theory of black-hole thermodynamics is that an exact quantitative measure of missing information is associated with the area of a black-hole event horizon. Given that one is allowed to assume that the choice of this area as a macroscopic coarse-grained parameter is unavoidable (as it is not possible for an observer outside a black hole to obtain a more detailed description of the state of matter and its gravitational field than is provided by a knowledge of this macroscopic parameter), then it becomes possible to conclude that an objective definition of the entropy of a system actually exists under such circumstances that is determined by the value of this unique parameter. Thus, any change to entropy which is reflected in a variation of the mass or the surface area of a black hole constitutes a non-subjective change which cannot be attributed merely to a particular choice of macroscopic parameters, as under such conditions no other macroscopic parameter is available to define an alternative measure of entropy. Something essential, therefore, differentiates the measure of information associated with the surface of a black hole from that which is associated with a general surface for which the Bekenstein bound applies and this distinction has to do with the availability of information. Yet the fact remains that the information loss which appears to be taking place when a black hole absorbs low entropy matter cannot be considered real, because, as I mentioned in section 3.10, it seems that the information about the microscopic state of the matter that was submitted to gravitational collapse is encoded in binary form in the microscopic degrees of freedom associated with the elementary units of area on the event horizon of the object and is released when the black hole decays through the emission of macroscopically thermal radiation.

Once it is recognized that no information needs to vanish from reality, even when a unique definition of coarse-graining exists that gives rise to an absolute measure of entropy, what must be understood is that the existence of a measure of information associated with the area of a black-hole event

horizon does not only mean that information is not lost when such an object absorbs low entropy matter, it also implies that there is a real growth in the amount of information that would be required to completely specify the state of all the microscopic, binary degrees of freedom on the surface of a black hole of increasing mass, which reflects the existence of a growing amount of microscopic structure in the gravitational field. Indeed, when one recognizes the appropriateness of the assumptions that allowed me to derive an exact measure for the entropy of elementary black holes, based on a knowledge of the relevant discrete variables that characterize the fundamental states of matter reaching a spacetime singularity, then it becomes clear that the amount of missing information which would be required to specify the exact microscopic state of all the matter particles which were captured by the gravitational field of a macroscopic black hole is not large enough to account for its entropy growth.

What my findings from section 3.10, regarding the existence of a relationship between black-hole entropy and the degrees of freedom associated with the discrete symmetry operations indicate, in effect, is that the amount of information which would be required to completely specify the microscopic state of matter particles is actually decreasing when matter is captured by the gravitational field of a black hole. But if this measure of information is not growing when the mass of a black hole is rising, while the total amount of missing information (the entropy) is growing faster than the mass of the object (which rises in proportion to its matter content), then one has no choice but to recognize that the amount of information which would be required to describe the microscopic state of the gravitational field is indeed rising when local gravitational fields grow stronger, at least in the presence of a macroscopic event horizon.

It should be clear, therefore, that when matter assembles into a macroscopic black hole, the number of microscopic degrees of freedom associated with the gravitational field grows larger, even while the number of microscopic degrees of freedom associated with matter particles is being reduced as a result of the constraints exerted on their states of motion by the gravitational field of the object. Thus, while information about the exact microscopic state of the matter that fell into a black hole is not provided by the macroscopic physical parameters that describe the object and may therefore appear to be lost, an even larger amount of information, concerning the exact microscopic state of the gravitational field, becomes missing, which also contributes to increase the entropy of the black hole. As a result, the amount

of missing information appears to be growing faster than we would expect if information was conserved, that is to say, it grows faster than the progression of our ignorance concerning the intricate details of the microscopic state of matter.

Now, given that I will later argue that the growth of inhomogeneities in the matter distribution, which is the source of stronger gravitational fields, provides the dominant contribution to entropy increase in our universe (because the entropy of matter itself does not change much as a consequence of expansion), then it would appear that temporal irreversibility actually arises mostly as a consequence of the growth of gravitational entropy. What is crucial to understand, under such conditions, is that the irreversible character of this evolution as well as the growth in the amount of missing information which is giving rise to it cannot be considered subjective features of reality, precisely because they can be associated with the presence of event horizons which constitute natural boundaries, enabling a unique definition of coarse-graining that is entirely determined by the strength of local gravitational fields.

It must, in effect, be recognized that a growing amount of information is required to describe in complete detail the structure that emerges in the gravitational field when overdensities develop in the matter distribution (or in the distribution of dark matter that is attributable to local variations of vacuum energy). I believe that it is merely because we do not benefit from the guidance of a fully developed quantum theory of gravitation that we haven't yet realized that the amount of missing information is actually growing faster than would appear to be allowed, when a gravitational field gains in strength as a consequence of a local increase in the energy density of matter (we often hear about people claiming that information may be lost when matter falls into a black hole, but I have never heard anyone complaining about the growth of missing gravitational information). One has no choice, however, but to recognize that when a gravitational field gains in strength as a result of a local growth in the density of positive- or negative-energy matter (either baryonic or dark), even in the absence of a macroscopic event horizon, a real increase in the amount of missing information, which is not due merely to increased ignorance about the exact microscopic state of matter particles, and a related growth of entropy, which is not dependent on any subjective definition, are taking place.

What one would normally object concerning this characterization of gravitational entropy is that the growth of missing information which can be

expected to occur when stronger gravitational fields develop, would appear to violate the constraint of conservation of information that is imposed by quantum theory. Yet, I have also come to understand that despite the fact that when the mass of a black hole is growing the amount of missing information appears to be rising faster than would be allowed as a mere consequence of growing ignorance concerning the microscopic state of matter, the total amount of information required to describe the exact microscopic state of our universe does not really change when gravitational fields gain in strength and therefore the requirement of conservation of information is not violated when such a transformation takes place. What must be clear, in any case, is that either information is always conserved, or else it never is, and given that the latter conclusion does not appear to be valid from a fundamental viewpoint, then it must be recognized that the additional measure of gravitational-field information which appears to be produced when a black hole forms already existed before it contributed to the measure of gravitational entropy associated with such an object, just like the information contributed by matter itself. In other words, if a larger than expected change in the amount of missing information is, in effect, occurring, when the density of matter grows locally, then it means that any such variation needs to be compensated somehow.

Indeed, what implies that the additional growth in the amount of missing gravitational field information which is associated with stronger local gravitational fields can be objectively characterized is merely the fact that it occurs as a result of adopting the natural definition of coarse-graining that is provided by the measure of spacetime curvature associated with the presence of macroscopic event horizons as natural boundaries with well-defined macroscopic physical parameters. If those considerations are appropriate, however, then it becomes necessary to recognize that the growth in the amount of missing gravitational field information that is taking place as a consequence of a local increase in the density of matter can only be compensated by a change in the amount of information which would itself be independent from any arbitrary choice regarding the coarse-graining and therefore we can already expect that the compensation would arise from additional changes in the strength of local gravitational fields.

The situation we face, therefore, is one in which the amount of information that is missing (and which determines the coarse-grained measure of entropy) is continuously rising, even though it is only in situations where stronger gravitational fields develop, due to a local increase in the density of

matter (of positive or negative energy sign), that this variation can be characterized as a non-subjective change attributable, in part, to a real growth in the amount of information required to completely describe the microscopic state of the gravitational field, rather than to our growing ignorance of the details of this exact microscopic state. Once this is recognized, then what remains to explain is how information can nevertheless be conserved, as would presumably be required in a quantum gravitational context. In fact, what allows me to conclude that the amount of missing information is growing faster than would appear possible, when stronger gravitational fields develop as a result of the formation of a matter overdensity, is not merely the results of my analysis of the nature of the microscopic degrees of freedom of matter constrained by the gravitational field of a black hole, but the very fact that it also appears necessary to assume that there is an opposite variation of the same kind that occurs when local gravitational fields grow stronger as a result of the formation of an *underdensity* in the large-scale matter distribution, which suggests that it is only as a consequence of the fact that there arises a compensation between those two variations that the measure of gravitational field information can be left invariant on a global scale regardless of how much it varies locally.

What I'm suggesting, more exactly, is that, given that a higher than average matter density appears to be associated with an additional amount of missing information which was not present initially, due to the fact that a larger amount of information is required to completely describe the detailed microscopic state of the gravitational field under such conditions, then it should necessarily be the case that a correspondingly smaller amount of information would be required in order to completely describe, with the same level of precision, the microscopic state of the gravitational field associated with a lower than average matter density. You may recall that, in section 2.6, I explained that a void in the cosmic distribution of positive-energy matter must actually be considered to exert a gravitational repulsion on the surrounding positive-energy matter, due to the fact that the presence of such a void implies an *absence* of gravitational attraction, which would otherwise compensate that which is attributable to the surrounding matter distribution, whose center of mass is always located in the exact position of the particle experiencing its gravitational influence. But if those gravitational forces are, in effect, attributable to an absence of gravitational interaction with the positive-energy matter that is missing in the void, then it means that it is entirely appropriate to assume that a lesser amount of information

would be required to describe the microscopic state of the gravitational field as a result of the presence of such a void.

In the context where it would be necessary to assume that the initial matter distribution was very uniform to begin with, this conclusion would imply that any additional increase in the amount of missing information necessary to describe the microscopic state of the gravitational field attributable to the formation of a local matter overdensity would be compensated by an exactly corresponding decrease of information attributable to the presence of the underdensity that must necessarily form in the surroundings of this overdense structure in order to allow it to grow. As a result, I can deduce that, despite the fact that real changes take place locally in the measure of information, when the matter distribution is growing more inhomogeneous, information is always rigorously conserved, even when this evolution involves an alteration of the macroscopic parameters associated with black hole event horizons. But it must be clear that those conclusions only apply in situations where it is gravitation that provides the dominant contribution to entropy change and where it is a variation in the amount of information required to describe the microscopic state of the gravitational field that compensates another variation of the same kind, because the entropy of matter itself does not diminish when its density decreases locally, just as it does not grow, in general, when its density increases and it comes to occupy a smaller volume of space.

Thus, when the density of matter grows larger than its average value locally, there is an increase in the amount of missing information that is not attributable merely to our growing ignorance concerning the exact microscopic state of matter, but which is due in part to an actual increase in the amount of information necessary to describe the exact microscopic state of the gravitational field. However, when the density of matter becomes smaller than its average value locally, there occurs a corresponding *decrease* in the amount of information necessary to describe the microscopic state of the gravitational field. It is the decrease of gravitational-field information, attributable to the formation of this underdensity in the macroscopically uniform matter distribution, that compensates the additional unaccounted increase in the amount of missing gravitational field information which is attributable to the formation of the corresponding overdensity and which would otherwise give rise to a violation of the condition of conservation of information. In other words, when the mass of an astronomical object is growing, more information than would appear to exist initially is required

to describe the exact, final, microscopic state of its gravitational field, whose higher strength is responsible for most of the entropy growth that occurs under such conditions. But in an originally smooth (positive- or negative-energy) matter distribution, the growth of mass in one place can only arise when a corresponding diminution of mass takes place in the surrounding area and from a certain perspective, less information would appear to be required to describe the exact *microscopic* state of the gravitational field when the matter density is reduced below its average cosmic value in such a way, even if a stronger repulsive gravitational field would seem to be produced locally as a result of such a change (from the viewpoint of an observer with the same sign of energy).

What is happening, therefore, is that, given that it is not necessary to describe the microscopic degrees of freedom which are absent in the gravitational field as a result of the absence of gravitational interaction with the matter that is missing inside an underdense region in the matter distribution, it follows that the microscopic state of the gravitational field can be completely specified using a smaller amount of information. In fact, as I will explain below, it is this particular dependence of gravitational-field information on the local density of matter that allows one to understand why it is that when the density of matter grows *larger* in a local region of space, more information than may appear to have existed initially is required to describe the microscopic configuration of the gravitational field, because there is no *a priori* motive for assuming even such an outcome, despite the fact that it appears to be required by the semi-classical theory of black-hole thermodynamics, in the context of my account of the constraints applying on the microscopic state of matter particles reaching a future singularity. What is crucial to understand, however, is that the decrease of gravitational-field information that occurs when the density of matter is being reduced below its average value locally does not translate in an overall reduction of gravitational entropy.

Indeed, if the density of matter is only allowed to decrease in a given region of space when a compensating increase takes place in its vicinity, it follows that the information loss that occurs as a result of the formation of an underdensity in the matter distribution only serves to increase the amount of information necessary to describe the exact state of the gravitational field associated with the corresponding matter overdensity. But this means that information which was available before the change took place, due to the absence of a macroscopic gravitational field or event horizon, would now be

missing and would merely contribute to the objective growth of entropy that is attributable to the formation of the overdensity whose mass grew at the expense of the formation of the underdense structure. Therefore, even if the amount of information necessary to describe the exact microscopic state of the gravitational field is allowed to diminish locally, it can be expected that the objectively defined measure of gravitational entropy will always be growing in the universe, as any measure of entropy does under ordinary conditions when information becomes unavailable, despite the fact that it does not vanish from reality.

This possibility for gravitational entropy to rise globally at the expense of a local decrease in the measure of available information contained in the same force field is reflected in the fact that the strength of local gravitational fields is actually growing, even when this growth is attributable to an absence of gravitational interaction consequent to the formation of a void in the matter distribution. As a consequence, the changes occurring when a void is forming in the matter distribution are still likely to take place when gravitation is predominant, which is certainly appropriate, given that they are actually favored from a thermodynamic viewpoint. It remains, however, that the gravitational fields attributable to the presence of voids in the positive-energy matter distribution do not have the exact same thermodynamic properties as the similar gravitational fields attributable to the presence of negative-energy matter overdensities, as will be emphasized below, and this is a reflection of the fact that, locally, there is always less information in the gravitational field following the formation of a void in the uniform matter distribution, while more information is contained in the gravitational field following the formation of a negative-energy matter overdensity.

Anyhow, once it is recognized that the amount of information required to describe the microscopic state of the gravitational field must be reduced when the density of matter diminishes below its average value locally, then it becomes possible to conclude that the total measure of information concerning the microscopic state of the gravitational field always remains constant globally, as required by quantum theory and despite the objective nature of the growth of entropy that occurs when gravitational fields gain in strength as a result of the development of inhomogeneities in the matter distribution. Thus, even though variations in the amount of missing information or entropy can be characterized in a more objective way when gravitation is involved, there is no fundamental difference between those changes and the ones that take place when there is no significant variation in the strength

of local gravitational fields¹⁶. A real diminution of information does occur when the density of positive-energy matter diminishes below its average cosmic value and the strength of the *repulsive* gravitational fields experienced by positive-energy matter grows locally and this is what allows to compensate the additional growth in the amount of missing information that occurs when a gravitational collapse is taking place in the positive-energy matter distribution.

What we fail to recognize, from a conventional viewpoint, is not only that a *local* increase in the density of positive-energy matter and a stronger attractive gravitational field give rise to an objective increase in the amount of information required to describe the microscopic state of the gravitational field that would appear to be larger than allowed by the conservation of information. We also fail to recognize that a *local* diminution of matter density below its cosmic average would actually give rise to a diminution in the amount of information required to specify the microscopic state of the gravitational field (given that the gravitational field attributable to such an underdensity in the matter distribution would arise from a local absence of gravitational interaction). The fact that such a compensation is required to take place if information is to be conserved can be considered to provide confirmation of the appropriateness of the results I derived in section 2.6 to the effect that, from the viewpoint of positive-energy observers, not only must voids in a uniform distribution of positive-energy matter be the source of repulsive gravitational fields, but those gravitational fields must originate from uncompensated gravitational attraction by the surrounding positive-energy matter distribution.

The significance of the conclusion that a local decrease in the density of matter must give rise to a local diminution in the amount of information required to describe the microscopic state of the gravitational field is much more profound than one may expect. Indeed, despite the fact that there is a certain equivalence between the gravitational field produced by the presence of a *negative-energy* matter overdensity and the gravitational

¹⁶In fact, it appears that it is always the case that when a large enough static force field develops, additional information is required to describe the exact microscopic state of this field of interaction and it may, therefore, always be required that a compensating contribution occurs in the environment. This is perhaps a desirable outcome given that, according to my analysis of black-hole entropy (discussed in section 3.10), the fields associated with other long-range interactions can actually be expected to carry their own specific measures of entropy.

field attributable to an underdensity in the *positive-energy* matter distribution, a clear distinction must nevertheless exist between those two situations with regards to thermodynamic properties. In section 3.11 I explained, in effect, that a negative-energy black hole in a vacuum would have to radiate negative-energy particles and would therefore have a negative temperature. Thus, if a void in the positive-energy matter distribution was deep enough over a sufficiently large region to produce a gravitational field equivalent to that of a negative-energy black hole, it would appear necessary to assume that it also has a negative temperature, given that its surface gravitational field is opposite that of a positive-energy black hole and similar to that of a negative-energy black hole. But temperature merely defines the relationship between thermal energy and entropy as Clausius' original definition of entropy change through the formula $dS = dQ/T$ clearly shows (in this equation dQ is the amount of heat absorbed by a system with a temperature T in the small time interval during which it evolves between two equilibrium states). Under such conditions, if a negative-energy black hole has a negative temperature, it must, in effect, radiate negative-energy particles, or negative heat (which is a positive change for the energy of such an object), so that its surface area and its entropy can diminish in the process, just as they do when a positive-energy black hole radiates positive-energy particles. One might, therefore, be tempted to assume that the thermodynamic properties of a sufficiently large void in the positive-energy matter distribution would be identical to those of a negative-energy black hole and that such a structure would radiate negative-energy particles. But that is not the case.

First of all, it must be clear that there is nothing wrong with the idea that the temperature associated with the thermal radiation of a negative-energy black hole is negative. Once it is understood that this radiation process arises as a consequence of the thermodynamic requirement that local energy differences be smoothed out, even in the presence of event horizons, then it is clear that a negative-energy black hole must lose negative energy if its mass is to decrease in the process. Given that positive-energy matter cannot cross a negative-energy black hole's event horizon and remain inside such an object, it follows that this loss of negative energy can only occur through the emission of negative-energy particles outside the event horizon. A negative-energy black hole would, therefore, *release* negative heat in its environment and in the process necessarily reduce its surface area and its entropy, which therefore requires the temperature of the object to be negative. But it is precisely here that the distinction between a negative-energy black hole and a sufficiently

large void in the positive-energy matter distribution would arise, because the thermodynamic tendency to reach equilibrium would not produce the same outcome in the case of the void in a positive-energy matter distribution, despite the similarity of the gravitational fields associated with both kinds of configuration.

Indeed, while the gravitational field produced by a sufficiently large void in the positive-energy matter distribution must be equivalent to that of a negative-energy black hole from an external viewpoint, in the case of the void, the uniformity of the distribution of energy cannot be re-established through the emission of negative-energy particles by the void, because there is no way that such a radiative process could allow the void to regain the lost positive energy that gave rise to its growth, even if negative-energy particles were present inside the structure and could surmount the growing gravitational attraction exerted on them as they would stray from the center of mass of the void. What would happen, therefore, is that equilibrium would be reached through the otherwise unlikely *absorption* of positive-energy particles from the surrounding matter distribution, which is not forbidden, as it would for a negative-energy black hole, because the strength of the repulsive gravitational field actually decreases as a positive-energy particle approaches the center of the structure, given that the equivalent mass of the object is not all concentrated in a central singularity, as is the case for an ordinary negative-energy black hole. Thus, even if the temperature of a sufficiently large void in the positive-energy matter distribution is negative, as long as the contribution of the negative vacuum-dark-matter energy, which must be present inside the structure, can be ignored, the void would not be expected to reach equilibrium through the emission of negative heat, but rather through the absorption of positive heat, which would actually allow gravitational entropy to be reduced in the process, not because the strength of its repulsive gravitational field would be reduced, but due to the fact that the positive energy that is absorbed would be gained at the expense of a reduction in the density of the surrounding positive-energy matter distribution.

This conclusion is actually a mere reflection of the fact that the temperature of a void in the positive-energy matter distribution, like that of a negative-energy black hole, must be assumed to be negative from the viewpoint of a positive-energy observer. In section 3.11 I have explained, in effect, that it is when an increase of energy produces a local decrease of entropy that the temperature of a system must be considered negative. From the preceding discussion, it should be clear that while a negative-energy

black hole would satisfy this condition as a consequence of the fact that a reduction of the magnitude of its negative energy, through the emission of negative-energy radiation, would give rise to a reduction of its surface area and therefore, also, of its gravitational entropy, a sufficiently large void in the positive-energy matter distribution would satisfy the same condition merely because, when it absorbs positive-energy particles (a process which is as unlikely as the release of positive-energy particles by a positive-energy black hole), a diminution of gravitational entropy must take place, due to the fact that more information becomes available about the exact microscopic state of positive-energy matter and its gravitation field, as the density of the surrounding distribution of ordinary positive-energy matter is being reduced. Thus, even if the surface gravitational field and the temperature of a sufficiently large void in the positive-energy matter distribution could actually be identical to those of a negative-energy black hole, one must conclude that the thermodynamic properties of those two kinds of matter inhomogeneities are not exactly the same and this is merely a reflection of the fact that, while the amount of information required to describe the microscopic state of the gravitational field grows when a negative-energy matter overdensity is forming, it must decrease when a void develops in the positive-energy matter distribution. Yet, it would appear that the requirement that the measure of *missing* information, or entropy, be rising overall, when a void is forming in the matter distribution, is not incompatible with this conclusion, as I have just explained.

But if the assumption that the information required to describe the microscopic state of the gravitational field decreases, locally, upon the formation of an underdensity in a homogeneous matter distribution is to be considered valid, it must be further justified from a more elementary perspective. I will now explain what justifies my conclusion that a local decrease in the density of matter is to be associated with a reduced amount of information in the gravitational field. What must be clear, once again, is that, despite the apparent similarity between the gravitational fields associated with voids in a positive- or negative-energy matter distribution and those attributable to overdensities of opposite energy sign, there nevertheless exists a fundamental difference between those two categories of objects, which arises from the fact that negative-energy matter does not consist of voids in a positive-energy matter distribution, but is rather equivalent to voids in the positive-energy portion of the *vacuum*, as I emphasized in section 2.8. It must be clear, also, that the conclusion that the formation of a void in the matter distribution

would give rise to a negative change in the amount of information required to describe the microscopic state of the gravitational field is also valid in the case of a negative-energy matter distribution. The sign of changes occurring in the measure of information required to describe the exact microscopic state of the gravitational field is not dependent on the sign of energy of the matter whose density varies. Based on the above discussed argument, it is, therefore, also necessary to conclude that a sufficiently large void in the negative-energy matter distribution should have a positive temperature, just like a positive-energy black hole, because its surface gravitational field is equivalent to that of such an object. However, such a void would not reach thermal equilibrium with its surrounding environment through the emission of positive-energy radiation, but really through the unlikely absorption of negative-energy particles, which would give rise to a local reduction of gravitational entropy, arising from a decrease in the density of the surrounding negative-energy matter distribution.

I believe that what explains that the formation of a void in the uniform positive-energy matter distribution would give rise to a negative change in the amount of information concerning the microscopic state of the gravitational field, while the formation of a void of similar magnitude in the positive-energy portion of the vacuum, which can be liken to the formation of a negative-energy matter overdensity, would produce a positive change in the measure of missing information concerning the gravitational field, is the fact that in the absence of vacuum dark matter (as local variations of vacuum energy density attributable to the curvature of space), the distribution of vacuum energy is really uniform on all scales, while the macroscopically homogeneous distribution of matter (baryonic or dark) in which a void may be produced is not really uniform on a microscopic scale. Indeed, in the absence of any matter there are no persistent density variations in the distribution of vacuum energy, such as those which would be associated with the presence of real particles, and removing energy from such a perfectly uniform distribution cannot be assumed to reduce the amount of structure that would initially have been present, on a microscopic scale, in the gravitational field which is attributable to the presence of this energy. This is unlike the situation we have when we are dealing with what would normally be considered a homogeneous matter distribution, in which there actually exist smaller scale variations in the density of energy, which create local gravitational fields which may not be apparent from a macroscopic viewpoint, but which can be as strong as the average density of matter is high.

Thus, when we locally reduce the density of matter particles in a *macroscopically* uniform matter distribution, we reduce the amount of microscopic structure contained in the gravitational fields which are present in this matter distribution as a result of its own small-scale inhomogeneity and this can only mean that, in such a case, we need less information to describe the exact microscopic configuration of the gravitational field, because we actually reduce the amount of structure that previously existed in this field as a result of the inhomogeneity of the microscopic distribution of matter particles. By contrast, when we increase the density of matter with an opposite energy sign in a local region of space (which amounts to add more voids in the positive or negative portion of the homogeneous distribution of vacuum energy), we produce microscopic gravitational fields that were not present beforehand in this region and it is only appropriate that, in such a case, the amount of missing information associated with the microscopic structure of the gravitational field is actually growing locally. This is all a consequence of the fact that the presence of more matter energy allows particles to exert stronger attractive gravitational forces on each other, so that removing positive or negative energy from a local portion of the *vacuum* really has for consequence to generate additional variations in the microscopic state of the gravitational field, given that it is equivalent to increasing, either the number of matter particles with an opposite energy sign, or the magnitude of their energy.

Unlike a local reduction in the density of positive vacuum energy that could be likened to the presence of negative-energy matter, a local reduction in the average density of positive-energy matter gives rise to a diminution in the amount of information necessary to describe the microscopic state of the gravitational field, and this is reflected in the fact that an underdensity in the positive-energy matter distribution does not have the exact same thermodynamic properties as an overdensity in the negative-energy matter distribution, despite the similarity of the gravitational fields produced by the presence of both kinds of astronomical structures, from an external viewpoint. In such a context, it becomes possible to actually explain, not only why it is that the amount of information contained in the gravitational field must diminish when a void forms in the uniform, large-scale matter distribution, but also why it is that the amount of missing information about the microscopic state of the gravitational field is actually growing when the density of matter is increasing locally.

This argument, concerning the distinction between those local diminu-

tions in vacuum energy density which are equivalent to the presence of additional matter energy and local diminutions in the density of matter energy itself, would also justify assuming that, even the gravitational field attributable to an apparently uniform matter distribution would contribute a certain measure of information, despite the fact it is traditionally assumed that only the gravitational fields associated with the presence of macroscopic inhomogeneities in the distribution of matter energy contain information. Indeed, if locally reducing the density of matter produces a decrease of information in the gravitational field, then it would seem that, even on a cosmic scale, a certain amount of information should be contained in the gravitational field produced by this uniformly distributed matter, which would be reduced as a result of expansion. This reduction would occur because the decrease in the average density of matter particles and the consequent diminution of their average kinetic energy that must take place as temperature goes down would reduce the magnitude of the microscopic inhomogeneities present in the distribution of matter energy (independently from any local density variation that would take place as a result of gravitational instability), which means that the strength of the gravitational fields which are present on a microscopic scale would also be reduced and would, therefore, contain a smaller amount of information. The situation here is similar to that which arises in the presence of macroscopic gravitational fields, only, in the present case we are dealing with additional degrees of freedom which are normally left out of a classical description of the gravitational field attributable to a uniform matter distribution. In fact, the same condition of conservation of information which imposes a compensation between the local variations of the different measures of gravitational-field information attributable to the formation of macroscopic matter inhomogeneities, appears to require that a certain measure of information be associated with the *microscopic* structure of the gravitational fields which are present in a macroscopically homogeneous matter distribution.

Indeed, if information must always be conserved, then, as expansion takes place and the average density of matter decreases, a reduction in gravitational-field information would take place which would need to be compensated by an increase of gravitational information of equal magnitude. Now, it has already been proposed that the expansion of space should perhaps be considered to produce an increase of information on the quantum-gravitational scale, given that it would appear to continuously produce additional elementary units of space in the vacuum, in apparent violation of the

theoretical requirement regarding the conservation of information. I believe that this suggestion is valid, because, according to the developments introduced in section 3.10, it appears that a larger number of elementary units of space would imply the existence of a larger number of fluctuating, elementary black holes in the vacuum and a complete determination of the state of the microscopic (quantum gravitational) degrees of freedom of the gravitational field on the surface of those objects would require additional binary units of information. But unlike those who previously discussed the issue, I do not believe that this information growth (which is actually a growth of gravitational-field information) constitutes a serious difficulty, because I know that this change is compensated by the diminution of gravitational-field information that takes place when the magnitude of the average densities of positive- and negative-energy matter is reduced as a consequence of the expansion of space¹⁷. For this to be a valid proposal, however, it must be recognized that a variation of the uniform portion of vacuum energy density, or the cosmological constant, does not contribute to alter the total amount of information contained in the microscopic state of the gravitational field, despite the fact that, like the varying average density of matter energy itself, the varying average density of vacuum energy provides a variable contribution to the gravitational fields that influence the rates of expansion determined by positive- and negative-energy observers.

The conclusion that a variation in the average value of vacuum energy density, or the cosmological ‘constant’ cannot result in a variation of gravitational-field information actually constitutes an essential requirement if information is to be conserved on a cosmological scale, because, if such a variation was taking place, gravitational-field information would vary as a result of changes occurring in the average density of both matter energy and vacuum energy and this would be problematic, because the variation of gravitational-field information which would occur over the entire lifetime of the universe, as a result of those changes, could not be compensated by the variation of information attributable to the growing volume of space. This would be a consequence of the fact that only the average density of matter necessarily

¹⁷In fact the expansion of space would also produce an increase in the number of microscopic degrees of freedom of the gravitational field associated with spacetime torsion, as well as in the number of microscopic degrees of freedom of the two components of the field of interaction associated with the unified charge, but those changes would need to be compensated independently (for similar reasons) by a decrease in the density of spin and a decrease in the density of non-gravitational charges, respectively.

diminishes as the scale factor grows and could be reduced to a minimum value if the volume of space was to become arbitrarily large, or be raised to a larger value if the volume of space was to be reduced to a smaller value (in either the past or the future direction of time). The magnitude of the average density of vacuum energy, on the other hand, could in principle become larger (at least temporarily) in a universe of growing size and diminishing matter density, or else become smaller in a collapsing universe with growing matter density and a diminishing volume of space, with no possible compensation of the changes taking place in the total amount of gravitational field information, thereby precluding information from being conserved as it must be, that is, independently for positive- and negative-energy observers who may experience different rates of expansion and different average matter densities.

However, once it is recognized that changes in the gravitational field attributable to a variation of the *average* density of vacuum energy or the cosmological constant do not contribute any changes to the amount of information necessary to describe the microscopic state of the gravitational field (given that all persistent inhomogeneities in the distribution of vacuum energy are actually equivalent to the presence of matter) and therefore need not be taken into account in balancing the rising amount of information associated with the growing volume of space produced by expansion, then those difficulties no longer exist. From the viewpoint of a positive-energy observer, the average density of matter (including the density of dark matter attributable to local variations of vacuum energy density) would, in effect, be continuously decreasing as a consequence of the expansion of space, along with the associated measure of information required to specify the microscopic state of the gravitational field attributable to the presence of this matter, which was originally maximum (this is allowed, given that a reduction of negative-energy matter density contributes like a reduction of positive-energy matter density to lower the measure of gravitational-field information, despite the opposite sign of the variation of energy density involved). But at the same time, the amount of information associated with the number of elementary units of space present within a co-moving volume (or, more accurately, the number of elementary units of area on the two-dimensional boundary of the same volume) would grow to some arbitrarily large value, thereby compensating the change to gravitational-field information that is associated with the diminishing matter density, while a similar compensation would occur in

the case of a collapsing space¹⁸.

I believe that this is the strongest argument which can be formulated to the effect that it is appropriate to consider that even the diminution of average matter density which is taking place on a global scale actually gives rise to a decrease in the amount of information contained in the microscopic state of the gravitational field, because the validity of this solution depends on the assumption that a variation of matter density occurring in a *macroscopically* uniform matter distribution must, in effect, be assumed to produce changes in the microscopic state of the gravitational field which are different from those we would expect to occur as the consequence of a variation in the uniform portion of vacuum energy density (which would also appear to confirm the validity of the hypothesis that matter energy is equivalent to missing vacuum energy). If I have properly conveyed the nature of the insights which have allowed me to arrive at such a conclusion, then it should be clear that there is no longer a problem with the fact that the expansion of space appears to produce information. From my viewpoint, even if this growth in the amount of information concerning the microscopic state of the gravitational field associated with empty space must indeed be considered real, it would not give rise to a net increase in the total amount of information required to completely specify the microscopic state of the gravitational field associated with both empty space and the matter distribution, for the universe as a whole.

Before concluding this section, I would like to add a few important remarks concerning the role played by voids in balancing the gravitational-field information budget. First of all, it must be emphasized, again, that the conclusion that there is less structure in the gravitational field following the formation of an underdensity in the positive-energy matter distribution (which gives rise to a repulsive gravitational field for positive-energy observers), would also be valid from the viewpoint of a negative-energy observer, despite the fact that such a configuration gives rise to a stronger *attractive* gravitational field for negative-energy matter. In fact, for negative-energy observers, just as is the

¹⁸If those conclusions are appropriate, it would mean that the idea proposed by certain authors that the size of the elementary units of space determined by the natural scale of quantum gravitational phenomena is perhaps growing with time, so that the amount of missing information associated with the total volume of space would be constant despite expansion, which should eventually give rise to a ‘Big Snap’ that would rip everything apart, can be considered unnecessary and this is certainly appropriate, given that no such an event seems to be occurring.

case for positive-energy observers, a local reduction in the density of positive-energy matter gives rise to a diminution of gravitational-field information, because the presence of positive-energy matter is equivalent to an absence of negative vacuum energy and when less *missing* negative vacuum energy is present, less inhomogeneity exists in the distribution of vacuum energy, just like when less negative-energy matter is present and therefore less information is necessary to describe the microscopic state of the gravitational field in both situations. In the limit where no matter would be present at all and only vacuum energy would remain, the energy distribution would become totally devoid of local structure that could be imparted to the gravitational field and it is only at this point that the amount of information contained in this field would be null. It is the presence of matter of any energy sign that produces more structure in the gravitational field on a microscopic scale, whether this gravitational field is that which is experienced by a positive-energy observer or that which is experienced by a negative-energy observer in the same situation and the fact that a local absence of positive-energy matter gives rise to a gravitational field similar to that produced by the presence of more negative-energy matter does not contradict this conclusion.

If gravitational-field information was not reduced, but rather increased, from the viewpoint of positive-energy observers, in the presence of a gravitationally attractive underdensity in the negative-energy matter distribution, a violation of the conservation of information would follow, unless one was willing to assume that the local growth of negative-energy matter density that takes place in the neighborhood of this underdensity as a result of its formation gives rise to a reduction of gravitational information into negative territory, from the viewpoint of such positive-energy observers. But this would imply that the gravitational entropy associated with such a gravitationally repulsive configuration would itself be negative. Indeed, less information also means less missing information and if the measure of entropy or missing information associated with the presence of a negative-energy matter overdensity was indeed negative, then, given that its magnitude must rise when a negative-energy black hole absorbs negative thermal energy or heat, one would have to conclude that the temperature of such an object cannot be negative, even though it must radiate negative energy particles or negative thermal energy. It is therefore possible to confirm that, despite the fact that the presence of an underdensity in the negative-energy matter distribution must give rise to an attractive gravitational field for positive-energy matter, the formation of such a structure does not produce an increase of

gravitational-field information, from the viewpoint of a positive-energy observer, because, if this hypothesis was valid, it would give rise to a contradiction concerning the sign of temperature of a negative-energy system.

It is actually the requirement that temperature be negative in the case of a negative-energy black hole that allows the same relationship between entropy variation and heat absorption to remain valid, even when negative-energy systems are involved and the heat can only be absorbed as negative. But given that temperature is merely an intensive measure of thermal energy, if we have no choice but to recognize that thermal energy itself is negative for a negative energy black hole or any other negative-energy system, then it only makes sense to conclude that the temperature of such a system is negative and that the relationship between the heat it absorbs and the variation of its entropy is preserved. Therefore, it is only the magnitude of temperature that becomes more evenly distributed when two opposite-energy systems (not black holes) exchange heat, in agreement with the previously noted observation (from section 2.13) that, from the viewpoint of negative-energy observers, kinetic energy, and therefore also thermal energy, is absorbed and released as a negative-definite quantity by negative-energy systems, even when they interact with positive energy systems (while it is absorbed and released as a positive-definite quantity by positive-energy systems).

Thus, even though, from the viewpoint of a positive-energy observer, more negative-energy matter means less positive vacuum energy and therefore less gravitational interaction of positive vacuum energy with the rest of the positive energy present in the universe, this does not mean that less information is required to describe the exact microscopic state of the gravitational field in the presence of more negative-energy matter, because what determines the amount of microscopic structure in the gravitational field is the inhomogeneity of the distribution of its sources in space, as well as the average density of this energy distribution itself. The truth is that, in the presence of negative-energy matter, positive vacuum energy is only missing locally and this means that more structure exists in the gravitational field when the density of negative-energy matter is growing. So, again, despite the fact that negative-energy matter overdensities produce a repulsive gravitational field, their presence does not give rise to a diminution of gravitational-field information, as would be the case in the presence of voids in the positive-energy matter distribution, and therefore entropy itself remains positive in all situations.

But this conclusion is valid only in the context where we recognize that

a macroscopically uniform distribution of positive- or negative-energy matter contains a non-vanishing value of gravitational-field information that is merely reduced toward the null value it would have in the absence of any matter, when a void forms in this homogeneous matter distribution. Thus, it would not be appropriate to define the measure of gravitational information as being zero in the absence of any inhomogeneities in a matter distribution, because, even in the absence of any macroscopic inhomogeneities, there is information in the gravitational field associated with such a matter distribution, due the fact that matter is not a perfectly uniform fluid. Indeed, even when matter particles are as homogeneously distributed as they can be in empty space, there remains enough energy inhomogeneities to produce local gravitational fields which would contain a certain measure of information and this information grows along with the average energy or temperature of the matter particles.

It should be clear, therefore, that underdensities in a macroscopically uniform matter distribution are not associated with a negative measure of gravitational-field information (whatever that would mean), but merely with a lower than average amount of information, because the gravitational-field information attributable to a macroscopically uniform matter distribution (which exists when the matter particles are distributed as uniformly as possible in the available space) is not null and is allowed to decrease as a result of the expansion of space on the cosmic scale. Indeed, while gravitational-field information is associated with the presence of overdensities in a matter distribution, a larger *average* density of matter energy implies that the gravitational fields attributable to the presence of microscopic matter inhomogeneities are stronger and therefore more variable in space, so that more information is required to describe their own microscopic states.

It is important to understand that the contribution by the average density of matter to gravitational-field information does not arise merely from the uniform portion of the matter distribution, but from the presence of inhomogeneities in this macroscopically uniform matter distribution and in the absence of microscopic inhomogeneities, this contribution would be null, even though it varies only as a function of the average matter density. If this was not the case, a similar variation of vacuum energy density would produce a similar change in gravitational-field information and this information would no longer be conserved in the course of expansion. But given that, even when positive-energy matter particles are as homogeneously distributed as they can be, local gravitational fields are still present, because

the particles only occupy a small portion of the space which is available to them, unlike is the case for the uniform portion of vacuum energy, which really is homogeneously distributed on a microscopic scale, then it is possible for gravitational information to exist, even in a macroscopically homogeneous matter distribution. It is this measure of information that decreases when the average density of matter decreases and the strength of microscopic gravitational fields is reduced, as a result of the decreasing magnitude of microscopic inhomogeneities which is attributable to the expansion of space and the diminution of the average kinetic energy of matter and radiation particles.

One must, therefore, distinguish the hypothesis that matter and radiation particles are homogeneously distributed (which can be satisfied when no macroscopic inhomogeneities are present in the matter distribution) and the hypothesis that the energy of matter itself is homogeneously distributed (which may not be satisfied even when no inhomogeneities are present in the matter distribution). Thus, even in what one would consider to be a perfectly homogeneous matter distribution, where no local variations would exist in the density of matter particles, above or below the average density, there would remain enough energy inhomogeneities to produce local gravitational fields which would contain a certain measure of information and this is also true for a positive-energy observer regarding the distribution of negative-energy matter, despite the fact that the homogeneous portion of the negative-energy matter distribution exerts no gravitational force on positive-energy matter. Indeed, it is the inhomogeneities which are present in a distribution of matter energy that contribute to the measure of gravitational-field information that can be reduced by expansion and inhomogeneities in the distribution of negative matter energy do exert a gravitational force on positive-energy matter (even though the force itself cancels out on the largest scale).

In case this is not yet completely clear, I must also mention that from a practical viewpoint, even for a local void in the positive-energy matter distribution, the amount of gravitational-field information would not really be null, because vacuum dark matter with negative energy would necessarily be present in the void (as a result of the presence of the repulsive gravitational force that is experienced by positive-energy observers) and this matter would contribute a measure of information whose sign is opposite that of the change which is attributable to the formation of the void. It is only the void itself which gives rise to a reduction of gravitational-field information, despite the fact that it gives rise to a stronger local gravitational field similar to that

which is produced by the vacuum dark matter it contains. It would be incorrect to assume that the information change associated with the formation of a void in the positive-energy matter distribution cannot be opposite that which is associated with the local growth of negative vacuum-dark-matter energy density that takes place as a consequence of the formation of this void. Anyhow, when the contribution of vacuum dark matter is taken into account, it follows that the gravitational entropy associated with the presence of a void in the positive-energy matter distribution would always remain larger than zero and would be continuously rising as more negative-energy matter concentrates in the structure.

4.8 The initial singularity

What emerges from the preceding reflection concerning the character of gravitational entropy is that, while the amount of information required to describe the microscopic state of the gravitational field is growing in those places where matter is becoming more densely packed, an equal amount of information is being lost at the same time in the gravitational field as a consequence of the resulting diminution of density which is taking place in the surrounding matter distribution. Yet, given that, ultimately (when a black hole forms), the information that is gained becomes missing information which is no longer accessible from an experimental viewpoint, while the information that is lost was in principle available (as the original matter distribution was not constrained by the presence of a macroscopic gravitational field or event horizon), then gravitational entropy must be assumed to rise whenever the matter distribution is becoming more inhomogeneous. What I will now explain is how significant this conclusion actually is, in the context where the initial distribution of matter energy at the Big Bang appears to have been one of inexplicably high uniformity. Thus, I will argue that, for what regards irreversibility, it is the measure of gravitational entropy that constitutes the significant difference between the state that emerged from the past Big Bang singularity and the state into which our universe will evolve in the far future (independently from whether it continues to expand or collapses back on itself). This discussion will set the stage for the more significant developments which will be introduced in the next section and which will provide the actual explanation for the existence of the thermodynamic arrow of time, as a cosmological phenomenon.

It is important to note, first of all, that there is no paradox associated with the fact that the universe still evolves irreversibly, while the initial state at the Big Bang was already one of near perfect thermal equilibrium, because, as Roger Penrose first pointed out [35], under such conditions it is only the portion of entropy which excludes the contribution of local gravitational fields that is maximum. In fact, what transpires from the developments introduced in the previous section, concerning gravitational entropy, is that it is precisely the smoothness of the initial distribution of matter energy (which is reflected in the uniformity of the temperature of the cosmic microwave background) that is responsible for having allowed the universe to evolve irreversibly at later times, because under such conditions it is the growth of matter inhomogeneities which must have provided the dominant contribution to irreversible entropy growth in our universe, since the epoch of matter-antimatter annihilation. What really needs to be explained, therefore, is not why the universe evolves irreversibly, despite the initial state of thermal equilibrium, but why the energy of matter was actually so homogeneously distributed initially that gravitational entropy was almost perfectly null, even if that would appear to be a highly unlikely configuration to begin with, in the context where a much larger number of possibilities exist for the microscopic state of matter and its gravitational field, which would not be characterized by such a uniform matter distribution and an absence of primordial black holes.

In order to explain those facts, one needs to identify the nature of the constraint imposed by the fundamental, time-symmetric physical laws on the boundary conditions at the Big Bang that is responsible for the very high level of homogeneity and the very low gravitational entropy that characterizes this initial state. We must, therefore, once again transcend our natural reluctance to apply the known principles of physics to the Big Bang if we are to avoid having to modify the laws themselves in order to achieve greater overall consistency. It would be incorrect to assume that proposing a solution to the problem of the origin of time asymmetry that relies on the application of certain constraints to the initial conditions at the Big Bang would be akin to requiring divine intervention. The most fundamental principles must be assumed to be valid under absolutely all conditions, including those that existed during the Big Bang. I believe that it is our failure to acknowledge the importance of this requirement that explains most of the difficulties we currently face in theoretical cosmology.

But before we can achieve some real progress in understanding why ir-

reversibility occurs, we must first recognize that the source of most changes to entropy that take place after the early annihilation of baryons with antibaryons is actually to be found in the growing strength of local gravitational fields. It is as a consequence of gravitational attraction that the stars, in particular, can form and are allowed to release their radiation and it is due to gravitational collapse that black holes, as the objects with the highest entropy density, can form and grow more massive at the expense of a local reduction of matter density in their environment, which is also the source of stronger gravitational fields. But one need not assume that this is due to the 'fact' that gravitation is always attractive, as all that is required is that it be attractive among particles with the same sign of energy, which allows gravitational energy and therefore also gravitational entropy to be proportional to the square of the mass of an object, instead of being merely proportional to its mass, as does matter entropy itself. In such a context, it may appear that a much larger number of possible initial states would be characterized by the presence of an abundance of black holes and other density fluctuations, while those initial conditions would not have had as much potential for allowing subsequent evolution to take place irreversibly. However, given that the presence of primordial black holes would have disturbed the process of structure formation in the initial matter distribution in ways which would have had observable consequences at the present epoch, then it seems necessary to assume that the initial Big Bang state was virtually free of black holes and therefore it remains to explain why the universe was in such an unlikely configuration at the Big Bang.

One thing that should be clear is that the weakness of the gravitational interaction in comparison with the other fundamental forces and the fact that it became predominant over those other interactions only during the matter-dominated era, does not mean that no constraint that would be imposed on the magnitude of local gravitational fields (or the curvature of space attributable to local matter inhomogeneities), could be involved in determining the early conditions which are responsible for the existence of the thermodynamic arrow of time. Indeed, if gravitational instability allowed structures to begin developing only at a relatively late time in the matter distribution, this is due precisely to the fact that the initial distribution of matter energy was so uniform to begin with and this is a condition that applies on the magnitude of local gravitational fields. But it would certainly be inappropriate to assume that a constraint which would be imposed on the initial magnitude of local gravitational fields would not have much impact as a consequence of

the very fact that the magnitude of local gravitational fields was, in effect, so small initially.

In section 4.5 I have explained why we can actually expect the universe to be expanding. But the fact that we are not instead observing it to be contracting at the present moment can only be explained as being the consequence of another fact, which is that the magnitude of local gravitational fields is decreasing continuously in this direction of time relative to which the universe is contracting. If we perceive the universe to be expanding, it is simply because, as thermodynamic processes, our memories are formed only in the direction of time in which the inhomogeneity of the matter distribution is growing, while, if the strength of local gravitational fields and the measure of gravitational entropy were growing in the direction of time relative to which the universe is contracting, then we would necessarily perceive the universe to be contracting. This is actually all that can be meant when we say that we experience the universe to be expanding, because in fact we do also ‘observe’ space to be contracting, but merely in the sense that we also have knowledge of the contraction of space that occurs in the past direction of time, as we may witness by watching a backward-running movie of the same events. Thus, what explains that the universe is observed to be expanding (what explains that the cosmological arrow of time is oriented in the same direction as the thermodynamic arrow of time) is the fact that gravitational entropy is practically null in the primordial Big Bang state, while it is allowed to grow to arbitrarily large values at later times and this means that if we want to explain why it is that we observe an expanding universe, then we must first explain why it is that its initial state was characterized by such a low gravitational entropy.

But, in the context where one must acknowledge the presence of negative-energy matter in our universe, the fact that the density of matter was much larger in the past does not make the initial smoothness of the matter distribution more unexpected, as one may be tempted to assume. Indeed, even in a hypothetical static universe with an arbitrarily large volume of space, an initial configuration characterized by a greater uniformity of the distribution of matter energy would not necessarily be more likely as a randomly chosen boundary condition for the universe, because even a diluted matter distribution could still contain inhomogeneities on the largest scale, as a result of the fact that negative-energy matter can be concentrated in regions of space distinct from those occupied by positive-energy matter, even if the average density of both types of matter is negligibly small. Thus, in the absence of

additional constraints, the most likely configuration for the initial distribution of matter energy would always be one of higher inhomogeneity, because there usually exist more microscopic configurations of matter and its gravitational field for which the portions of positive- and negative-energy matter that survive the annihilation of matter with antimatter are not mixed up in a perfectly uniform manner.

In any case, the hot Big Bang did occur and the distribution of matter and radiation energy was, in effect, homogeneous to an inexplicably high degree initially. If this hadn't been the case, then macroscopic event horizons would abound in the primordial state and even if the magnitude of density fluctuations was not large enough to prevent the universe from expanding in most locations, what we would observe (either around us at the present moment or in the cosmic microwave background at the epoch of last scattering) would be a much different world. The problem, therefore, is that, it seems that the observable universe should have begun its evolution in a state where the energy of matter would already be highly inhomogeneously distributed and strong local gravitational fields would be present, with which would be associated an arbitrarily large measure of entropy. But if the initial state was not of such a nature, then it means that something must have constrained the universe to have a much lower gravitational entropy initially, because this does not appear to be a natural configuration to begin with when all possibilities are allowed.

It should be clear, however, that the simple fact that the universe must be expanding locally, if an observer is to be present to witness an absence of inhomogeneities, does not provide strong enough a constraint to explain that the initial distribution of matter energy was as smooth as it is observed to be, even if the presence of event horizons would indeed prevent space from expanding locally. The energy of matter could be much more inhomogeneously distributed than it currently is and expansion would still be allowed to proceed unaffected in most locations, even if a large number of primordial black holes had been present initially. It is merely the fact that the inhomogeneity is not as pronounced as it *could* have been that is unexplained.

What must be understood is that the homogeneity of the initial distribution of matter energy is not merely apparent in the low magnitude of local variations in the density of positive matter energy (which can be compensated by *local* variations in gravitational energy, as I explained in section 4.5), but must also be apparent in the near absence of large-scale disparities in the spatial distribution of positive- and negative-energy matter particles.

This particularity is especially significant in the context where one of the only differences which would exist during a recollapsing phase of the universe's history (if such an evolution was actually allowed to happen), would have to do with the fact that, in the recollapsing phase, the dissociation of the positive- and negative-energy matter distributions would actually be much more pronounced, as a result of the gradual polarization of the matter distribution along energy sign which can be expected to occur in the context where particles with the same sign of energy are submitted to mutual gravitational attraction, while concentrations of matter with opposite energy signs gravitationally repel one another.

For the present discussion to be meaningful, however, one must also understand that there are strong motives for believing that, even in the presence of negative-energy matter, it is still appropriate to consider that there arises a state in the past which, from a classical viewpoint, would be characterized as consisting in a spacetime singularity. Indeed, what should be clear, based on the developments introduced in section 2.6, is that a globally homogeneous distribution of negative-energy matter (regardless of what its current average density might be) would exert no influence on the rate of expansion of positive-energy matter and would not diminish the strength of the gravitational field attributable to the presence of this matter, despite the fact that negative-energy matter would in general exert a repulsive gravitational force on positive-energy matter. This conclusion follows from the description of negative-energy matter as being equivalent to the presence of voids in the positive-energy portion of the vacuum and the acknowledgment that the void of cosmic proportion that must be associated with a homogeneous distribution of negative-energy matter cannot give rise to uncompensated gravitational attraction from a surrounding distribution of positive vacuum energy, which would otherwise be the source of the gravitational repulsion that would arise from the presence of such a void.

Once this is recognized, it becomes possible to predict that if the initial matter distribution is sufficiently homogeneous on the largest scale, then nothing can prevent the formation of the trapped surface which according to classical theorems would give rise to a past singularity, even if one of the axioms of the theorems is that matter must always have positive energy. It is only the inappropriateness of the traditional description of negative-energy matter as being the source of absolutely repulsive gravitational fields that makes it seem like the presence of such matter could prevent the formation of a past singularity (or the occurrence of a state of maximum positive and

negative matter energy densities, as one would rather need to assume in a quantum gravitational context). It would, therefore, appear that the very uniformity of the matter distribution which is responsible for giving rise to the existence of a thermodynamic arrow of time is actually required in order that the existence of a past singularity, or an initial state of maximum matter density, be considered unavoidable. This is a decisive observation whose significance will be made more explicit in the following section. I must emphasize, however, that what I have in mind when I'm referring to an initial singularity is not a state where the laws of physics would actually break down, but simply a state where the average, positive and negative densities of matter energy have reached the maximum theoretical values determined by the natural vacuum-stress-energy tensors which enter the generalized gravitational field equations introduced in section 2.15.

In any case, if we are to assume that there must, in effect, be a singularity, or a state of maximum positive and negative matter energy densities at the beginning of (unidirectional) time, then it seems necessary to assume that this singularity is also different, in certain respects, from an ordinary black-hole singularity. First of all, even if the initial state that emerged from the past singularity at the Big Bang had been highly inhomogeneous, it would not be expected to have given rise to the same evolution as that through which a future Big Crunch singularity would go from a backward-in-time viewpoint, because, whereas the state that would emerge from a Big Crunch singularity, in a universe like ours, would evolve back to a more homogeneous state, the state emerging from an initial Big Bang singularity with the same level of matter-energy inhomogeneity would not evolve toward a more homogeneous state, because in our universe future evolution is unconstrained. Therefore, a highly inhomogeneous distribution of matter energy emerging from a past singularity could only evolve toward an even more inhomogeneous state (if there does not exist any limit to gravitational entropy growth), otherwise it would not evolve at all, from a thermodynamic viewpoint, as it would already be in one of its most likely maximum gravitational entropy states. As a result, no reversed gravitational collapse or white hole phenomenon would occur that would release objects of lower entropy, as we would expect to happen in the course of a time-reversed, generic Big Crunch.

But while the Big Bang is not the time-reverse of a Big Crunch, or black-hole gravitational collapse, it also appears that the initial singularity is different from a future singularity owing to the fact that it does not give rise to an initial state characterized by large fluctuations in matter energy density

with which would be associated a very large gravitational entropy, such as would be the case for the final state of a generic future singularity. This observation makes it even more apparent that what is occurring in the past direction of time in our universe is not what one would expect to happen as a mere consequence of the contraction of space. The initial singularity was of such a nature that it could not constitute the outcome of a gravitational collapse of the kind that would occur in a universe in which local gravitational fields are growing in strength. The universe is changing as it collapses in the past direction of time, but not in the way one would expect in the absence of a constraint that operates a continuous decrease in the magnitude of the inhomogeneities present in the distribution of matter energy.

What's significant, as well, is that the presence of past singularities appears to be restricted to the one known initial singularity from which the Big Bang emerged, even if there does exist solutions of the gravitational field equations that would appear to describe processes which would be the time-reverses of a black-hole gravitational collapse. All the evidence indicates that the hypothetical white hole processes which could be described using those solutions never occur in our universe. I believe that if those solutions do not represent processes that can be observed in the forward direction of time in our universe, it is because they would allow gravitational entropy to decrease in this direction of time, while such an evolution is thermodynamically unlikely in the absence of a specific constraint. Indeed, white holes would expel low entropy matter at an arbitrarily high rate, which would reduce their masses and the area of their event horizons faster than would be allowed as a consequence of the emission of macroscopically thermal radiation (this has nothing to do with negative-energy black holes expelling positive-energy matter), so that the processes would involve a decrease of gravitational entropy in the future. It should be clear, therefore, that the Big Bang does not constitute a generic white hole, even though it originates from a past singularity.

The only motive one might have to assume that generic white holes could exist would be that, in all likeliness, the gravitational entropy of a black hole should rise in the past, just like it does in the future, so that, from the forward-in-time viewpoint, the evolution taking place during the same period of time would actually appear as a fluctuation involving a decrease of gravitational entropy that would persist until the present moment is reached. But the problem is that, even if our present state would seem to require the occurrence of such a phenomenon, there appears to be something that con-

strain evolution in the past direction of time to take place with continuously decreasing gravitational entropy, despite the apparent improbability of this evolution, and this is precisely what remains unexplained. If white holes are never observed, therefore, it is simply because such processes would require a decrease of gravitational entropy in the future (which is unlikely), or equivalently, a continuous increase of gravitational entropy in the past direction of time (which for some reason appears to be forbidden). Therefore, if we can understand why the state that emerged from the initial Big Bang singularity had minimum gravitational entropy, then we may also be allowed to understand why there is only one such past singularity.

At this point it should be clear that, even though black holes are the objects associated with the highest possible density of gravitational entropy, it would not make sense to simply assume that the most likely initial state for the universe would be one for which all matter would be contained in one giant black hole, because even a closed universe with a highly homogeneous distribution of matter energy could be considered to satisfy this condition. What is required for gravitational entropy to be maximum is that matter energy be as inhomogeneously distributed as possible, even while the universe is in the process of collapsing into a higher density state from a backward-in-time viewpoint. The relevance of this remark is made more obvious when we are considering a universe that contains both positive- and negative-energy matter. Indeed, in such a context, the state with the highest gravitational entropy would necessarily be one for which the initial distributions of positive- and negative-energy matter would be completely polarized, in such a way that all the matter would be contained in opposite-energy black holes with arbitrarily large masses, whose magnitude would be limited solely by the amount of matter in the universe and the time available for the inhomogeneities to form (if they are not already present to begin with). What must be understood, therefore, is that there is no *a priori* motive for assuming that a high level of polarization of the positive- and negative-energy matter distributions could not also apply to the initial Big Bang state (regardless of the fact that the matter density is then maximum) if such a configuration is, in effect, favored from a thermodynamic viewpoint, because a universe that would evolve without constraint, as space is contracting in the past direction of time, would have more chances to reach such a configuration, not merely despite gravitational repulsion, but as a result of it.

Now, it was once suggested that the smoothness of the initial distribution of matter energy might only be apparent and that a state of higher inhom-

geneity might have existed initially, that was later made uniform through various smoothing processes. But given that such processes would have released a large amount of heat that would have modified the temperature of the cosmic microwave background to an extent that appears to be incompatible with measurements, then it appears that, even if the smoothing could occur at the appropriate time and on the appropriate scale, its outcome would not agree with observational constraints. Furthermore, if the distribution of matter energy had been highly inhomogeneous before any such process could smooth it out, the magnitude of those inhomogeneities would have rapidly been amplified under the effect of the gravitational interaction and it would have become even more difficult to give rise to the homogeneous distribution that is revealed by measurements of the temperature of cosmic microwave background radiation. Indeed, the same argument implies that the initial state cannot have been *perfectly* uniform, otherwise the universe could not have evolved into its present state early enough to allow for the existence of stars, galaxies and other large-scale structures, which means that the constraint responsible for the high level of homogeneity of the initial state must not be so restrictive that it would imply a complete absence of energy fluctuations.

Of course, in the presence of negative energy matter, an additional contribution to gravitational instability exists that more readily triggers the formation of inhomogeneities. But given that most of the baryonic negative-energy matter vanished following the early annihilation of matter with antimatter, then smaller-scale fluctuations in the negative-energy matter distribution were not allowed to grow as much as they otherwise would by the time the cosmic microwave background was released, while larger-scale fluctuations in the distribution of negative vacuum-dark-matter energy only began to grow when they were encompassed by the particle horizon. Therefore there weren't much additional perturbations to the temperature of CMB radiation as a result of the presence of negative-energy matter. But it remains that energy fluctuations, even if only those attributable to local variations of the sign of energy of matter particles, would need to be present initially if inhomogeneities are to develop at a later time, on a larger scale.

What constitutes the most significant difficulty for the smoothing hypothesis, however, is the fact that the existence of cosmic horizons would have forbidden any such process from ironing out inhomogeneities above the scale determined by the size of the horizon, at the time when the CMB was released, and therefore we should not observe uniformity on the largest scale

if the homogeneity of the distribution of matter energy is attributable to smoothing processes obeying the requirement of local causality. An intrinsic limit is actually imposed on such processes, that would prevent them from producing the kind of homogeneous state which emerged from the Big Bang and therefore it appears appropriate to conclude that, regardless of any other difficulty, conventional smoothing processes should probably not be considered a viable explanation for the homogeneity of the initial distribution of matter energy.

As a consequence of the clear inadequacy of conventional smoothing processes and in the absence of a better alternative, it is still widely believed that inflation may be the cause of the very high homogeneity of the universe's distribution of matter energy which is reflected in the small amplitude of cosmic microwave background temperature fluctuations. However, I think that the occurrence of this hypothetical process of accelerated expansion would not be of much help in explaining the observed time asymmetry that characterizes cosmic evolution, because there is no reason to expect that a contracting universe would evolve toward a more homogeneous configuration during the epoch that would precede a hypothetical phase of exponentially accelerated *contraction*, which would then take the universe back to a more likely state of maximum inhomogeneity. If inflation could perhaps explain why the universe evolves in an otherwise unnatural way (from the viewpoint of the growth of gravitational entropy), between the moment when matter emerges from the initial singularity and the instant at which inflation ceases, it could not explain why it evolves toward greater homogeneity from far in the future and back toward the time at which the universe would presumably begin to contract at an exponentially accelerated rate into the initial singularity, now with naturally growing inhomogeneity.

Even if inflation may give rise to a homogeneous universe forward in time, a Big Crunch would not be expected to occur with decreasing inhomogeneity forward in time, unless the state immediately preceding the exponentially accelerated contraction into the final singularity would be required to be as smooth as the state which was produced in the past following ordinary inflation. But assuming that this would occur would amount to require that causality operates backward in time from the final singularity, instead of forward in time from the initial singularity, because, from the viewpoint where causality operates from the past toward the future, a Big Crunch would be more likely to occur with increasing inhomogeneity in the future, right up to the moment when inflation would perhaps take place in reverse and merely

increase the inhomogeneity that would already exist even further and produce an inhomogeneous final singularity. Assuming that this is not what occurs would amount to postulate without motive that classical (unidirectional) causality must rather operate backward, from the instant at which matter emerges from the future Big Crunch singularity and until the moment when the universe would begin recollapsing, after having reached its maximum volume, so that the period of inflation that would occur backward in time from the instant at which matter emerges from the final singularity would give rise to a homogeneous state *after* inflation, in the past direction of time. But there is no *a priori* reason not to assume, instead, that it is a highly inhomogeneous final state existing *before* the phase of exponentially accelerated contraction that gives rise to the inhomogeneous state that would occur in the future direction of time following this phase of exponentially accelerated contraction, as we may expect based on the hypothesis that causality still operates forward in time.

The problem is that the hypothesis that classical causality operates forward in time from the past singularity is necessary for the conclusion that inflation would necessarily produce a homogeneous state, because if it was assumed that it is the events in the future that can influence what occurs backward in time until the moment when matter would start contracting at an exponentially accelerated rate back into the initial singularity, then the state we would expect to obtain following inflation, from the forward-in-time viewpoint, would still be a state of maximum inhomogeneity, while this does not correspond to reality. What must be understood is that, even if we simply interchange future and past, we are still facing a mystery, because, if it is the future that influences the past and if inflation operates backward in time, so as to smooth out the state emerging from a future Big Crunch singularity, instead of giving rise to a homogeneous state forward in time, beginning from the inhomogeneous state that emerged from the past singularity, then we simply reverse the direction in which irreversible evolution would take place and we still have no explanation for why causality, in effect, operates in this particular direction of time (which we would then call the future) and not in the opposite one, while this is precisely what we are trying to explain. Indeed, classical causality, or the rule that past events always have an influence on future events and not the opposite, is simply a consequence of thermodynamic time asymmetry or irreversibility, and if this property is assumed to characterize our universe without question, then it cannot be used to explain irreversibility itself. Therefore, assuming that inflation necessar-

ily produces a highly homogeneous state, from a more likely inhomogeneous state, amounts to assume without justification the very outcome we want to derive, which means that inflation is not valid as an explanation of the origin of time asymmetry that would arise from the necessity of a homogeneous initial state (following inflation).

The fact that inflation itself requires quite unique initial conditions to occur does not even need to be taken into account in order to conclude that it is insufficient to explain the observed asymmetry of the evolution of gravitational entropy. To actually explain the unlikely homogeneous state that emerged from inflation during the Big Bang, using the hypothesis of inflation itself, we would have to predict that this process operates in both the future and the past directions of time (in the same portion of history) to produce a homogeneous state out of the generic inhomogeneous initial states that would emerge from both the initial and the (hypothetical) final singularity and this would require that the direction of time relative to which inhomogeneities are growing mysteriously reverses when the universe starts contracting, when its volume would be maximum (or at any, arbitrarily-chosen, intermediary time, indeed) and as I previously explained there is absolutely no reason to expect that a reversal of the thermodynamic arrow of time associated with the growth of gravitational entropy would occur when space would begin contracting on a global scale.

At this point it is necessary to mention that a variation of the more conventional attempt at explaining cosmological time asymmetry by making use of inflation theory which was proposed more recently [36] postulates that it is through the process of creation out of ‘nothing’ that symmetry with respect to the direction of time can be reintroduced in our description of cosmic history. What is proposed is that the initial state of our universe is actually an extended, low-density, vacuum state, which we may perhaps consider to be a likely state from a thermodynamic viewpoint, given that under such conditions and when only positive-energy matter is allowed to exist, the entropy of matter itself would seem to be maximum. Of course this is not the state in which our universe began, according to observations, but it might be possible to assume that what happened is that the universe emerged out of a local fluctuation in this extended vacuum and that it is inflation that is responsible for having allowed the high-density state so produced to start expanding at a critical rate and if this is indeed the case then the universe could perhaps be considered to necessarily begin in a state (preceding inflation) that is not so unlikely from the viewpoint of gravitational entropy, even if this would

otherwise be unexpected. The idea is that this kind of process could take place in both the past and the future, starting from this initial, extended vacuum state, and in such a way produce a globally time-symmetric history for the universe.

The problem I have with this description, however, is not merely that the proposed solution is dependent on the validity of the hypothesis that there occurred an early phase of inflationary expansion, which can only produce the desired outcome under highly unlikely initial conditions of a distinct nature. The more unavoidable difficulty has to do with the fact that, as a tentative explanation of time asymmetry, it would suffer from the same reliance shared by more traditional approaches on the implicit assumption that there is already a favored direction of time. Indeed, despite what is usually assumed, an extended vacuum state can only arise out of a prior phase of expansion that would occur in the future direction of time. If a large volume of space is to remain nearly empty for a sufficiently long time that a localized fluctuation of vacuum energy is perhaps allowed to give rise to the creation of an entire universe, then this space must have been expanding prior to the creation event and this expansion can only take place in one direction of time at once. In the following section I will explain, in effect, that it is not possible to simply assume the existence of an expanding low-density universe without assuming that it has emerged out of a state of maximum matter density at some point in the past and if this is what happened, then the most likely possibility is still that the initial state was an inhomogeneous state of high gravitational entropy. But it will also become clear later on that, in the presence of negative-energy matter, it is not possible for inflation, alone, to produce a homogeneous state out of a heterogeneous distribution of positive- and negative-energy matter inhomogeneities. In fact, even if one assumed that our universe emerged out of a fluctuation in the extended vacuum state that would follow this initial expansion phase, there would still be no reason to assume that the state that would be produced by this second phase of inflationary expansion would itself have a low entropy, unless we assume that classical (unidirectional) causality operates forward in time (as I explained above), which, again, amounts to simply assume the validity of the result we are seeking to derive.

If we believe that the initial conditions at the Big Bang must be subjected to the same constraint of likeliness as applies to the configurations of matter which are reached through random evolution under more general circumstances, then the fact that it does not appear that this initial state

could have been produced by chance alone, means that there must be an explanation for this anomaly, but this explanation cannot be found in the traditionally favored cosmological models based on inflation theory. It must be noted, again, that the anthropic principle would be of no use in trying to achieve such a goal, because, if the initial conditions are freely determined, they would not be required to be so highly constrained as they appear to have been when matter emerged from its initial state of maximum positive and negative energy densities. Indeed, a state as homogeneous as that which appears to have existed in the remote past is so unlikely to have arisen randomly that, even the chance occurrence of an observer in a universe with a less thermodynamically favorable initial state would be a more likely phenomenon, in comparison. If the universe was initially characterized by such a low gravitational entropy, it is because it necessarily had to go through such a constrained state at least once in its lifetime. What I will now explain is why this conclusion should have been expected all along.

4.9 The horizon problem and irreversibility

So, here we are, having actually ruled out the possibility that the high degree of homogeneity of the matter distribution in the primordial universe could be due to any conventional or inflationary smoothing processes, but with apparently no option left to explain this remarkable fact. Although this outcome may be quite perplexing, the attentive reader may already have perceived a glimmer of light on the cosmic horizon. Indeed, when one carefully looks at all those failed attempts, I believe that one cannot avoid getting the feeling that it is the very fact that there exists a state of maximum positive and negative matter energy densities in the remote past that must constitute the basis of a consistent explanation of the origin of the anti-thermodynamic evolution that is taking place in the past direction of time in our universe, because this is the only aspect of our universe which is correlated with the state of minimum gravitational entropy. What I will now explain is that there is actually a requirement for the magnitude of the densities of positive and negative matter energy to be maximum at a certain point in the history of the universe that does not just follow from the fact that space must expand or contract and this actually allows to explain why it is that the universe did not come into existence in an extended vacuum state with negligible positive and negative matter energy densities. But, quite remarkably, this same re-

quirement is also responsible for having produced a maximum-density state so exceptionally configured that it guarantees that all future evolution will take place irreversibly.

Before embarking on an explanation of how it can be that an arrow of time was allowed to emerge as a result of the existence of a past singularity, however, I would like to first recall my earlier discussion from section 4.5 concerning the requirement for the universe to have a null energy and the fact that it is possible in principle for matter energy to be compensated by the energy of the gravitational field through a variation of the kinetic energy of expansion. There, I mentioned that the zero-energy condition alone does not require that the energy of matter be uniformly distributed, because, even if this constraint may require gravitational energy itself to be null globally, it would not prevent the energy of matter and that of the gravitational field from varying in opposite ways *locally*, as long as there is an overall compensation between all positive and negative contributions to the energy of the gravitational field. But if the energy of the matter that was present in the very first instants of the Big Bang had been as inhomogeneously distributed as it can be, macroscopic black-hole event horizons (any event horizon produced by a black hole more massive than an elementary black hole) would abound in the early universe. So, why is it that observations rather seem to show that the distribution of matter energy was highly uniform, with very few black holes, in the very first instants of the Big Bang?

One can only begin to understand the cause of the homogeneity of the matter distribution that emerged out of the past singularity when one acknowledges that what is significant with our current description of the physics of the early universe is the explicit assumption that the cosmic horizon (sometimes called the particle horizon) begins to grow at the exact moment when the magnitude of the densities of matter energy is maximum. But why should causality have anything to do with the magnitude of the average densities of positive and negative matter energy? I must admit that I always had difficulty accepting the very validity of the notion that the universe could have come into existence as a set of disconnected entities, not causally related to one another, due to the presence of non-overlapping cosmic horizons in the primordial state. The conclusion that the limited velocity of causal signals would forbid interactions between sufficiently distant regions of the universe, however, appears unavoidable. Yet how could such an assortment of disconnected parts as is usually assumed to exist at the Big Bang be considered to form a single universe if its elements are not even related to one another in

any way? How could they even have been allowed to come into contact with each other in a well-defined manner later on, if they weren't part of the same causally interrelated ensemble initially?

This situation would be particularly puzzling in the context where we would consider that there was no quantum bounce and that the Big Bang really constitutes the beginning of time, as there would then be no prior state at which causal relationships could have been established between the initially disconnected regions. But it seems that it is precisely the hypothesis that time begins with the past singularity that motivates the assumption that the cosmic horizon must begin to grow at the precise moment when the density of matter and radiation is maximum. Here, again, I just couldn't understand the appropriateness of a picture that most people accepted as valid without a second thought. But this led me to develop a better understanding of the conditions imposed by the principle of local causality on the initial Big Bang state that turned out to be crucial for explaining the high degree of homogeneity of the primordial matter distribution that is responsible for the existence of the thermodynamic arrow of time.

First of all, I think that it is important to mention that the notion that the size of the cosmic horizon increases with time, as the universe itself expands, contains an implicit assumption that is not always recognized for what it is. Indeed, when one considers that the horizon encompasses an increasingly larger portion of space in the future, one is actually presuming the validity of the classical principle of causality, that is, of the idea that causes always precede their effects. But it is actually always *past* causes that produce *future* effects. It is never assumed that a future cause could produce an effect in the past. This is usually appropriate, as we experience time in a unidirectional way as a consequence of the fact that the thermodynamic arrow of time always operates from past to future and never in the opposite direction. But when we are considering that no signal was allowed to propagate farther than the distance reached by the cosmic horizon at any given time after the Big Bang, we are implicitly assuming that it is only the past that can influence the future and that effects propagate in the future direction of time, from causes which originate in the initial singularity. In other words, we are assuming the existence of a preferred direction in time (the future) and a preferred instant (that of the past singularity) at which causes begin to propagate. But it must be clear that this is an assumption and that there is no *a priori* reason not to assume that classical causality instead operates toward the past from the instant at which a hypothetical future Big Crunch singularity would be

formed, in which case the size of the horizon would already encompass all of space, or at least a very large portion of it, at the Big Bang.

The truth is that causality could begin to operate at any given instant of time, even the present time, and the cosmic horizon would then spread from there in either the future or the past direction of time. Therefore, if what we are seeking to explain is the existence of a preferred direction in time, then we cannot simply assume the validity of the classical concept of a horizon expanding from the Big Bang in the future direction of time. We cannot claim that there is a problem with the homogeneity of the large-scale matter distribution, if this problem arises as a consequence of assumptions concerning the size of the horizon which are only meaningful in the context where there is a preferred direction to causal signals which originates from this very same homogeneity. What we must provide is a consistent justification for the very validity of this particular choice of a horizon concept. We must explain why this particular state in the past was configured in such a way that it allowed classical (unidirectional) causality to be a meaningful concept that came into effect at the exact moment when matter emerged from the past singularity.

But even apart from those considerations, the cosmic horizon concept, as it is currently understood, is somewhat problematic in that, quite ironically, it does not provide any specific requirement concerning the conditions which would make it possible for causal relationships to exist among the various elements of the universe, in the context where it would need to be assumed that the most elementary particles were once separated by distances larger than the horizon. Despite those difficulties, I came to recognize the validity of the limitations imposed by the existence of cosmic horizons. I believe that what allows this concept to be acceptably formulated is simply the fact that, ultimately, as we consider increasingly earlier times, the size of the causal horizon would actually reach the limit imposed by quantum theory on the classical definiteness of any measure of spatial distance. When the size of the cosmic horizon passes below the limit at which the uncertainty that is intrinsic to quantum phenomena would apply to spacetime relationships themselves, it is certainly no longer appropriate to assume that the limited velocity of signal propagation forbids the existence of causal relationships between regions of space separated by distances larger than the size of the horizon, but smaller than this characteristic scale of quantum gravitational phenomena, as there are no classically well-defined relationships of distance and duration below that scale.

In such a context, it would be suitable to assume that there may, after all, exist causal relationships between all physical elements of the universe which were in contact with one another to within an elementary quantum gravitational unit of distance at the Planck time, if we also have good reasons to expect that the size of the cosmic horizon was then equal to this elementary unit of distance, within which the gravitational field and the metric properties of space were submitted to quantum indefiniteness. In fact, from a quantum gravitational viewpoint, it may be preferable to simply recognize that there is nothing smaller than the elementary units of distance or surface associated with this particular scale. But given that causality is a feature of the classical spacetime structure, this means that there would be no sense in imposing limitations on signal propagation below that scale. Therefore, when the size of the cosmic horizon reaches the natural limit imposed by quantum gravitation, as it contracts in the past, if the most elementary particles (whose characteristic size is also that of quantum gravitational phenomena) are allowed to be in contact with one another to within such a unit of distance, then no smaller components would remain causally unrelated in the initial Big Bang state, which is probably sufficient a condition to impose, regarding the necessity for the universe to form a global entity whose elements (the elementary particles) are all causally interrelated as a result of having been in direct contact with one another at least once in the history of the universe, before becoming separated by large spatial distances.

Now, this simple formulation of the requirement which I believe allows the universe to exist as the ensemble of all those things which are physically related to one another and to nothing else may appear benign, even if adequate, but in fact it can be attributed the most amazing consequences, in the context where it is recognized that a density of negative-energy matter as large as that of positive-energy matter must have existed in the initial Big Bang state (for reasons I have explained in section 4.5). Thus, I would suggest that all the elementary particles originally present in our universe at the Big Bang be required to have been in contact with at least one other particle to within a quantum gravitational unit of distance at the Planck time. More specifically, I propose that the following condition must apply.

Global entanglement constraint: There must exist an event at one particular moment of cosmic time, when all the elementary particles which are then present in the universe, regardless of their energy sign, are in contact with at least one neighboring

elementary particle of either positive or negative energy sign to within a Planck unit of distance, in a state of maximum matter density.

If this condition is fulfilled, then any particle that is present in the universe would have once been in contact with a particle that was in contact with another particle and so on, which means that at no time could a physical element of the universe exist that would be causally unrelated to the other elements which are considered to be part of the same ensemble, even if the particles which were initially present in the state of maximum positive and negative matter energy densities later become separated by space-like intervals and are no longer in contact with one another. If this requirement was not fulfilled, there would be no reason to expect that when the cosmic horizon grows in the future, particles which were causally unrelated initially could begin to influence one another through long-range interactions, because those particles would not even be elements of the same universe.

Of course, the existence of such a smallest, physically significant causal horizon does not mean that the limits usually imposed by the size of cosmic horizons on the propagation of causal signals no longer apply, but merely that they need not apply at times earlier than the Planck time. It must be clear that there would be no sense in speaking about the ultimate horizon as being that which would be associated with the epoch at which the whole visible universe would be contained within a single Planck surface, because, once the average densities of positive and negative matter energies is maximum and there is a matter particle with a Planck energy in every elementary unit of area, then no further contraction is possible, as all tentative quantum theories of gravitation appear to confirm. What this means is that it wouldn't even make sense to impose a condition of causal contact on a state that would be reached at an even earlier time. But even if the constraint of global entanglement concerns the state of the universe at the Planck time, it would be incorrect to assume that only the detailed knowledge of a fully developed quantum theory of gravitation would allow us to say anything meaningful regarding the state of the universe at such an early time. In any case, we still need to explain why it is that the matter distribution was almost perfectly smooth, on a scale larger than the size of the horizon at the time when the magnitude of the densities of positive and negative matter energies was maximum, as required if the growth of this cosmic horizon, as a unidirectional phenomenon, is to actually begin at that particular instant of time. This is a

particularly difficult question given that large-scale homogeneity is precisely what would appear to be forbidden by the existence of such a horizon.

The implications of the global entanglement constraint only emerge in the context where it is recognized that event horizons (such as those associated with black holes) can, under certain conditions, constitute potential barriers which are impossible to overcome. It must be clear, first of all, that even though certain positive-energy particles could be prevented from coming into contact with other positive-energy particles in the initial state of maximum matter density, as a consequence of being contained within the macroscopic event horizon of a positive-energy black hole, if only positive-energy matter existed this would not allow to justify imposing a limit on the amplitude of primordial density fluctuations. Indeed, in such a case, regardless of the presence of macroscopic event horizons, all matter particles would eventually end up being in contact with their neighbors, because the contraction of space that takes place backward in time, toward the initial Big Bang state, would lead to the merger of all the event horizons which were originally present and their spacetime singularities, as in a generic Big Crunch process. Under such conditions, all the particles which may now be isolated by the presence of event horizons would nevertheless merge into one initial state of maximum positive matter energy density, where every matter particle with a Planck energy would occupy an elementary unit of space and be in contact with the surrounding particles present in this initial singularity. Thus, if only positive-energy matter was present in our universe, it would seem that the global entanglement constraint could be satisfied in the initial state without gravitational entropy being minimal, because even if strong local gravitational fields and macroscopic event horizons existed in the instants immediately preceding the formation of the singularity (in the past direction of time), all elementary particles would nevertheless be allowed to come into contact with their neighbors in the maximum-density state, because those are attractive gravitational fields.

When negative-energy matter is present, however, things become more complicated. Indeed, if the constraint of global entanglement imposes contact between neighboring elementary particles at the Planck time, regardless of their energy sign, then given that gravitational repulsion, unlike gravitational attraction, may forbid local contacts, by giving rise to insurmountable potential barriers for particles located within black-hole singularities of opposite energy signs, it follows that event horizons can be expected to be absent initially on all but the smallest scale, even if macroscopic black holes

are allowed to form at later times. If this was not the case, then certain particles could exist in our universe that would not be causally related to the rest of it, which I believe would involve a contradiction. In the absence of a condition of global entanglement, the most likely initial state, from a purely statistical viewpoint, would be one for which all the matter in the universe would be concentrated in the smallest possible number of opposite-energy black holes with arbitrarily large masses which would already be in a state of maximum gravitational entropy. But this was not allowed to constitute our boundary conditions at the Big Bang simply because, under such conditions, the singularities at the center of the objects could never come into contact with one another in the maximum-density state, while this is required by the global entanglement constraint.

In the presence of negative-energy matter, global entanglement actually constitutes a very constraining requirement, because any sufficiently large fluctuation in the initial density of positive or negative matter energy would give rise to the presence of an event horizon that would forbid the condition from applying. Therefore, such large fluctuations in the density of positive- and negative-energy matter must be completely absent in the first instants of the Big Bang and can only develop gradually at later times, in an initially smooth and homogeneous matter distribution. The mass of any black hole that is now present in the universe must, therefore, diminish continuously in the past direction of time, as we approach the initial singularity, so as to allow the condition of homogeneity imposed on the initial matter distribution to be satisfied, despite the fact that it is actually the past condition that gives rise to the future configuration, in the context where the condition that applies on the initial Big Bang state is, in effect, one of minimum gravitational entropy, from which the classical (unidirectional) principle of causality itself can be expected to emerge.

What was so puzzling about the previously unexplained fact that an ever smaller number of microscopic configurations seems to be available for matter evolving in the past direction of time under the influence of the gravitational interaction was that no such a decrease in the number of allowed microscopic configurations is observed in the future direction of time. As a consequence of this limitation, predictions of a statistical nature, such as those made using quantum theory, are always valid only for evolution toward the future, while evolution toward the past cannot, in general, be accurately predicted (the probability of prior events cannot be determined from that of subsequent events, while the probability of future events can usually be determined from

that of past events), which is annoying, given that the equations of the theory are symmetric under a reversal of the direction of time. But this is not a consequence of the fact that information concerning the state in which a system will evolve is only available for the past and not the future, because it is possible to recognize, retrospectively, the absence of statistically significant constraints that would apply to future evolution by considering the future of an initial state at a time when this future is now itself in the past. This is in contrast with the evolution that can be observed to take place at the same time toward the past and which reveals that systems can only come to occupy a subset of their theoretically allowed microscopic configurations whose only distinctive property is its lower entropy.

What remained unexplained, therefore, is the fact that an ensemble of systems started in the same macroscopic state evolve to occupy all available microscopic states in the future, while a similar ensemble, started in the same macroscopic state, usually evolves only to past states characterized by a lower entropy and more particularly, a lower gravitational entropy. But I have now explained that this diminution in the number of available microscopic states toward the past originates from the necessity that all the elementary particles present in the state of maximum positive and negative matter energy densities at the Big Bang come into contact with their neighbors of any energy sign, in order that there exist causal relationships between all independently evolving components of the universe. The unnatural evolution that takes place in the past direction of time is the direct consequence of the limitation imposed on the initial state by the condition of global entanglement in the presence of negative-energy matter and it would not merely characterize a small portion of all possible universes, but really all universes governed by the known, fundamental principles of physics, in which a minimum amount of negative-energy matter is actually present. Remarkably enough, this unrecognized, but necessary condition allows to explain why it is that only the gravitational component of entropy was not maximum at the Big Bang, while the entropy of matter and radiation was allowed to be arbitrarily large, which is certainly appropriate given that the universe was then already in a state of thermal equilibrium. The constraint of global entanglement only limits the magnitude of entropy attributable to the gravitational field and this is exactly what we need.

It must be clear that the fact that a perfectly uniform distribution of negative-energy matter exerts no gravitational influence on positive-energy matter does not allow one to assume that it is not necessary to take into ac-

count the presence of negative-energy matter in trying to identify the origin of the constraints that give rise to the homogeneous, initial distribution of positive matter energy, because it is precisely the magnitude of local inhomogeneities in the distribution of positive and negative matter energy which needs to be constrained and negative-energy matter inhomogeneities do have an effect on positive-energy matter. In fact, negative-energy matter always exerts an influence on positive-energy matter under conditions of maximum average matter density, because, locally, elementary black holes are necessarily present (as I have mentioned while discussing the problem of black-hole entropy in section 3.10) and the energy distribution is never perfectly smooth and homogeneous, especially in the context where it is understood that two particles with maximum opposite energies cannot be located in the same elementary unit of space, due to the insurmountable gravitational repulsion they would exert on one another.

The constraint of global entanglement, therefore, merely imposes that the positive and negative matter energy present in the initial, maximum-density state be as homogeneously distributed as necessary for an absence of *macroscopic* event horizons associated with black-hole masses larger than the Planck mass, because it is only under such conditions that the most elementary particles of matter (with the highest possible positive and negative energies), submitted to the gravitational fields of the most elementary black holes (with the smallest possible surface areas), can be in direct contact with one another, regardless of their energy sign, and thus be part of the same universe. This conclusion remains valid even in the context where it must be assumed that there are no direct interactions between positive- and negative-energy particles, because we are not really dealing here with an interaction propagated across space and time, but with the existence of a minimum distance below which there need not even be an exchange of energy (mediated by an interaction boson) for a causal relationship to exist.

While the event horizon of a macroscopic, negative-energy black hole may prevent local contact between the particles that reached its singularity and neighboring positive-energy particles, which cannot cross this event horizon, from a quantum gravitational viewpoint, an elementary black hole, by virtue of its minimum size, would merely constitute the surface of the one and only elementary particle whose motion it constrains, which means that this particle would be allowed to come into contact with particles which are under the influence of the gravitational fields of other elementary black holes in the state of maximum positive and negative matter energy densities, regardless

of their energy signs. This is what justifies assuming that the condition of global entanglement only imposes an absence of *macroscopic* event horizon. If all the matter in the universe was initially concentrated in two macroscopic black holes of opposite energy signs, the particles contained in the singularity of one of the object would remain isolated from those contained in the singularity of the other black hole, even if the event horizons of the two objects were as close as they can be from one another, and this is what must be considered forbidden, for an initial Big Bang state, by the constraint of global entanglement. The very meaningfulness of this condition is in fact dependent on the hypothesis that there exists a minimum, physically significant spatial scale, below which no causal signal needs to have propagated and which corresponds precisely to the size of an elementary black hole, which is associated with the state of at most one elementary particle of maximum positive or negative energy.

What's interesting is that, contrarily to the situation we would have if inflation was assumed to be responsible for the smoothness of the initial matter distribution, it is now possible to explain why it is that the constraint that gives rise to a homogeneous initial state is necessarily effective in only one direction of time. Thus, gravitational entropy can be expected to decrease continuously in the past direction of time from its current intermediary value, even if this would appear to be a very unlikely evolution for the universe to go through from a statistical viewpoint, because if the present inhomogeneity is not reduced, then the smooth merger of the positive- and negative-energy matter distributions that is required for the global entanglement of all particles to take place would not happen. This reduction of gravitational entropy can now be understood to occur regardless of whether space is expanding or contracting, as long as we are, in effect, approaching the instant at which is formed the unique singularity on which the condition of global entanglement is to be imposed.

It is, therefore, simply the fact that the condition that applies to the initial singularity is precisely one of minimum gravitational entropy, from which can emerge a phenomenon of classical (unidirectional) causality that operates toward the future from that particular instant of time, that requires the evolution that takes place at all later times to be such that it allows an initial state obeying this condition to be reached in the past direction of time, because under such circumstances it is, in effect, past conditions that determine future evolution. If quantum theory only works for predicting future events it is because all possibilities are indeed allowed for evolution toward

the future, while only a limited subset of potentialities can be actualized in the past, as a consequence of the constraint that is continuously being exerted on past evolution by the requirement of global entanglement, which imposes a minimum gravitational entropy on the state which existed in the past, when the density of matter was maximum. It is quite remarkable that this apparent backward teleology can be shown to arise from the existence of an inescapable constraint that applies on one particular state only, but even more surprising is the fact that this can be achieved despite the time-symmetry of the physical laws involved, which gave no hint of having the potential to produce such a manifestly irreversible evolution.

It is important to emphasize that, in the context of this explanation of temporal irreversibility, all physical systems, regardless of how isolated they may have become at the present time, must evolve with continuously decreasing gravitational entropy in the same past direction of time, because they are all submitted to the same unavoidable constraint applying to the same unique state of maximum positive and negative matter energy densities in the past. This constraint, therefore, is stronger than any condition that would be imposed independently on the present state of one or another system in order to favor an evolution to lower entropy states in a given arbitrarily-chosen direction of time. Indeed, in the context where all processes are fundamentally unpredictable, a constraint that would apply merely to the present state of a non-equilibrium system could not, alone, impose on this system that it evolves with decreasing gravitational entropy over a very long period of time in either the past or the future, regardless of how carefully the system is prepared. This is in contrast with the condition imposed by the global entanglement constraint which, by its very nature, necessarily and unavoidably applies to all physical systems which are part of the same universe and of no other and which exerts its influence incessantly in the same unique direction of time (toward the initial singularity) and in such a way gives rise to an asymmetry which is actually shared by all systems, including any branch systems which may no longer be in contact with their environment. In the present context, this temporal parallelism is a simple consequence of the fact that all physical systems in the universe are led by a common condition which applies to the state they occupied when the cosmic horizon began to spread and which originates from the requirement that they actually be part of the same universe.

If the initial or final conditions applying on current states cannot alone explain the temporal parallelism of branch systems, it is because, even if this

would be possible for the evolution that takes place in the future direction of time, the fact that, for all practical purpose, such isolated systems never evolve toward a state of higher gravitational entropy in the past direction of time, like they do in the future, but rather *always* evolve to even lower entropy states in the past, means that it is not the conditions applying on current states which alone determine their past evolution. It is precisely the fact that the requirement of global entanglement must, as a matter of consistency, apply to all particles in the universe that guarantees that all branch systems, without any exception, must obey the same constraint of decreasing gravitational entropy in the same direction of time toward the initial singularity. The parallel thermodynamic behavior of isolated branch systems can be expected to occur as a result of the fact that any system that is part of a given universe, regardless of how isolated it might have become, must have been entangled with the rest of the matter in this universe at the Big Bang in order that causal relationships be established between all components of the universe and this implies that even those portions of the universe which are now isolated must follow the same kind of gravitational entropy decreasing evolution that is necessary for achieving this global entanglement at some point in the past.

The parallelism of the asymmetry of thermodynamic evolution can only be explained if there exists a constraint that requires the diminution of the gravitational entropy of all systems in the past direction of time, regardless of how isolated they have become and independently from what their initial or final states are at the present time and the fact that such parallelism is actually observed under all circumstances clearly shows the validity of the arguments that allowed me to determine the nature of this constraint. If those considerations are appropriate, then it would mean that the assumption that the initial conditions at the Big Bang should always be fixed arbitrarily, which would appear to conflict with the assumption that thermodynamics is fundamental and irreducible, is not really incompatible with the notion that there exists a constraint applying on those initial conditions, that gives rise to irreversible thermodynamic evolution as a derived property.

What is important to understand is that a maximum-density state must necessarily occur at one time or another for the global entanglement of all elementary particles to be satisfied and given that such a state would not likely be characterized by an absence of macroscopic event horizons unless it constitutes the mandatory unique event at which global entanglement is enforced on the universe, then one must conclude that our Big Bang really

is this unique event. In such a context, the presence of an initial singularity would no longer be a mere fortuitous consequence of the fact that space is expanding, but would be an essential requirement for the existence of any universe obeying the known principles of physics. I believe that it is the widespread ignorance of this fact that explains that it took so much time for all the consequences of the presence of a Big Bang to be properly understood and appreciated. To the usual three pieces of evidence in favor of the Big Bang which are the observation that space is expanding, the accuracy of the prediction of light element abundances, and the detection of the cosmic microwave background, I would therefore suggest that one adds the theoretical argument concerning the very necessity of a maximum-density state, which is made conspicuous by the undeniable character of our experience of a thermodynamic arrow of time.

Once we recognize that there actually exists an independent requirement for the presence of an initial state of maximum positive and negative matter energy densities in the history of our universe, then any attempt at explaining the apparent unlikeliness of the homogeneous distribution of matter energy in this initial state, by assuming that it is the consequence of an initial phase of inflationary expansion that originated from a fluctuation that would have taken place in an extended empty space, can no longer be considered satisfactory. Indeed, if it is governed by the principle of local causality (as an unavoidable feature of relativity theory), such an extended space would need to have emerged out of a prior state of maximum positive and negative matter energy densities at which global entanglement would have been allowed to take place, which means that this state would necessarily have minimum gravitational entropy, regardless of whether inflation happened or not. Thus, comments to the effect that it would become impossible to explain the existence of an arrow of time if there only existed one single universe in all of space and one single Big Bang in all of time appear to be misguided, because, in fact, it seems that the truth is to be found in the exact opposite statement. It is as a result of having tried very hard to understand why it is that there should, in effect, be a unique initial state of maximum matter density for the universe, by first acknowledging that this is a perfectly legitimate hypothesis, that I was allowed to achieve progress in identifying the cause of the homogeneity of the initial distribution of matter and energy that gave rise to temporal irreversibility.

If gravitational entropy does indeed rise in only one particular direction of time, it is because only evolution away from the initial singularity, either

in the future or in the past, can be expected to be left unconstrained by the condition of global entanglement, which actually gives rise to a well-defined thermodynamic arrow of time, independent of whether space is expanding or contracting. It is, therefore, possible to understand why it is that classical causality operates from past to future in the portion of history that follows the initial singularity and also why it is that the cosmic horizon only begins to spread outward at the Big Bang. It is the fact that the condition of global entanglement would only be required to apply once, even if the universe was to return to a state of maximum matter density at some point in the future, that explains that the evolution that takes place from the moment at which this condition is enforced is not symmetric in time. Thus, it is incorrect to argue, as certain authors do, that in order not to assume the very outcome we are seeking to derive (the temporal irreversibility), it is required that any condition that applies to some initial state should also apply to a final state. Once it is understood that there need only be one state of high matter density and low gravitational entropy in any given universe, then the kind of evolution which can be expected to take place in the direction of time toward that unique state, either in the past or in the future direction of time, would necessarily be different from that occurring in the opposite direction and this allows to explain time asymmetry without assuming it in the first place. What I'm assuming here is not the asymmetry itself, but merely the uniqueness of the state which allows it to arise. I'm not picking up a unique direction of time, I'm merely identifying the necessary conditions that must apply to the distribution of matter energy at one single moment of time and it just happens that those conditions are so unlikely to ever be satisfied again by chance alone that any later or earlier evolution can be expected to take place irreversibly.

Now, if *bidirectional* time does extend past the 'initial' singularity, following a quantum bounce, we can expect space to be expanding and the density of matter to be decreasing immediately after the event (in the past direction of time), while the inhomogeneity of the matter distribution would still need to be minimum if there is to be any continuity in the evolution of the microscopic state of matter and its gravitational field as we pass the point of maximum positive and negative energy densities. But this means that, even for this portion of history, the thermodynamic arrow of time would (initially at least) have the same direction as the cosmological arrow of time associated with expansion and would actually be opposite that we observe on our side in time of the initial singularity. As a result, the area of black-hole event

horizons and the associated gravitational entropy would be growing toward the past (which any observer then present would consider to be her future), which means that, in the future direction of time, the same objects would evolve as white holes emerging from generic (high entropy) past singularities. Thus, it would be inappropriate to simply propose that it is because a condition of low gravitational entropy applies to *all* past singularities that the energy of the matter that emerged from the Big Bang was so uniformly distributed. Anyhow, it must be clear that, if the thermodynamic arrow of time is indeed reversed as soon as the instant of the initial singularity is reached, then whatever occurred during the portion of history that preceded the Big Bang would remain unknowable to observers in the current portion of history. This would be true for the exact same reason that events located in our future cannot be known in advance, which is attributable to the fact that classical (unidirectional) causality and the formation of mutually consistent records of events only take place in the direction of time relative to which entropy is rising.

Still regarding the possibility for bidirectional time to extend past the initial singularity, I believe that it would be inappropriate to assume that, if this hypothesis is valid, it would become impossible to explain the low gravitational entropy of the Big Bang state by imposing a condition on the initial singularity. It was argued, in effect, that if there was a history prior to the Big Bang, the final singularity which would be produced in the future direction of time (which would constitute our initial singularity) would likely be in a high gravitational entropy state (as any future state reached after a long period of random evolution), which would require the state following it (our initial state) to have a similar configuration. But in fact, it is exactly the opposite which is true and the state preceding the initial singularity must actually be very homogeneous, because the constraint of global entanglement applies to the singularity itself, while it is the evolution away from it, in *any* direction of time, which is unconstrained. Continuity merely imposes that the configuration be similar on both sides of the initial singularity, but it does not allow one to determine what this configuration actually is. It is, in effect, only in the absence of an appropriate constraint to be imposed on the initial singularity that gravitational entropy would have to be maximum in both the immediate past and the immediate future of the initial singularity and indeed at all times. Not recognizing this would, again, amount to favor one particular direction of time (that relative to which entropy would be assumed to grow before the initial singularity) without justification, instead

of explaining why such a preferred direction naturally emerges, as I have done.

What's interesting is that, when time is actually unfolding toward states of higher gravitational entropy past the 'initial' Big Bang singularity, then the universe is allowed to be completely symmetric with respect to past and future, not just because the thermodynamic arrow of time is reversed in that portion of history which is unfolding past the initial singularity, but also because it can be expected that the unequal degrees of violation of time-reversal symmetry which explain the matter-antimatter asymmetry that is observed on our side in time of the initial singularity, will have an opposite effect in the portion of history that precedes the initial singularity (in the sense that there will be more baryonic negative-energy matter than baryonic positive-energy matter following the early annihilation of matter with antimatter in the time-reversed portion of history), because the direction of propagation in time that is favored by the violation of time-reversal symmetry remains the same for that portion of history taking place before the initial singularity, while that relative to which entropy is growing and annihilation is happening is reversed. Thus, it is perhaps not just possible, but actually compulsory, to assume that there is, in effect, a history, not so unlike our own, that unfolds in the past direction of time before the instant of the initial singularity. It must be clear, however, that when I say that the thermodynamic arrow of time reverses when the state of maximum matter density is reached in the past, I do not mean that an observer living in this portion of history would be allowed to witness white hole phenomena and other violations of the rule of entropy increase. In the present context, the thermodynamic arrow of time associated with the variation of gravitational entropy would reverse for all processes, without any exception, at the exact moment when the universe would begin expanding, following the quantum bounce, and therefore no entropy-violating processes could actually be observed.

The picture that develops, therefore, is that of a universe for which gravitational entropy is growing continuously in both the future and the past of a state of maximum positive and negative matter energy densities, characterized by a highly homogeneous spatial distribution of positive- and negative-energy matter particles. This irreversible evolution can be expected to continue regardless of whether space is expanding ever more rapidly or eventually begins to recontract. But in the context where gravitational entropy is continuously growing as a consequence of the polarization of the positive- and negative-energy matter distributions, it follows that, if there is an in-

finite amount of matter in the universe, there may never arise a state of maximum stability, equivalent to thermal equilibrium, where gravitational entropy would become maximum and would no longer rise. Under such conditions, it cannot be expected that the universe will ever evolve back to a state similar in every respects to the state in which it was at the Big Bang, because the probability that such a universal Poincaré return would occur is not merely low, it is decreasing all the time (in section 5.12 I will provide a decisive argument to the effect that there is not enough time for such a return to occur). We may, therefore, be justified in describing the evolution that takes place on a cosmic scale, in both the very far future and the very far past, as truly irreversible. The ancient view of a universe reaching its heat death and remaining in this sterile randomly fluctuating state forever may well be incompatible with the most basic theoretical constraints governing its birth process and later evolution, which rather bespeak of its potential for eternal vitality.

At this point it is necessary to mention that even though the positive- and negative-energy matter particles which were present in the first instants of the Big Bang must be as homogeneously distributed as required to avoid the presence of macroscopic event horizons, the constraint of global entanglement does not impose on matter energy that it be perfectly homogeneously distributed (which would be impossible anyhow, given that, even in a perfectly smooth state of maximum positive- and negative-energy matter densities, the sign of energy would need to vary from one elementary unit of space to another). Thus, there can still be fluctuations in the densities of positive and negative matter energy on all scales in the initial state that emerged from the Big Bang singularity and it is those fluctuations that would give rise to present-day structures.

But in the context where both positive and negative matter energies were, in effect, very uniformly distributed in the first instants of the known Big Bang, as a consequence of the requirement of global entanglement, if the kinetic energies of expansion experienced by positive- and negative-energy observers must everywhere compensate any non-zero energy of matter attributable to local variations of matter energy density, in order that the universe have a null energy, as I proposed in section 4.5, then it becomes possible to conclude that the universe must expand isotropically to a very good degree of precision, even in the absence of an initial phase of inflationary expansion, because under such conditions the expansion rate must be

roughly the same everywhere, as the amplitude of fluctuations in matter energy density is itself very low. Furthermore, if it is actually the case that the expansion is nearly isotropic around every point, due to the requirement that this expansion rate be fixed by the matter density, then it follows that the matter distribution must remain highly homogeneous on the largest scale, as expansion proceeds.

The uniformity of the expansion rate also allows one to deduce that the temperature of the cosmic microwave background should be homogeneous even on a scale larger than the size of the cosmic horizon at the epoch of last scattering, because the absence of macroscopic event horizons is required on all scales and this imposes very stringent conditions on the fluctuations of matter energy density that could be observed, even on the largest scale. In fact, the condition of global entanglement can be expected to exert an even greater constraint on the magnitude of fluctuations in the density of matter energy occurring on a larger scale, given that an overdensity of lower magnitude would be required to produce a macroscopic event horizon on such a scale, as witness the fact that larger black holes have lower mass densities. This means that no smoothing process is required to make the temperature of the cosmic microwave background uniform, because the distribution of matter energy and the expansion rate were mostly uniform on all scales right from the beginning, even if the size of the cosmic horizon decreases more rapidly than the scale factor as we approach the initial Big Bang singularity in the past direction of time, so that regions which are now in contact must have been causally disjoint at the epoch of decoupling (despite the existence of causal relationships between all elementary particles which were present in the universe at this epoch).

When one properly recognizes the limitations imposed by the global entanglement constraint on the initial state at the Big Bang, the horizon problem simply no longer exists and no independent assumption is required to confirm the relevance of the cosmological principle for a description of the early universe. There is no longer any mystery associated with the fact that only one parameter (the scale factor) is required to describe the state of the universe at all but the most recent epoch. In fact, it would now appear that the cosmological principle must be obeyed as accurately as we are considering increasingly larger regions of space, at times increasingly closer to the initial singularity. Despite the enormous densities and the extreme temperature that characterizes the Big Bang, it would therefore be possible to determine the general properties of the initial state with much greater pre-

cision than is usually assumed, despite an absence of knowledge of the exact unified theory that would apply under such conditions. To be sure, the usual assumption that, in order to obtain a homogeneous matter distribution on the cosmic scale it is necessary for the entire observable universe to have been contained within the cosmic horizon at some point in the past (which would be impossible without inflation), can now be recognized as inappropriate and unnecessary.

In the context where it is indeed the magnitude of local fluctuations in the primordial distribution of matter energy which is restricted by the global entanglement constraint, it would also follow that arguments to the effect that topological defects should have been abundantly produced at the Big Bang may no longer be as significant as they used to be. Of course, even from a conventional viewpoint, one must be careful when considering the prediction that there should have formed in the universe a large number of magnetic monopoles or cosmic strings, because the validity of the Grand Unified Theories on which those deductions are based hasn't yet been experimentally confirmed. However, some of those predictions appear to be largely independent from the details of the theories from which they are derived and therefore cannot be ignored. What I have realized is that the relatively low abundance of topological defects may simply be a consequence of the fact that they are very-high-energy objects, similar, in certain respects, to naked singularities of the future kind. The presence, in the early universe, of compact objects that would concentrate such large amounts of positive or negative energy in such small volumes of space may simply be incompatible with the requirement of smoothness of the initial distribution of matter energy which arises from the condition of absence of event horizons that is imposed by the constraint of global entanglement. Indeed, magnetic monopoles are sometimes described as magnetically charged black holes and if this characterization is appropriate, then we should certainly not expect such an object to have developed out of the highly homogeneous distribution of matter and radiation energy that is required to exist in the very first instants of the Big Bang by the condition of global entanglement.

The rarity of topological defects at the end of the GUT era may, therefore, simply be a consequence of the fact that the amplitude of fluctuations in the density of positive and negative matter energies is too small initially to allow the production of highly dense topological defects at such an early epoch. I believe that this constraint, which simply does not exist from a traditional viewpoint, imposes strong limits on the presence of topological

defects. In any case, the fact that the vacuum itself is a much different phenomenon, in the context where its natural, average energy density is actually null, certainly contributes to explain why it is that traditional expectations regarding the cosmological consequences of symmetry-breaking phase transitions are not reflected in what we already have of experimental evidence. Those explanations may not be as satisfactory as the solutions I have provided to other aspects of the inflation problem, but given that, according to the most knowledgeable experts, there is only a very small chance that, in the absence of additional constraints on the initial, pre-inflationary state, the conditions could be met that would allow inflation to occur and to last for a sufficiently long time that it could actually reduce the density of topological defects to acceptable levels, then we may have no choice but to recognize that the constraint discussed above provides a more solid foundation for explaining the rarity of those theoretical objects. In any case, the fact that the physics of topological defects is still relatively uncertain means that the tentative explanation provided here cannot be rejected, even if at this point it is not itself entirely conclusive.

Now, it is sometimes argued that the distribution of matter energy was so uniform at the time when the cosmic microwave background was released that what remains unexplained is really that the temperature was not perfectly smooth and free of any fluctuations initially. But I believe that this smoothness problem is a mere consequence of the fact that we do not properly understand what gives rise to the high level of uniformity of the initial distribution of matter energy. It is only when we assume that perfect smoothness is produced by default that we must invoke a cause in order to explain the fact that there actually existed fluctuations in the density of matter energy as far back in time as the epoch of last scattering. Given that, in the context of the above discussed solution to the horizon problem, it is merely the upper bound of fluctuations in matter energy density which is constrained, then it is to be expected that certain local variations in matter energy density would necessarily be present, as the absence of macroscopic event horizons can be satisfied even when some fluctuations are present. Therefore, if the initial state is still chosen randomly, as it should, it would likely not be perfectly smooth. But we do not need local perturbations that propagate across vast distances in order to give rise to the correlations observed on the largest scales in the temperature of the cosmic microwave background. The cause of the existence of correlations on scales larger than the cosmic horizon is the constraint of global entanglement, which requires some level of smoothness

even in the presence of inhomogeneities, thereby giving rise to a certain uniformity which need not have involved the propagation of effects, as it must have been already present initially.

The usually favored approach to the problem of the origin of primordial inhomogeneities in the distribution of matter energy, which involves assuming that they arise as a consequence of the presence of irreducible quantum fluctuations in the initial distribution of vacuum energy, only makes sense in the context where inflation is assumed to generate an otherwise perfectly homogeneous ‘initial’ state out of a much smaller volume of space. From my viewpoint, what must be explained is not the presence of inhomogeneities, but the overall uniformity of the initial distribution of matter energy, which, in the absence of a specific constraint, should not be observed. The natural configuration for the initial state is not one of perfect smoothness and there is no need to invoke a particular effect to generate the observed fluctuations, which are allowed to be present on any scale, as long as they do not violate the condition of global entanglement. What is truly remarkable is that the spectrum of fluctuations in the density of matter energy which is actually deduced from observations of cosmic microwave background temperature fluctuations is a (nearly) scale-independent spectrum of the Harrison-Zel’dovich type (for fluctuations larger than the scale of the horizon at the epoch of decoupling) while this is the only spectrum which, according to specialists, does not allow the magnitude of fluctuations to diverge on either large or small scales and which, therefore, does not give rise to the creation of a large number of primordial black holes, while those are precisely the conditions which are required by the theoretical constraint of global entanglement.

In this context, it must be clear that the requirement that the spectrum of density fluctuations be scale invariant is not a unique property of inflationary cosmology. Also, the idea that only inflation allows the initial perturbations in the distribution of matter energy to begin oscillating (between compressions and rarefactions) at the same epoch of cosmic time, as required to explain the existence of harmonics in the spectrum of cosmic microwave background temperature fluctuations, may not be justified, because any scale-invariant spectrum of fluctuations in the distribution of matter and radiation energy which would be present initially on scales larger than the horizon, would have the same consequences. But this is precisely what we can expect to occur as a consequence of the constraint of global entanglement, in the presence of negative-energy matter. Thus, even in the absence of inflation, the required fluctuations in matter density would already exist

in the ‘initial’ state and those fluctuations would all begin oscillating as soon as they are encompassed by the cosmic horizon, which means that on a given angular scale all maximum compressions and rarefactions would be reached at the same time.

The fact that the power spectrum of density fluctuations appears not to be exactly scale invariant, on the other hand, may not be a consequence of inflation either. It is usually assumed, in effect, that it is because the gravitational repulsion driving inflation becomes weaker as the process is occurring, that the smaller-scale fluctuations, which are produced later by the inflation process, have a smaller amplitude. But what may be happening, instead, is that, given that there remained only a negligible portion of baryonic negative-energy matter following the early annihilation of matter with antimatter, then one can expect that inhomogeneities in the negative-energy matter distribution were less developed on smaller scales, because, even though there existed as much negative vacuum-dark-matter energy before and after the annihilation process, no significant small-scale fluctuations could develop in this matter distribution in the absence of baryonic negative-energy matter (as I have explained in section 4.4). Thus, while negative vacuum-dark-matter energy would be inhomogeneously distributed on larger scales at the epoch of last scattering, it would be more homogeneously distributed at the same epoch on smaller scales. But given that negative-energy matter inhomogeneities do contribute to the presence of temperature fluctuations in the cosmic microwave background at the epoch of last scattering, then there naturally arises a deficit of fluctuations on smaller scales, even when inflation is not assumed to be responsible for producing the inhomogeneities that were present on both smaller and larger scales at the epoch of last scattering. I believe that this is the true reason why the spectrum of density fluctuations is slightly tilted (with a power-law index smaller than one) compared to a pure Harrison-Zel’dovich spectrum.

It should be clear, therefore, that negative-energy matter does have an effect on the observed properties of cosmic microwave background temperature fluctuations. Of all the measurements concerning the spectrum of CMB temperature fluctuations, the only ones which would remain mostly unaffected by the presence of negative-energy matter are those which regard a determination of the angular scale of fluctuations from which are derived the average density of positive matter and vacuum energy, because the trajectories of positive-energy photons are not affected by the presence of negative-energy matter on the largest scale (particularly when its density is as uniform as it

must be on smaller scales, following the initial annihilation of matter with antimatter). The appropriateness of this conclusion appears to be confirmed by the fact that estimates of the density of positive-energy matter (both visible and dark) based on measurements of the spectrum of CMB temperature fluctuations produce a value largely equivalent to that which is derived using more direct methods. I must acknowledge, however, that discrepancies have emerged more recently (see in particular Ref. [37]) between the average density of matter derived from CMB measurements and that which is inferred from weak gravitational lensing of nearer structures and it is possible that we will only be able to resolve those difficulties once the various effects of the presence of negative-energy matter on the process of structure formation, in both the early universe (before decoupling) and at a more recent epoch, are properly taken into account and the consequences of a variation of the magnitude of the cosmological constant are fully worked out.

As I explained above, we cannot expect that the inhomogeneities which were present on smaller scales in the early distribution of negative-energy matter were as developed as those which were present at the same epoch in the positive-energy matter distribution. But given that in the presence of negative vacuum-dark-matter energy there should nevertheless be more inhomogeneities in the matter distribution than we could attribute to positive-energy matter, at least on a very large scale, and therefore more fluctuations in the temperature of the cosmic microwave background, then it would appear necessary to revise the magnitude of density fluctuations attributable to positive-energy matter downward, as I mentioned in section 4.4. The presence of negative-energy matter inhomogeneities would, therefore, allow to ease the tension that emerged from even more recent observations [38] which indicate that the current magnitude of fluctuations in the positive-energy matter distribution is smaller than the value deduced from CMB fluctuations extrapolated to the present epoch (a problem known technically as the S_8 tension). Indeed, if the actual magnitude of density fluctuations in the distribution of positive-energy matter at the epoch of last scattering can be assumed to be smaller than the value currently deduced from observations of CMB temperature fluctuations, due to the fact that some of the observed fluctuations are attributable to the presence of negative-energy matter, then it is natural to expect that the current magnitude of fluctuations in the positive-energy matter distribution, obtained by evolving the fluctuations which existed at the epoch of last scattering forward in time to their present state, should be smaller than that which is actually observed, even when one

recognizes that the presence of negative-energy matter must have accelerated the process of structure formation in the early universe (on smaller scales) and at later times (on larger scales).

It is certainly possible that certain characteristic features of the spectrum of cosmic microwave background temperature fluctuations, distinct from those mentioned above, could also be explained by taking into account the effects that would be attributable to the presence of very large scale inhomogeneities in the distribution of negative vacuum-dark-matter energy, either in the initial state, or in the space through which radiation propagates before reaching our detectors. One very clear implication of a cosmological model based on the ideas developed in this report, however, is that it is no longer necessary to assume that gravitational waves must have been produced as a result of the stretching of space produced by inflation, starting from a highly inhomogeneous initial state, and therefore it may be futile to search for an unmistakable sign of the existence of those gravitational waves in the polarization of cosmic microwave background radiation.

4.10 A criticism of inflation theory

Now that I have provided alternative solutions to all aspects of the inflation problem, I would like to offer a constructive criticism of inflation theory itself and explain why it may no longer constitute an appropriate response to the most enduring difficulties facing theoretical cosmology. It must be clear, however, that I do not claim to have proven that inflation theory is wrong or that the phenomenon it describes did not occur. Indeed, what I have shown is simply that inflation is no longer *necessary* to solve the flatness problem and that there may be no substance to the related problem of matter creation, outside of inflation theory itself. I then proceeded to explain that an alternative solution to the horizon problem and the related problem of smoothness can be formulated that may also go some way in solving the related problem of topological defects. But this does not mean that the hypothesis that there occurred an early phase of accelerated expansion is not valid and that we should no longer expect something like inflation to have happened, only that the existence of such a phenomenon may not be required for explaining the puzzling features of the universe which are giving rise to the inflation problem. I find it significant, however, that of all the major difficulties facing cosmology, the cosmological-constant problem is the one for

which inflation theory was never allowed to provide an appropriate answer, as this is definitely an issue that can only be addressed in the context of the generalized gravitational theory proposed in chapter 2. There may have been truth, then, to the long forgotten suggestion that the same insights which would turn out to be required in order to solve the cosmological-constant problem may allow to do away with the other outstanding difficulties of cosmology, which would otherwise need to be addressed by resorting to inflation theory.

Thus, while inflation may not be invalidated as a theory, it appears that there are more natural solutions, based on more unavoidable theoretical and observational constraints, not only to the inflation problem itself, but to certain other important issues as well. In fact, I'm now in position to provide satisfactory answers to practically all the remaining outstanding problems of theoretical cosmology, including the problem of the origin of the arrow of time and that of the nature of dark matter. But what should motivate one to recognize the necessity for an alternative approach to cosmology such as the one I have proposed in the preceding sections, is the fact that even some of the originators of inflation theory have, more recently, expressed doubts concerning the usefulness of the theory for solving any of the problems to which it was originally believed to provide a satisfactory answer, because inflation itself requires very unlikely initial conditions to be initiated and to give rise to the desired outcome. Those criticisms, however, are usually overlooked, because of what appears to be the overwhelming evidence in favor of inflation that was provided by the discovery that the universe is flat and that the distribution of matter and radiation energy is homogeneous above the scale of the horizon (as revealed by observations of the cosmic microwave background).

Indeed, it is definitely the fact that a universe with a density parameter $\Omega_0 = 1$ was always favored by inflation, even at a time when it appeared that lower values of Ω_0 were favored by observations, that is responsible for having transformed inflation theory into the paradigm it is today, when it was later found that this parameter is, in effect, equal to unity and space actually is perfectly flat on the largest scale. But given that I have shown that, in the presence of negative-energy matter, when we appropriately require the universe to have zero energy, the specific expansion rates of positive- and negative-energy matter are required to be critical to an arbitrarily high degree of precision, based merely on the assumption that an observer must be present in the universe to measure those parameters, then it would appear

that flatness is not valid as a confirmation of the validity of inflation theory, but is actually a generic property of any observable universe obeying the known principles of physics. In this context, what was wrong was really the early expectation that, by default (in the absence of inflation), space should be observed to be highly curved at the present time, given that perfect flatness requires very unlikely initial conditions on which there appeared to be no constraint. In fact space must be perfectly flat at all times in an observable zero-energy universe, but what I have tried to explain is that inflation has nothing to do with that and therefore flatness does not provide an unmistakable confirmation of the validity of inflation theory.

Of course, if all I had done was to show that the flatness problem does not occur, even in the absence of inflation, when the universe is required to have a null energy and observer selection effects are taken into account, then it would not be possible to conclude that inflation is unnecessary, because there would still be a problem associated with the observed large-scale homogeneity of the early matter distribution. But given that, when negative-energy matter is present in the primordial state and one recognizes the necessity for all elementary particles in the universe to be causally related to one another, it actually becomes necessary for the initial matter distribution that existed in the very first instants of the Big Bang to be smooth enough that no macroscopic event horizon is present, then it follows that the overall homogeneity of the temperature of CMB radiation is no longer a fact in need of some explanation.

The presence of a high initial density of positive-energy matter, on the other hand, no longer requires a hypothetical, post-inflation reheating process, dependent on very specific conditions, to have occurred, once one recognizes that there is no need for matter to be created out of nothing during the Big Bang, given that it must have already been present prior to the initial maximum density state, before being submitted to a quantum bounce. It is actually only when we assume that there occurred an early phase of accelerated expansion that we need matter to be created at some point during the Big Bang and therefore it would certainly be wrong to assume that inflation constitutes an absolute requirement for the existence of matter in our universe. In fact, when all the dust has settled, it appears that not much evidence remains to possibly confirm that inflation really occurred. But again, that does not mean that none of the theoretical motivations behind inflation were justified, merely that inflation is not required to produce the apparently unlikely ‘initial’ conditions which were previously thought to remain

unexplained outside the realm of this theory.

Concerning the flatness problem, however, it transpires that if the *specific* expansion rates of positive- and negative-energy matter were fixed to their critical value by inflation alone, while only the magnitudes of the initial, average densities of positive and negative matter energy were required to be equal by the condition of null energy (so that the gravitational potential energies and the kinetic energies of expansion would be left unconstrained by the same condition), then it would be difficult to explain how the average, specific densities of the two opposite-energy matter distributions could remain mostly the same following inflation (but before the early annihilation of baryons with antibaryons), as required if the cosmological constant is to not be much larger than it appears to have been at this remote epoch (which would also imply a larger present value). There is no reason, in effect, to assume that the outcome of inflation would be exactly the same from both the viewpoint of positive-energy observers and that of negative-energy observers, while this is required if the specific densities are to be of similar magnitude following inflation. Thus, it would appear that the explanation of flatness provided in section 4.5 is actually an absolute requirement, under such conditions, and cannot merely be considered an alternative possibility.

In fact, the difficulty discussed here would be even worse if one was to assume that the early phase of inflationary expansion was not merely produced by the presence of some hypothetical inflaton field, but was instead the outcome of the existence of a large average value of vacuum energy density in the initial Big Bang state, because in such a case the process would necessarily have opposite effects of considerable magnitudes on the expansion rates measured by observers with opposite energy signs. Thus, while the space experienced by a positive-energy observer could be driven to inflate exponentially, the space experienced by a negative-energy observer may actually be made to collapse back into a singularity, which means that the ‘initial’ densities of positive- and negative-energy matter measured by such an observer following inflation would remain maximum, while those measured by a positive-energy observer (before reheating) would become negligible, which once again is not quite compatible with observations, which indicate that the expansion rates and the spatial volumes experienced by opposite-energy observers are still similar at the present epoch, due precisely to the small value of the cosmological constant. If inflation is a product of the energy of zero-point vacuum fluctuations, it would therefore appear that additional fine-tuning, of a kind that hasn’t even been considered yet, would be required

to make the theory viable.

What must be clear is that there is only one (positive or negative) value for the average density of vacuum energy at any given time and if the magnitude of this value is too large for too long a period of time initially, then there may be conflict with observations, even independently from whether matter is present initially (as I'm assuming) or is created through reheating after inflation (as must be assumed from a traditional viewpoint). In the context of the approach I favor, this problem does not exist, because the average value of vacuum energy density is not affected by the changes which are taking place in the vacuum and depends only on the difference between the scale factors experienced by opposite-energy observers. What's more, for a zero-energy universe, when it is recognized that negative-energy matter must be present in the initial state, there is a requirement for space to be perfectly flat and vacuum energy to be null initially, on the global scale, from both the viewpoint of positive-energy observers and that of negative-energy observers, and under such conditions, even if a non-zero average density of vacuum energy may be required to compensate any difference that may develop at later times between the average density of positive matter energy and that of negative matter energy, the anthropic principle provides sufficiently strong a constraint to allow one to expect that the current value of the cosmological constant should nevertheless be as small as it is observed to be.

Concerning the solution potentially offered by inflation theory to the horizon problem and more specifically to the unexplained uniformity of the initial matter distribution, it was already pointed out by Roger Penrose that the usual assumption, to the effect that inflation would take the universe from a highly inhomogeneous state to a perfectly smooth one, appears doubtful in the context where the initial state would, in effect, be characterized by a maximum gravitational entropy. But those remarks were made before we even had a theory of gravitation that allowed for the presence of negative-energy matter in the initial state. From the viewpoint of the developments introduced in chapter 2, it would seem that it is definitely impossible to assume that a universe with an arbitrarily large gravitational entropy could be rendered homogeneous through accelerated expansion, as the potential for ever more polarized positive- and negative-energy matter distributions is unlimited, just like the initial amount of matter itself. The opposite-energy black holes that could be present in the initial state, if it was not for the limitation exerted by the constraint of global entanglement on density fluctuations, could be as massive as the radius of curvature of the universe is

large and would concentrate all the matter in the universe in their gravitationally repelling singularities, which means that no amount of expansion could ever result in a homogeneous matter distribution. Thus, if negative-energy matter does exist, it seems that inflation alone could not prevent the initial distribution of matter energy from being highly inhomogeneous and this provides additional motive to believe that the process is not necessary, even if it still cannot be ruled out that it might have occurred.

One particular aspect of the horizon problem which is currently believed to have been solved by inflation theory has to do with predicting the spectrum of density fluctuations in the initial positive-energy matter distribution. Indeed, the fact that observations of the cosmic microwave background have revealed a nearly scale-invariant spectrum of density fluctuations, of the kind that is predicted by inflation theory, is often considered to provide the strongest empirical evidence of the validity of this theory. But given that a scale-invariant spectrum of density fluctuations of the Harrison-Zel'dovich type would be typical of any theory according to which the presence of macroscopic black hole event horizons is forbidden on all scales and space itself does not have a characteristic scale, as is the case for a universe with a critical energy density and an infinitely large radius of curvature, then it seems that a well-behaved cosmological model, that would describe a universe with null energy, would also predict a scale-invariant spectrum of density fluctuations, given that it would require black holes to be absent in the initial Big Bang state and space to be flat on a global scale. Therefore, once again, the observations cannot be assumed to provide a definitive confirmation of inflation theory and this conclusion is reinforced in the context where, even the fact that the spectrum of density fluctuations is not perfectly scale-invariant can be explained as being a consequence of the absence of baryonic negative-energy matter, as I explained in the preceding section.

Now, it has been hailed that the fact that certain versions of inflation theory may allow the 'true' universe to be comprised of many regions like the known universe, separated by arbitrarily large portions of inflating space in which new 'universes' like our own are born all the time, could be a positive development, given that it seems increasingly more likely that some properties of our universe are constrained by the weak anthropic principle. Indeed, one of the implications of the existence of such otherwise unexplained properties is that it makes plausible the idea that there must be more than one possible instance of physical reality, so that the anthropically constrained universe we observe can exist as a mere possibility whose improbable nature

need not be explained by appealing to divine intervention. Some of us, however, appear to favor, for some mysterious reason, that all of those realities, instead of just existing on their own, be somehow tied (however loosely) to the universe we do experience, as if this was a requirement of the multiverse concept. This conception of the multiverse has been appropriately renamed the ‘megaverse’ by Leonard Susskind and now enjoys respectable status as if it had been proved right by the ‘successes’ of inflation theory. But given that we may now have to reconsider the degree of inevitability of the phenomenon of inflation, it would appear that all this extraneous amount of inflating vacuum may no longer be as appealing as it once was.

In any case, one thing should be clear and it is that eternal inflation is not necessary for making the multiverse concept a viable notion and in fact the emergence of a megaverse concept in inflation theory may actually constitute a problem for this approach to cosmology, given that it may indefinitely postpone the moment at which the global entanglement of this whole enlarged universe would occur in the past, while this must be considered a necessity, as I explained at length in the preceding section (if global entanglement is required only within the bubble universes, then the megaverse itself could not be assumed to exist as an ensemble of causally related parts in the first place). This remark is all the more relevant in the context where, unlike the megaverse, the existence of an arrow of time (understood as being a consequence of global entanglement) is an observable fact with undeniably real consequences.

The most enduring problem facing inflation theory, however, remains the fact that it is still as difficult today as it was back when the model was introduced decades ago to identify what is the deep principle from which it would emerge as an unavoidable aspect of physical reality. If such a foundation cannot be developed, we will perhaps eventually need to recognize that what was provided by inflation was a solution that was useful merely because of the absence of a better alternative. It is not appropriate, in the context where a better explanation of facts is available, to just keep adjusting the free parameters of a theory which is supposed to determine the very boundary conditions applying to the universe as a whole. At this point it is not just questionable whether inflation can actually solve any of the outstanding problems of Big Bang cosmology, it is even uncertain whether it is still possible to assume that the phenomenon occurred at all. Under such circumstances, only our inherent resistance to paradigm change may prevent us from acknowledging the eventual failure of the theory. But if there is any reason to believe that

inflation, in effect, did not occur, it would have to be the fact that it is not merely a single one of the difficulties originally assumed to be solved by the theory that can be explained away in the context of negative-energy-matter cosmology, but nearly all aspects of what was once the inflation problem.

Chapter 5

Quantum Theory and Causality

5.1 The problem of interpretation

In the preceding chapters of this report, I offered original solutions to several outstanding problems in the fields of gravitational physics and cosmology, which were all based on an alternative interpretation of the concepts of time reversal and negative energy. First of all, I introduced a generalized, classical theory of gravitation that is consistent with the possibility that elementary particles could exist that would propagate negative energy forward in time. Based on the understanding that it is necessary to distinguish between a fundamental, bidirectional concept of time direction associated with the propagation of elementary particles and the classical, unidirectional concept of time direction associated with thermodynamic irreversibility, I was then led to introduce a more consistent formulation of the time-reversal symmetry operation that was shown to be relevant to a description of the fundamental states of matter particles on the quantum gravitational scale. I also showed that the hypothesis that negative-energy matter was present alongside positive-energy matter in the first instants of the Big Bang allows the formulation of a satisfactory solution to the problem of the origin of thermodynamic time asymmetry as being the outcome of a certain condition that must be imposed on this initial state in order that all the elementary particles present in the universe be causally related to one another. But given that the bidirectional concept of time that underlies this approach constitutes a challenge to our traditional conception of causality and the idea that causes always precede their effects, then it becomes necessary to introduce a revised concept of

causality that would allow to take into account the time-symmetric nature of elementary particle processes. This is one of the objective I'm seeking to achieve in this chapter.

Part of the progress I have accomplished concerning this issue actually emerged from an investigation of the significance of certain puzzling aspects of the currently favored interpretation of quantum theory which did not appear to be connected with the issue of time directionality, but which were of interest in their own right. Yet, some of the results I have obtained regarding the issue of time directionality in gravitational physics turned out to be necessary for developing a solution to the remaining problems that still stand in the way of a truly consistent interpretation of quantum theory. Thus, in the present chapter, I would like to address not only the question of time directionality as it arises in a quantum mechanical context, but also, and more specifically, the important problem of the interpretation of quantum theory itself. I will show, in particular, that it is now possible to provide a realist picture of quantum phenomena that does not violate the principle of local causality, even though it is not incompatible with the consequences of quantum entanglement and the existence of non-local correlations. This improved understanding will then be used to provide a complete and definitive solution to the quantum measurement problem that allows to explain the emergence and the persistence of a quasiclassical world.

But some of the most significant contributions I will offer in this chapter consist in showing that there really is a problem with some aspects of our current understanding regarding quantum physics. Two categories of questions I will try to address more specifically are not always distinguished from one another and together constitute the problem of the interpretation of quantum theory. I will explain, however, that they must be considered as independent questions in need of separate answers. There is thus, in effect, a problem of interpretation concerning the mathematical framework of quantum theory in general and the distinctive features of quantum physics, which are mainly the use of probability amplitudes instead of classical probabilities (a remark which becomes significant once it is recognized that, in a time-symmetric context, quantum field theory is simply a more fundamental instance of statistical mechanics) and the existence of entangled states, which are both unavoidable features of physical reality. It is not clear, from a traditional perspective, how entanglement, in particular, and the kind of non-locality it makes possible, is to be understood in a manner that is consistent with the classical principle of local causality. It is commonly believed

that the problem here does not have to do with the inappropriateness of our interpretation, but with the inappropriateness of our traditional approach to understanding reality. Yet, if this posture may be legitimate for what concerns probability amplitudes and the existence of quantum interferences in general, I will show that it is not justified for what concerns the problem of quantum non-locality, which is actually a question in need of an answer. Thus, answers will be offered to a problem which I call the *quantum reality problem* and which includes the problem of quantum non-locality as a particular aspect.

This problem must be distinguished from the associated problem that is usually referred to as the *quantum measurement problem*. Those who are not actively working on this particular problem often believe that it too may not be real, or, alternatively, that it was entirely solved by more recent developments that showed how the evolution of a quantum system is affected as soon as it becomes entangled with certain irreversible processes taking place in its environment. But, as a handful of researchers have already understood, this opinion is not warranted and even though real progress has indeed been achieved in trying to solve the quantum measurement problem and more generally the problem of the emergence of ‘quasiclassicality’, some related questions remain unanswered and it is precisely those I will address. However, if you happen to be among those who are convinced that there is no longer a problem with quantum measurement, then I would ask you a very simple question: what is the cause of the irreversibility that characterizes the evolution of the environment degrees of freedom with which a quantum system becomes entangled during measurement and which is necessary for explaining decoherence? Clearly, an appropriate answer to that question must be provided before the problem may be considered to have been solved and this is what I have tried to achieve in the previous chapter. But, as I will explain, that is not the only difficulty. In order to clarify this complex situation I will, therefore, need to draw on the insights I have gained while solving the problem of the origin of time asymmetry, but I will also have to build on the insights I have gained while solving the quantum reality problem, which illustrates how important it was, in effect, to first solve that perhaps more intangible problem.

It is quite remarkable that, in order to answer those two categories of questions, it is possible to rely on the most appropriate of the already existing mathematical frameworks within which quantum physics is currently formulated and it is not necessary to alter the foundations of the theory. I

must immediately point out, however, that there is something terribly wrong with the often met remark to the effect that, choosing which of the existing interpretations of quantum theory is the correct one is a mere matter of taste, given that they are all mathematically equivalent and therefore all constitute equally valid proposals, which all agree with observations. Before any progress can be achieved, what needs to be understood is that most interpretations are not equally appropriate, but rather all equally incorrect or incomplete. It would be misleading, therefore, to argue that the problem is that there are too many viable candidates for an interpretation of quantum theory, because in fact *none* of the currently available proposals is fully consistent, either from a logical viewpoint, or regarding the requirement that the obtained theory be compatible with all observable aspects of physical reality. This state of affairs can only mean one thing and it is that further progress is required to formulate the one interpretation that will meet both of those requirements. I believe that the original results I have unveiled in the preceding chapters of this report provide some of the missing elements which are required to achieve just such a leap forward in our understanding of quantum physics, which will, at last, allow it to become a fully coherent theory.

5.2 A simple analogy

One particular event from my early years in elementary school contributed more than any other in developing my awareness of the deep structure of physical reality. I do not remember much about the many events that happened during the period of my life when I was acquiring many of the skills which I'm still using today (like writing and calculating), but I still remember perfectly well that when I was about eight years old my teacher once gave me and each of the other children in my class a few copper wires, an electrical battery, and a tiny light bulb, with as a mission to figure out how to produce light using only those components. This may seem like an easy task and most of the kids did, in effect, manage to achieve the assigned objective quite rapidly. Yet, even though I was usually considered a fast learner in most traditional academic disciplines, I really had trouble finding out how to obtain the desired result. I believe that this is because, even as a kid, I always preferred to actually understand things, rather than simply be satisfied with learning about the finished answers I was proposed. So, rather

than just trying to combine the elements in every possible way and be satisfied once I had accomplished my homework, like the other kids, I tried very hard to understand what the rule could be that would justify that a certain arrangement does, in effect, produce light. I don't know why I had such an inclination, but it has remained with me all my life as I began to develop an interest in the sciences and learned about the unsolved mysteries of modern physics. What I have since realized, with retrospect, is that somewhere in this simple laboratory experiment was hidden the answer to some of the most enduring problems facing fundamental theoretical physics.

The first lesson I have retained from this experience with the light bulb and the battery is that there is a polarity to all relevant physical properties. The battery has a positive and a negative pole, and so does any light bulb, and it is only by taking this aspect into consideration that one is allowed to understand what constitutes a successful configuration for producing light. In the first two chapters of this report I have discussed at length how this aspect is relevant even in a gravitational context, where the sign of action is the decisive physical property that is involved in determining the attractive or repulsive nature of the gravitational interaction between two particles. I also emphasized the purely relational nature of any polarity, whether it regards the sign of electric charge, or the direction of propagation in time of a particle. Only the difference or the identity of any such property of a system with respect to that of another system has a physical significance. But the most difficult part in devising an electrical setup that works consisted in understanding the role of the wires. What is essential to learn, in effect, is that the experiment can only work if the wires are arranged so as to form a circuit that goes from one pole of the battery into one pole of the light bulb and then back from the opposite pole of same light bulb into the unconnected pole of the same battery. I only realized that this must be so when I carefully examined the light bulb and saw that there is a special kind of wire inside of it that connects the two poles from within, thereby suggesting that, for some reason, the setup must be closed on itself.

But after I came to understand the requirement for the setup to form a closed circuit, I was not only faced with the problem of understanding why only such a configuration would produce the desired outcome, but also with the difficulty of understanding what was the role of the battery in allowing light to be produced by the light bulb. In other words, I had trouble understanding in which way the role of the battery could be distinct from that of the light bulb, despite the fact that both components were connected along

the same circuit. It is only much later in my life that I learned that what the light bulb does, from a fundamental viewpoint, is simply dissipate the energy that is stored in a useful form in the battery, as a result of the friction that is exerted on the electrical current that flows through the part of the circuit located inside this light bulb. This is an expression of the second law of thermodynamics in its purest form. The very objective of producing the circuit is to allow the current to dissipate the energy contained in the battery, by producing an enormous number of light particles that expand out of the system irreversibly. The whole mystery associated with the second law of thermodynamics and the irreversibility of time is contained in this little experiment and with it the solution to the quantum measurement problem. Any circuit that produces a useful outcome, i.e. one that is observable and which has an effect on its surroundings (the light turns on), must dissipate energy that was originally present in the universe in a well-ordered configuration.

I will eventually explain what is the essential role of irreversibility in allowing the emergence of the quasiclassical character of reality that is revealed by any process of quantum measurement, but first, I would like to point out the profound significance of the property of closure that is imposed on any operational electric circuit. For anyone who works as an electrician, the notion of a closed circuit is omnipresent, but it is also forgotten somehow, in a practical context, where one always works with pairs of polarized wires in which the two branches of a circuit are always contained in a single cable that invariably goes from energy source to appliance, over large distances, as if what was involved was one single flow from source to sink, similar to the flow of water inside a pipe. Thus, it is easy to forget that one is always dealing with a closed circuit, however complicated it might be. I believe that what explains some of the difficulties we encounter in quantum physics is that we have always learned to work only with pairs of ‘polarized wires’ and this is why we fail to understand that what we are dealing with, in general, is not a single process that unfolds from initial conditions to final measurement, from ‘source’ to ‘sink’, but really a *closed causal chain*, similar to the closed electrical circuit of my childhood experiment.

It is the fact that, for some reason which will be discussed later, we are always working with portions of a ‘closed circuit’ which are highly stretched and extended along the unidirectional dimension of time and whose polarized components are constrained to evolve along very similar trajectories, that explains that we have been allowed to ignore the fact that we are actually always dealing with two processes which are the oppositely polarized portions of a

causal chain that closes up on itself, like a functioning electrical circuit. We always model very long causal chains, similar to electric cables, that extend not along a distance in space, but along the unidirectional dimension of time and for that reason we have never realized that what the polarized character of this causal chain really means is that we are dealing with a ‘closed circuit’. I will explain how the simple analogy discussed here can be developed into a rigorous interpretation of experimental facts that allows to provide, not only a consistent explanation of the persistence of quasiclassicality, but also a realist and fully intuitive picture of quantum phenomena that reproduces their non-local character without violating the principle of local causality. When we will reach that point, it will become possible to actually understand why it is, in effect, that the causal chain associated with the history of the universe as a whole closes up on itself like any electrical circuit that produces light.

5.3 Time-symmetric causality

It is somewhat strange that it is Richard Feynman himself who once remarked that one question he believes to be unanswerable or unscientific is the one that asks why it is that we are allowed to guess from one part of the universe what the rest of it will do? Indeed, it has become pretty clear to me that, if this is possible, it is simply because things propagate, not just in space, but also in time and as Feynman himself was one of the first to understand, not just forward in time, but also backward, from the future toward the past. In fact, this is the essence of causality. The events that form the universe are all related to one another and to nothing else by the network of causal relationships that is established by the propagation of elementary particles between those events, across spatial distances and through time, both forward and backward. The results I have discussed in section 4.9, regarding the role of the constraint of global entanglement in giving rise to thermodynamic time asymmetry appear to confirm that the existence of causal relationships, established through local contact and propagation, is, in effect, essential for a consistent description of physical reality.

Another significant conclusion from the preceding chapters of this report, that concerns time directionality more specifically, is that a distinction must be made between the traditional concept of time direction associated with the thermodynamic arrow of time and a more fundamental concept of time direc-

tion, which has to do with the direction of propagation in time of elementary particles and which merely distinguishes two opposite directions without favoring one of them in any particular way (under most circumstances). Thus, even at a semi-classical level of description, there already emerges a notion of bidirectional causality more fundamental than the classical, unidirectional concept of causality, according to which causes always precede their clearly distinguishable effects in the same unique direction of time. From this more fundamental viewpoint, there is no longer an absolute distinction between causes and effects and all that one can meaningfully ask is whether a certain event does exert an influence on another event taking place at a different time, either earlier in the past or later in the future (which would affect the probability that one of the events is observed when the other is).

Those who have seriously examined the question usually recognize that the idea that the present can influence the future, but not the past, is not entirely correct. Indeed, when calculating correlation probabilities, we must take into account the effect of the future on the present whenever there are antiparticles in the final state, because antiparticles are most appropriately described as particles propagating from the future. In fact, there appears to be no real distinction, at the *elementary* level of description, between the past and the future, and despite the fact that the future remains unknown to present observers, it is as unique as the past and we are merely discovering what that future is as we progress irreversibly towards it. What we consider to be our control over the future is not a complete illusion, of course, because there are correlations between what we do now and what happens later, but those correlations are no more real than the more subtle correlations which are predicted to occur by quantum theory and which arise from the effects exerted by certain present events on certain events in their past. As I will explain below, it is merely the fact that entropy is always growing only in the future that makes it look like we only have control over the future, while the future itself appears to exert no influence on the present and the past, because this is what explains that the present can exert multiple recognizable effects on the future, while the opposite is so unlikely as to be virtually impossible. But this is also why we can have no information about the future, while every future outcome appears to be possible, which makes it seem like we have a certain freedom over what future outcome we choose to generate that does not exist for the past, as if causality could only operate from the present toward the future.

What I'm suggesting, therefore, is that, at the level of elementary par-

ticles, where thermodynamic time asymmetry is not a meaningful concept, causality is not constrained to always operate from past to future, which means that causes and effects cannot be distinguished based merely on the sign of the time interval between the events they relate. Thus, while it may still be necessary to assume that causes precede their effects, this can actually be achieved in any of the two directions of time in which the particles conveying the effects are propagating. At a fundamental level of description there simply is no restriction regarding the direction in which causality operates and this means that, from the unidirectional-time viewpoint, effects can actually precede their causes. But instead of saying that, under certain circumstances causes may actually follow from their effects, while effects would give rise to their causes, it is more appropriate to define causes and effects in a more fundamental bidirectional way, so that effects can propagate either forward or backward in time, but always in the direction in which the particle that generates the effect is propagating in time.

The absence of an *absolute* distinction between causes and effects does not mean that the relativistic concept of a future light cone, clearly distinct from its past equivalent and defining the causal structure of spacetime, is wrong. But it does mean that there is no *a priori* reason to differentiate the structure that arises as a consequence of the limits imposed on the propagation of causal signals in the future from that which would arise as a consequence of the constraints imposed on the possible propagation of causal signals in the past. Yet, it would be incorrect to argue that only correlations exist between various past and future events at a fundamental level of description and that causes cannot be distinguished from their effects in any way, because what the bidirectional, or time-symmetric nature of causality implies is merely that there is no absolute (non-relationally defined) distinction between the past and the future for what concerns the propagation of effects and that only a relatively defined notion of time directionality is still involved, from a semi-classical perspective, that allows to distinguish a direction in the propagation of effects.

What must be understood is that the only invariably true notion of causality is the time-symmetric, or bidirectional one, while classical, thermodynamic causality, or unidirectional causality, is valid only as a consequence of the existence of the constraint of low entropy that applies on the initial conditions at the Big Bang and is not a fundamental property of nature. In fact, what I have shown in section 4.9 is that a certain condition of local causality that is not *a priori* asymmetric in time can be used to explain the

observed thermodynamic time asymmetry from which unidirectional time and classical causality emerge. Later on, I will discuss the role played by the constraint of global entanglement (which I previously identified as the ultimate cause of thermodynamic time asymmetry) in constraining classical causality to operate only in the future direction of time. It is already possible, however, to appreciate the fact that the direction of time relative to which entropy grows and information flows is independent from the direction of propagation in time of the particles involved in producing such a change. The direction of propagation in time of an elementary particle, which determines its particle or antiparticle nature, merely allows to assess whether the particle propagates an effect in the past or in the future, while the flow of information is a higher level property that is fixed merely by the macroscopic boundary conditions imposed on a process, regardless of whether it involves particles or antiparticles.

Thus, a classical, unidirectional causal chain can be differentiated by the fact that it invariably involves a unique event in the past exerting a recognizable effect on multiple spatially separated events in the future¹. In this particular sense, it transpires that unidirectional causality does, in effect, always operate from past to future in our universe, as no single event in the future has ever been observed to exert a unique recognizable influence on multiple physically separated events in the past that would actually involve a flow of information from that future toward its past. But, again, this does not mean that a future event cannot influence a past event, merely that this cannot occur in a way that would allow the formation of mutually consistent records of the future. It is not causality, in the fundamental sense, that is asymmetric in time, but the making of records by which it is usually made manifest. This asymmetry has already been recognized as arising from the existence of a thermodynamic arrow of time associated with the continuous decrease of entropy that takes place in the past direction of time.

¹In fact, as emphasized by Lawrence Sklar [39], it is not always possible to associate the asymmetry that distinguishes causes from effects with thermodynamic time asymmetry, other than by relying on the property of parallelism of the direction of time that allows to ‘project’ the thermodynamic asymmetry characterizing certain causally related pairs of events onto other pairs of causally related events where the cause is not so easily differentiated from the effect (as in the case of certain mechanical or astronomical processes where friction and dissipation are not manifestly apparent). But this only strengthens the validity of the position I’m defending, to the effect that causality is not, by necessity, an asymmetric property.

Therefore, the only mystery regarding the apparent absence of causal chains that would run from the future toward the past does *not* have to do with a real absence of such phenomena, but with the fact that future causes do not exert the same kind of recognizable consequences on the past, as past events exert on future events. The absence of recognizable effects, at an earlier time, from an event that would have taken place at a later time, can always be attributed, not to the absence of backward-in-time causality and to a *fundamental* character of unidirectionality, but to the fact that entropy always increases only in the future direction of time, while records can only be formed in the direction of time relative to which entropy actually increases. Given that, in chapter 4, I have explained that the thermodynamic arrow of time is not a fundamental property of nature, but arises from a condition regarding the homogeneity of the initial distribution of matter and radiation energy at the Big Bang (in the presence of anomalously gravitating negative-energy matter), then it clearly follows that there is absolutely no rational motive to argue that backward-in-time causation is forbidden in our universe. Indeed, all the observable properties of naturally occurring processes can be explained without relying on this assumption, while the requirement of a relational description of the direction of propagation in time implies that backward causation must exist at a fundamental level, given that it cannot be distinguished in any absolute way from forward-in-time causation.

I believe that it is, again, our failure to recognize the full significance of Feynman's (or maybe Stückelberg's) description of antiparticles as particles propagating backward in time that is responsible for our ignorance of the necessity (and not just the possibility) of a time-symmetric description of causality in the quantum realm. Once it is understood that there is a requirement for causality to be described in a time-symmetric way (due to the existence of backward-in-time propagation), then what we are facing is no longer merely the problem of understanding how the future can influence the known past, but really how it can be that a unique future may itself be causally related to this unique past, when there is obviously more than one way it can be influenced by those past events. Indeed, the relative nature of the order in time of space-like separated events, which is implied by the special theory of relativity, means that what appears to be the future, for a certain observer, is actually the past, for a different observer in a different state of motion, and therefore if a unique past is causally related to the experienced present, then the future can only be similarly characterized. In a time-symmetric context, it is not just the teleological character of backward

causation that would need to be justified, but the equivalent teleological character of ordinary forward-in-time causality. If one insists that there is a problem with the possibility of causally influencing the known past, then one must at least also admit that this problem could not be distinguished from that which would arise from the fact that the past also influences the future, while a unique (even though unpredictable) future is associated with the known unique present.

What makes it look like the present state of the universe is not causally related to one particular future is simply the fact that all future states compatible with current boundary conditions are, in principle, allowed to be the outcome of random evolution from a given present state (any outcome is possible for future measurements), while not all past states compatible with current conditions are allowed as ‘final’ states for backward-in-time evolution, in the context where the constraint of global entanglement discussed in section 4.9 exerts a limit on entropy growth in the past. Thus, from a practical viewpoint, regardless of what happens, certain constraints are always present for evolution toward the past and it is the fact that those constraints are precisely such as to prevent one present event from having recognizable consequences on multiple past events that makes it seem like current conditions have no consequences on past evolution. In fact, this is precisely the nature of the difference between what we usually call causality and which relates to the thermodynamic asymmetry between causes and effects and the kind of causality that is involved in a time-symmetric context. But once it is realized that the past and the future are not distinct, from a more general perspective, then it clearly follows that if we are willing to accept that the future can be influenced by what happens in the present, as confirmed by our direct experience of reality, then it is also necessary to recognize that the present can itself affect that past, in a certain way, so that imposing final conditions is no less appropriate than imposing initial conditions, as long as these conditions are not those responsible for the observed thermodynamic time asymmetry itself.

One commonly encountered misconception regarding backward causation in general, as well as in a quantum context, is that if the future is allowed to causally affect the past, then we can no longer be confident that the past is what it seems to be, because it can be altered by future events. This is usually provided as an argument against backward-in-time causation because, as everyone knows, the past is unique and unalterable and therefore any approach that would allow the future to ‘change’ the past is certainly

based on incorrect assumptions. But what is incorrect with this apparently logically unassailable conclusion is the idea that an alteration of the past would involve changing an observable fact from the past which we already know has occurred, just like we are allowed by our apparent free-will to alter the course of future history. In fact, that is just an inappropriate understanding of the meaning of backward causation, because if an event in the future changes the outcome of an observation in the past, this change has already taken place at the moment in the past at which it was observed and the fact is not changed ‘again’, from an alternative counterfactual, at the moment when the effect of this future cause reaches the event in the past which it contributes to determine. In other words, it is not possible to change history using the kind of time-reversed causal chain that is allowed by fundamental theories. History is the outcome of all effects, from both the past and the future, and is experienced only once as such a globally consistent whole. No known fact is altered or changed by effects propagating from the future, as any change that would be produced by future causes would need to have already taken effect at the time at which the fact first occurred. The ‘effects’ that the future may exert on the past would always be made conspicuous merely through the influence they would have on correlation probabilities established after the fact (when that future itself becomes a known past).

Now, it must be clear that the condition imposed on special-relativistic transformations that they should preserve the direction of all time-like intervals (the causal ordering postulate) is not incompatible with the conclusion that backward-in-time causation must be allowed at a certain level, because all that is required by relativity is that, if a causal chain operates from the past toward the future (as we usually assume to always be the case), then the same causal chain cannot be found to operate from the future toward the past as a result of such a transformation. But that does not mean that a distinct causal chain cannot operate from the future onto the past, merely that, if such a backward-propagating causal chain exists, it too cannot be turned into a causal chain propagating in the opposite direction of time, which in this case would be toward the future. In any case, the fact that the causal time of relativity theory is not unidirectional does not itself constitute a problem, because, even in such a context, unidirectionality is allowed to emerge from the global entanglement constraint which imposes a condition of low gravitational entropy at the Big Bang, as I have explained in section 4.9. What explains time asymmetry is not a fundamental property of unidirec-

tionality applying to causal chains, but the particular boundary conditions which apply at a certain time in the history of our universe. What makes a flow of information from the future toward the past impossible is not a limitation that would be arbitrarily imposed on the direction in which causality operates, but a distinct constraint that limits the growth of entropy in the direction of time toward the initial Big Bang state. Thus, as long as an effect does not propagate faster than the relativistic speed limit, it cannot give rise to a violation of the classical, unidirectional principle of causality, whether the effect propagates forward or backward in time.

It is not true that the scientific method excludes the possibility that ‘final causes’ of any kind might exist, despite the fact that time-symmetric causality appears to be allowed by relativity, because what can be scientifically demonstrated is merely that entropy does not increase in the past, not that there is no backward-in-time propagation of effects. What explains that we have become naturally suspicious, regarding the possibility that effects could propagate backward in time, is only the fact that from a classical perspective it never appeared necessary, or even possible to describe an object or a component of an object as propagating backward in time, while the time asymmetry that characterizes the history of macroscopic systems was always observed to operate from past to future (which made it look like a fundamental requirement). This prejudice remained in effect, even when it became clear that backward causation is a necessary assumption, in the context where one must recognize that certain particles do propagate backward in time (even though, from a unidirectional-time viewpoint, they are observed as oppositely charged particles propagating forward in time which are involved in the same entropy increasing processes as their ordinary matter counterparts). The teleological problem of time, which is often believed to arise in the context where present conditions are allowed to influence the known past, or even when one recognizes that a unique future is associated with the known present, is not a true problem, but merely follows from the psychological expectation of unidirectional causality that we inherited from our thermodynamically constrained experience of reality and which does not reflect any fundamental limitation on the propagation of effects. There is no other explanation for the widespread belief that causality must always operate forward in time.

What must be understood, then, is that it is not merely the order in which time flows that varies for an antiparticle, but really the fundamental direction in which effects propagate in time. When a particle propagates backward

in time, the direction in which it may come to influence other particles is actually reversed and it is merely the thermodynamic arrow of time and the direction in which classical causality operates that remain unchanged. If there is any meaning to be associated with a concept of causality, from a fundamental viewpoint, then a reversal of the order in which causality operates must be allowed to occur and future events must be allowed to exert an influence on past events. There cannot be a distinction between causal order and the direction of time in which an elementary particle propagates an effect, even though this direction is a relatively defined property and is only significant in relation to the direction of time in which another such particle is propagating a similar effect. This requirement may perhaps appear doubtful, in the context where it seems that the concept of particle trajectory may no longer be relevant in a quantum mechanical context, where multiple distinct histories would be required to take place all at once in order to account for the statistics of quantum processes. But I will explain, later in this chapter, why this apparent lack of uniqueness of particle trajectories is not an obstacle to a proper understanding of causality as actually depending on the direction of propagation in time of elementary particles. Causal order may be a locally-variable property, but it is not arbitrary, even when backward-in-time causation is allowed to take place.

5.4 Closed causal chains and time travel

As is already apparent from the viewpoint of a semi-classical description, the time-symmetric nature of causality does not merely imply that there is no absolute distinction between causes and effects, it also means that a certain event can, all at once, influence another event and be influenced by the very same event. In other words, not only is there no absolute difference between causes and effects, but the cause of a certain event can also be an effect of the same event, although this circularity can only be appropriately described in the context of a purely quantum-mechanical model of reality such as the one I will propose in a latter portion of this chapter. It must be clear, though, that the possibility that such closed causal chains may occur does not constitute a valid motive to reject the whole concept of time-symmetric causality and backward-in-time causation, because, as I will explain, it is possible to provide a consistent description of such phenomena without encountering logical contradictions.

Reichenbach's insistence [21] (p. 135) that one must be able to differentiate causes and effects independently from their temporal order, if we are to avoid the occurrence of closed causal chains, is not totally inappropriate, however, because, as I mentioned in the previous section, the direction in which causal chains operate is determined locally by the direction of propagation in time of the particles involved and therefore it is not fixed merely by the global time order of causally related events. But it is precisely the fact that one can distinguish causes from effects in such a way that allows backward-in-time causation to occur, while this is what makes closed causal chains unavoidable, thereby requiring such phenomena to be properly described and interpreted at the most fundamental level. Yet, even at the level where closed causal chains may occur, it is certainly necessary to require that no inconsistencies can arise, which would involve an incompatibility between some known present and some known past. What I will eventually explain is that there is actually a requirement for histories not to be self-contradictory and this condition can be satisfied, not merely despite the fact that causality also operates backward in time, but as a very consequence of the reality of backward causation.

In any case, it is certainly incorrect to argue that there is empirical evidence to the effect that closed causal chains are forbidden, because, in the course of elementary particle interactions, particle-antiparticle loops are often encountered that constitute just such a phenomenon, which can be adequately described, even from a semi-classical perspective. Once again, I believe that the problem here does not have to do with the possibility that closed causal chains themselves may occur, but rather concerns the hypothesis that classical, unidirectional causality could operate in both the future and the past directions of time along a closed causal chain. I will soon return to this question, but what should be clear already is that it is, in effect, only at the level where unidirectional causality operates that the order of events in time should be absolutely distinguishable and that no closed causal chains would be likely to arise. But, as I have mentioned already, this is a distinct issue, because at the fundamental level, where time-symmetric causality operates, thermodynamic time asymmetry is ineffective and any restriction that would be imposed by the existence of the thermodynamic arrow of time would be irrelevant.

One significant outcome of the existence of closed causal chains is that it is not always possible to establish the time order of events in an absolute way, because one event occurring along such a causal chain can be considered

to occur both before and after another event occurring along the same chain, even if the events are uniquely ordered from the macroscopic viewpoint of thermodynamic time. The topological order of time is always clear locally, along a particle world-line, but globally (even on a small scale) it must be determined in a purely relational way (as dependent on an arbitrarily-chosen reference point along a given circular trajectory), like any physically significant property.

It is important to realize that the existence of closed causal chains does not introduce additional unpredictability above that which is already assumed to characterize quantum evolution, because the current framework already involves some backward-in-time propagation (I will further explain, in section 5.8, what motivates the idea that backward causation is involved in determining correlation probabilities in quantum mechanics). But even in a deterministic context, the fact that the present cause of a certain future event could itself be affected by this very same event would imply that it is not possible for the future to be determined by the past alone, because the future itself would be involved in determining the past that determines this very future. Thus, it seems that even in the context of a hidden-variables model, backward-in-time causation would imply that reality must remain fundamentally random and not merely unpredictable (as when we have insufficient knowledge concerning the exact present state of a system, from a conventional viewpoint). Of course, the simple fact that the cause of a future event can be located not in its past, but in its own future, also implies that even when a complete knowledge of the present quantum state of a system and its environment is available, it is not possible to identify all the causes which exert an influence on its future evolution, which means that a certain measure of unpredictability is unavoidable that would not be present from a conventional viewpoint.

It is usually recognized that the problem which would be raised by what might be called a time travel experience has to do with the fact that such a phenomenon may allow the kind of closed causal chain in which the classical, *unidirectional* principle of causality would be violated. More specifically, given the assumption that we are free to decide how we influence the future, in the context where our evolution is taking place irreversibly toward what would normally be the unknown future, it may appear that a time-traveler, arriving from the future, would be able to alter the course of a known history in the same way a normally evolving person is allowed to influence the un-

known future. The problem here does not only have to do with the fact that we don't know why such unidirectional-causality-violating evolution is never observed, it also concerns the fact that if we are, in effect, free to influence the course of events taking place along the direction in which our thought processes are functioning, then it would appear that by 'traveling' back in time we might be able to alter a known future and to modify the course of events in a way that would be incompatible with the very possibility that the experience itself might have occurred, thereby giving rise to a time travel paradox.

Although time travel has never been observed to occur and therefore remains a purely hypothetical problem for physics, the standard answer to the questions it raises is often believed to be David Deutsch's proposal [40], based on the many-worlds interpretation of quantum mechanics. What Deutsch suggested, basically, is that every time a paradox would be expected to occur that would involve a particle arriving from the future and altering the past conditions that gave rise to the future state that allowed the process to happen, the universe would 'split' into alternative branches, where both the initial history (in which the backward-propagating particle did not change the future) and its modified version (were the particle does effect a change that would prevent it from having effected this very same change) do occur, but cease to interfere quantum mechanically with one another. Thus, it is proposed that there is an alternative future for every possibility that might be produced as a result of the influence exerted by a future event on a past event through backward causation.

I'm not sure what most people make of this description, but the problem I have with it is that I just can't figure out how it actually makes things any better. If we say that a particle arrives from the future and changes the past, then this past must be assumed to have already taken the 'effect' into account and must be such that it allows the said future to occur, as I previously explained. So, how could this future be made different by such an altered past? Clearly the problem with the hypothetical problem of a future 'cause' influencing a past 'effect' only occurs when we assume that there can actually arise inconsistencies or contradictions in the observed historical description of events. But when it is assumed that a particle can arrive from the future and change the past to which it was causally related, it is not possible to say that the future is merely altered from what it 'originally' was by the presence of the particle, because the particle itself could not even have arrived from the future in such an altered version of history. How could one

possibly argue that a new future is written in an alternative branch of the universe's history, as a result of the arrival of a particle from the future, if the backward-propagating influence of that particle did not even occur, in this alternative branch of history?

What the many-worlds approach purports to show is that inconsistencies and contradictions can actually arise in our historical description of facts, but that this is acceptable, because the future always adapts to the inconsistencies it itself generates. But this is just non-sense, because if a future is such that it influences a given past, then this past must be such that it necessarily gives rise to this unique future and this is not made to 'happen' by some hypothetical splitting process taking place at a given arbitrarily-chosen moment, it is just how things actually are all along, in both the past and the future. Also, if we are to allow for the existence of other universes, then by definition those universes should be causally independent from one another and things happening in one universe should not be allowed to influence what is taking place in another universe. I believe that what is missing from our current understanding is an acknowledgment of the fact that a universe, by necessity, actually consists of a unique ensemble of events causally related to one another and to nothing else (as a consequence of the requirement of relational definition of physical attributes which was discussed in the preceding chapters of this report). From such a viewpoint, if an event in the past is influenced by the presence of an event in the future then this past event cannot be causally related to a different future, but only to the future that actually influenced it. Thus, it becomes a fundamental requirement for the universe to form a consistent whole, free of internal contradictions.

Of course, we never experience time travel, so this issue only has to do with elementary particles propagating backward in time and in this realm quantum field theory already does a very good job of consistently describing physical reality and predicting facts. In this particular sense, Deutsch's proposal is a solution to a problem that does not exist and this becomes especially obvious in the context where, as I will later explain, the many-worlds interpretation of quantum theory is not required to make sense of the quantum measurement process and can even be understood to have consistency problems of its own (which does not mean that the multiverse concept, as a distinct hypothesis, cannot be considered valid and fully justified). I believe that the strange and convoluted reality that emerges from such a description merely illustrates the kind of complications we would run into if we adopted an interpretation of quantum theory involving such multiple splitting reali-

ties, present all at once in the same universe. If the many-worlds approach cannot even be made to work in a quantum-mechanical context, what motive do we have to invoke it in order to explain problems occurring at a classical level? Consistency requires that if a process is allowed to exert an influence backward in time, then this process must, in effect, evolve toward the very same past that gave rise to the future which is causally influencing this past, because the process must at all times remain causally related to the same external reality, otherwise nothing at all could be assumed to be causally related to anything else. But, then why is it, in effect, that we never experience time travel, if backward-in-time causation must be allowed to happen? Does this prohibition have to do with the fact that if it did not apply, then real contradictions might be made to occur whenever such closed causal chains would form?

To answer those questions, I must first point out that what would really differentiate time travel experiences from the backward-in-time propagation of elementary particles that is routinely observed in laboratories is the fact that, with time travel, a macroscopic and thermodynamically constrained system, such as a living human being, would need to evolve not just in the past direction of time, but with its thermodynamic arrow of time reversed and pointing toward the past instead of the future. From the viewpoint of an observer not part of the process, this evolution would be seen as a local violation of the second law of thermodynamics, or the principle that entropy never decreases in the future, because, indeed, if the time traveler really travels back in time, then, as he does, he would not just remember what happened at his past *destination*, but also what happened in the future (which to him would appear to also be a past), thereby allowing information to flow from the future toward the actual past. But this means that, from the viewpoint of a normal observer, the processes of memory formation and all the other irreversible processes usually involved in allowing a person to experience time as a unidirectional phenomenon would all appear to function backward for the time traveler, as the process is taking place, even if the time traveler is not composed of particles (like ordinary antiparticles) usually considered to be propagating backward in time.

At this point, it is necessary to recall the discussion from section 4.9 concerning the origin of thermodynamic time asymmetry in a universe like ours. There, I explained that it is the inescapable nature of the constraint of global entanglement (which must be imposed in order to allow relationships of causality to be established between all elementary particles in the universe)

that explains the parallel nature of thermodynamic time asymmetry (there does not coexist opposite thermodynamic arrows of time in different regions of the same universe), even for temporarily isolated branch systems. Thus, in a universe in which negative-energy matter was originally present with a density similar to that of positive-energy matter (for reasons I explained in section 4.5), gravitational entropy (and therefore all entropy) must be continuously decreasing as we approach the instant in the past (which corresponds with the very first instant of the Big Bang) where the global entanglement of all elementary particles is effected, because otherwise certain particles would not be allowed to be in contact with other particles present in the universe at the Planck time. As a result, all macroscopic systems must evolve with decreasing entropy in the past direction of time, because all matter particles, without any possible exception, must become entangled with the rest of the matter in the universe if they actually constitute elements of that universe.

Of course, the point here is that if time travel is never experienced or observed, it is not because backward-in-time causation is impossible at a fundamental level, but merely because entropy must be continuously decreasing in the past and only in the past (because no global entanglement constraint applies to the future), which means that the conditions necessary for the thermodynamic arrow of time to be experienced backward are not merely unlikely, they are actually forbidden, for all practical purpose. Unidirectional causality only operates from the past toward the future, because it would take a very significant fluctuation for entropy to temporarily decrease in the future, from a present state of non-maximum entropy, but given that this would be required for time travel to occur, then it is possible to understand why we never experience time backward. Indeed, classical, unidirectional causality is reflected in the fact that it would take only a little change in the past to allow a present event not to have occurred, while, in general, enormous changes would be required to take place in the future for some present event not to have occurred. This asymmetry is precisely what is enforced by the global entanglement constraint, when it is assumed that causal relationships must exist between all elementary particles in the universe. For time travel to be possible, this thermodynamic time asymmetry would need to be reversed locally for the whole duration of the process and the unlikelihood of such an evolution is responsible for the fact that time travel is virtually impossible, at least as a controlled phenomenon.

Therefore, it is not possible, in practice, to be involved in a closed causal chain while remembering what occurred at a later time (no information can

be transferred from the future toward the past), even if this restriction does not affect the possibility for microscopic systems to be involved in such causal chains, as long as global consistency is preserved (I will explain in section 5.8 how this condition is enforced at the fundamental level). This means that so-called ‘knowledge’ paradoxes are also unlikely to occur. It was suggested, in effect, that, if time travel was possible, there could arise situations where some valuable piece of information (say a beautiful treatise about the physics of time directionality) would be brought from the future that would not have existed before it arrived from that future, but which would nevertheless become available as a result of the process, so that it can later be brought back to the present, thereby raising questions as to its origin. But given that what would be required for such a paradox to occur is a sustained local increase of entropy toward the past (or equivalently a sustained local decrease of entropy toward the future), then it follows that the creation of information out of nothing in such a way would be as unlikely as the possibility that it materializes out of chaos by pure chance alone, which again is not fundamentally impossible, but merely ridiculously unlikely. Thus, from my viewpoint, despite the fact that backward causation is allowed to occur, it is not possible for information to be created out of nothing, which certainly agrees with what I have written concerning the conservation of information in chapters 3 and 4.

It is, therefore, possible to understand that, even in the absence of closed causal chains, what prevents violations of the classical, unidirectional principle of causality is actually the global entanglement constraint that restricts the growth of entropy in the past direction of time and this is clearly a constraint of irreversibility that is not imposed at a fundamental level, but that emerges from the particular boundary conditions which apply to the initial Big Bang state. The frequently encountered remark, to the effect that objects can move in any direction of space, but not in any direction of time (at least when they are restricted to not move faster than light in a vacuum), is only true in the sense that it is not possible to reverse a macroscopic system’s thermodynamic arrow of time; it does not mean that an object cannot propagate backward in time under appropriate conditions. If it was not for the constraint that is responsible for the diminution of entropy in the past, all evolution would be symmetric with respect to the direction of time, at all levels, and there would be no way for information to flow from either the past or the future, as all systems would remain in a state of thermal equilibrium (if that was possible) and no record making process would ever be allowed

to take place. Only under such conditions would it be possible to directly appreciate the fact that the future is not fundamentally different from the past.

What is important to understand is that, not only would entropy be observed to decrease in the future during a hypothetical time travel phenomenon (which would require it to increase in the past), but the process would remain observable all along as an entropy diminishing process taking place forward in time, even after a hypothetical time travel paradox would have been produced. Indeed, an observer which would be evolving backward in time, from a thermodynamic viewpoint, would still be causally influenced by the events taking place at the moment of unidirectional time which would appear to her as the present, so that she would observe the same sequence of events about which she already had knowledge, only now those events would take place in the reverse order. But she would also be allowed to causally influence those events at all times and those events would immediately influence her future, even before the process stops in the remote past. However, if a time traveler does remember the future as she evolves backward in time, then this future can only be the one which she already witnessed, even though it may seem that she would be free to alter this particular future, because if the process is allowed to occur it is because the future is such that it allows the process to actually arise and this means that *at all times* the present itself is necessarily such that it allows that particular future (about which information would be available) to happen. If that was not the case then the time traveler would not have advanced knowledge of the original future she wants to alter, but would rather remember the alternate future she intends to create, which means that she would not be better off than an ordinary observer at telling what that future is that she would try to alter and therefore we have no reason to believe that she would be allowed to *voluntarily* effect such a change.

What makes the paradoxes themselves impossible, however, is not the fact that they would require entropy to grow in the past, but the very same constraints that would forbid the occurrence of a contradiction from a fundamental viewpoint, as when elementary particles are propagating backward in time, without being involved in anti-thermodynamic evolution. In this particular sense, it is true that the problem of time travel can be fully resolved only in a quantum-mechanical context, but as I previously indicated (and for reasons that will be discussed only later) this does not mean that one must invoke the hypothetical splitting branches of a many-worlds interpreta-

tion of quantum theory. In any case, it is now possible to appreciate that what makes time travel, itself, impossible is not the fact that it may allow forbidden contradictions to occur, but really the improbability of observing processes for which entropy decreases in the future.

But, even if one was allowed to travel back in time, as a result of a phenomenal anti-thermodynamic fluctuation, one would not be allowed to alter one's own future, even if that is unexpected from our everyday viewpoint. Even under such conditions, there would necessarily occur events that would enforce global consistency and this would happen despite the fact that, under normal conditions, we seem to be free to modify the future at will. It is simply the fact that we are used to experience the future as unknowable in advance that explains that it appears doubtful that we would *not* be able to alter the course of reality². We usually have no factual knowledge about the future and this is why we never run into the possibility of being prevented from making a decision that we would expect to alter a known fact about the future. The requirement that reality be globally consistent appears to have unexpected consequences merely because we are not used to experience a reality in which we would have available information about what has not yet occurred. We are accustomed to observe that present actions exert an influence on the probability that such or such a future occurs, but this is only a reflection of the fact that there exist correlations between the past and the future, which are the result of both forward- and backward-in-time propagated effects and it does not mean that it is impossible for a unique future to be causally related to the unique past. It is merely the fact that all possibilities are usually allowed for the unknown future, while only a subset of them is allowed for the known past (due to the constraint responsible for the diminution of gravitational entropy) that justifies the impression we all share of being able to exert a certain control over the future which does not apply for the past, because it is precisely those limitations which imply that we can obtain knowledge about the past (which, therefore, appears unalterable), but not about the future (which, therefore, appears modifiable).

From a fundamental viewpoint, the future is not different from the past (even if it cannot be determined in advance) and we do know that the past cannot be changed from what it already is. If we never remember the fu-

²In order to understand how global consistency can be obeyed, even under such circumstances, it may help to notice that if knowledge about some future was to become available to a given observer, a prediction of her actions would have to take into account the fact that the prediction itself can influence the outcome.

ture and if we are never confronted with the limitations to free-will which exist as a result of the global consistency requirement, it is simply due to the fact that information does not usually flow from the future toward the past. This is probably the most important lesson that can be learned from the study of hypothetical time travel experiences, in the context of time-symmetric causality: we are causally related to one unique past. But this is also true for the future. We live in an unpredictable universe and while it is certainly true that what we choose to do now has an effect on what will happen tomorrow, based on the most rational explanation of both classical and quantum-mechanical phenomena, it is necessary to recognize that we are causally related to only one such future and even if we were to obtain, in advance, knowledge about what this unique future actually is, events would have to unfold in such a way that the consistency of history would remain inviolable.

5.5 Advanced waves and time asymmetry

Since Maxwell introduced his electromagnetic wave equations, more than a hundred fifty years ago, it has been known that there exist both retarded and advanced solutions to those equations (this is equivalent to say that Maxwell's equations do not distinguish the future from the past). The retarded solutions describe the propagation of positive-energy electromagnetic waves leaving a point source and spreading into a growing volume of space as time passes. The usually rejected advanced solutions, on the other hand, would describe the propagation of electromagnetic waves of opposite energy sign leaving a point source and spreading into a growing volume of space in the past direction of time. This is usually described as the hypothetical phenomenon of a spherical and concentric positive-energy light wave converging on a point source in the future direction of time³. From this equivalent viewpoint, it is obvious that the advanced solutions represent a kind of process that cannot occur, because, from the unidirectional-time viewpoint, one never observes light waves, or indeed any kind of waves, converging on a 'source' just to be

³The positive value of the energy of this converging wave, which allows the 'source' to gain energy as a result of the absorption process, arises from the fact that, as I explained in chapter 2, a negative-energy photon propagating backward in time is always observed as a positive-energy photon propagating forward in time, while a negative-energy photon propagating forward in time would not even be allowed to interact with ordinary matter.

absorbed by this source.

But while this observation reassures our commonsense expectations, the fact that the phenomenon described here never occurs, while there is no *a priori* reason why it couldn't happen, still constitutes a profound mystery from a theoretical perspective. It is usually recognized, in effect, that if a valid theory describes a certain phenomenon and there is no good motive to assume that this phenomenon should be forbidden, then its occurrence is compulsory. It is not enough to argue that what prevents the hypothetical phenomenon of a radio wave produced by multiple sources in the environment converging in perfect spherical symmetry and with perfectly correlated phases onto a transmitter, where it would be absorbed, is the unlikeliness of the phenomenon, because, as I emphasized in chapter 4, this is precisely what we 'observe' to occur in the past direction of time and this evolution is clearly not the outcome of the singular nature of present conditions. Given arbitrary initial conditions, what we should expect to observe are waves that would be diverging in the past, just like they do in the future, because this is in fact the most likely evolution when only the present conditions are fixed, even if, from the unidirectional-time viewpoint, such a process would appear unlikely.

If it is considered natural for electromagnetic waves to spread outward in the future, despite the fact that this means that they converge on their source in the past, then it should also be expected that electromagnetic waves would spread outward in the past, even if that means they would converge on their source in the future. Therefore, what remains unexplained is the asymmetry of the situation in which waves do not spread outward in the past, while they do so in the future of some arbitrarily-chosen initial state. The problem discussed here is all the more significant, given that it is not restricted to Maxwell's theory. Indeed, there exist advanced solutions to all relativistically invariant wave equations, including the equations that describes the propagation of electrons in quantum field theory.

Once again, this is a problem that Feynman visited, although it appears that he failed to resolve the issue. What he and John Wheeler proposed was a theory [41] that would have allowed advanced electromagnetic waves to be produced on an equal basis with retarded waves, just to be canceled out through destructive interference, as a consequence of the difference in opacity that seems to characterize the far past and the far future of our universe. According to this model, retarded and advanced electromagnetic waves are always produced together in equal proportions and propagate in the future

and the past respectively. But when the retarded wave is absorbed in the future, the absorbing process itself triggers the emission of an additional retarded wave of identical amplitude which is completely out of phase with the original retarded wave, thereby erasing all traces of this additional retarded wave. At the same time, the absorber also produces an advanced wave and if certain conditions are met, this advanced wave only serves to strengthen the retarded wave produced by the source through constructive interference, while it also conspires to cancel out the advanced wave originally emitted by the same source through destructive interference, which may allow to explain the fact that it is not observed. The problem is that this theory requires that there is more absorption in the future than in the past, while this possibility would appear unlikely in the context where the universe is expanding in the future direction of time and the matter density is maximum at the Big Bang.

Other theories, based on similar assumptions (see for example Refs. [42] [43] [44]), and which tried to overcome the problems encountered by Feynman and Wheeler through various alternative hypotheses (for example by assuming that the Big Bang acts as a reflector of all advanced radiation), have apparently also failed to produce a satisfactory solution to the problem of advanced waves. It seems that, whenever it is not independently assumed that, for some unknown reason, a fundamental asymmetry exists in the interaction of matter with radiation that would differentiate the far past from the far future, the desired outcome is never obtained. In other words, the only way to reproduce the observed time asymmetry that characterizes wavelike processes in our universe, using such a model, is by postulating that some asymmetry exists, which is responsible for reducing or increasing the amount of interference that takes place either in the past or in the future. But given that no convincing explanation exists that would justify this assumption, then it is apparent that it merely amounts to assume the very outcome we would like to explain. From the difficulties encountered with this kind of approach, it has become pretty clear that it is not possible to explain the absence of advanced waves as being a mere consequence of hypothetical interference effects.

I was only able to understand what explains the absence of advanced waves when I began considering the quantum aspect of this hypothetical phenomenon. Indeed, I already knew that backward-in-time propagation was possible for elementary particles and therefore it seemed to me that what was not allowed was not really backward propagation itself, but merely the spreading of a backward-propagating wave into an increasingly larger

region of space. I also knew that there was a requirement, imposed by the constraint of global entanglement which I had recently uncovered, that backward-in-time evolution be such that it gives rise to a continuous decrease of gravitational entropy in the past. But, as elegantly explained by Olivier Costa de Beauregard [45], there is a certain equivalence between entropy increase and wave retardation, which is implied by Planck's definition of entropy and which arises from the quantized nature of electromagnetic radiation. Thus, in a quantum-mechanical context, entropy necessarily rises when an electromagnetic wave spreads into a larger volume of space, because, at any given time, the photons associated with an expanding wave front can be detected anywhere on its growing surface. In fact, given that, from the viewpoint of relativistic quantum field theory, any wavelike phenomenon is associated with the propagation of some elementary particle, it follows that entropy increase is always associated with wave retardation, while the observation of advanced waves would always imply that a decrease of entropy has taken place in the future direction of time. From my bidirectional-time viewpoint, this is equivalent to say that entropy would need to increase in the past for an advanced wave to spread as it propagates backward. But this is precisely what is forbidden by the constraint I have previously identified as being responsible for thermodynamic time asymmetry.

What is also unexpected, from a thermodynamic viewpoint, is the fact that, from the unidirectional-time viewpoint, the existence of advanced waves would seem to allow work to be generated out of nothing, when radiative energy would converge on a 'source'. But the existence of advanced waves would also make possible the transmission of information from the present or the future toward the past. It is natural to expect, therefore, that this kind of process should be prevented from occurring by the same condition that explains thermodynamic time asymmetry. It must be clear, however, that simply invoking the classical (unidirectional) principle of causality does not allow to solve the problem of the absence of advanced waves, because, in the above discussed context, saying that there always exists a unique preferred direction in time for the propagation of effects merely amounts to restate the problem of advanced waves (which is also known as the problem of the electromagnetic arrow of time) without explaining why such a restriction is indeed observed to apply. In fact, the hypothetical phenomenon of time travel I have described in the preceding section would be one particular instance of backward-in-time communication, of the kind that would be allowed by the existence of advanced electromagnetic waves, and therefore a solution to the

problem of advanced waves would definitely rule out time travel.

Now, I mentioned in section 5.3 that the causal structure of spacetime is not incompatible with the concept of backward-in-time causation, given that with every event is associated both a future and a past light cone, which reflect the existence of a speed limit imposed on the propagation of causal signals in either the future or the past. But it should also be clear by now that there is a difference between the kind of backward-in-time causation that may occur as a consequence of the propagation of an elementary particle backward in time and the kind of causality we experience in a purely classical context and which is known to operate only forward in time. Thus, while it is not observationally forbidden for an electron to propagate backward in time, an explanation of cosmological time asymmetry based on the global entanglement constraint would not allow this propagation to occur in such a way that the surface over which the presence of the electron could be detected at an *earlier* time would be growing continuously along with the two-dimensional boundary of the past light cone. But this is precisely the kind of evolution that an advanced wave would describe from a quantum-mechanical viewpoint and therefore what explains that advanced waves are absent is the constraint of global entanglement I have identified in section 4.9, which enforces a continuous decrease of entropy in the past, as a consequence of the requirement that there exist causal relationships between all the elementary particles which are present in the expanding universe.

Our failure to observe advanced waves must not, therefore, be interpreted as an indication that backward-in-time propagation, or backward-in-time causation are forbidden, but rather as evidence that only a small subset of potentially available states is available as ‘final’ conditions for backward-propagating particles. Such particles are not only prevented from propagating faster than the speed of light in the past direction of time by the causal structure of spacetime and the existence of a past light cone, they are also prevented from propagating in all possible directions of space in ways that would allow entropy to grow in the past. This means that the statistical predictions, obtained using quantum theory, for the evolution of a large number of identically prepared physical systems are not valid in the past direction of time and this is what explains that electromagnetic waves, as particular instances of wave functions, are never observed in their advanced form.

In such a context, it becomes apparent that the only true virtue of the Feynman-Wheeler absorber theory (aside from the fact that it was one of the first models which actually took the problem of advanced waves seriously)

is that it sought to deduce the absence of advanced waves from boundary conditions imposed on the universe at large, instead of requiring that time-asymmetry be imposed at a fundamental level, which could only be satisfied by assuming that backward-in-time propagation is impossible. In any case, even if absorber theory had conveniently solved the problem of advanced waves, this solution would have remained problematic, because it would not have allowed to explain the origin of thermodynamic time asymmetry in a more general context (when quantum interferences are absent). From my viewpoint, the fact that there also exist advanced solutions to Dirac's relativistic equation for the electron allows to confirm the validity of the conclusion that the absence of advanced waves does not preclude backward-in-time propagation, because, while it is not possible to assess whether a given photon propagates forward or backward in time, in the case of electrons it is possible to differentiate a forward-in-time-propagating particle from a backward-in-time-propagating particle, given that, from a unidirectional-time perspective, such an electron is observed as a positron with its positive electric charge. Therefore, if we do observe positrons, it means that the irrelevance of advanced solutions cannot arise from the nonphysical nature of backward-in-time-propagating particles and must, in effect, be the outcome of the global entanglement constraint.

5.6 Early interpretations

To begin the portion of this chapter that deals with quantum aspects of reality more specifically, I would like to first describe what constitutes the distinctive characteristic of the revised interpretation of quantum theory I will propose. What I had already understood, even before I was able to solve the problem of advanced waves, is that the processes that constitute the essence of our experience of reality are all mirrored by similar processes which obey the same observable macroscopic conditions, but which take place in the opposite chronological order, in a portion of history that must be assumed *independent* from the viewpoint of local causality. The hypothesis that history does not occur only once, but must happen a second time in the reverse order may appear arbitrary and unnecessary, given that we know of only one history, but, as I will explain, this proposition is actually made unavoidable by some of the most fundamental principles of physics and also reflects the basic mathematical structure of quantum theory. Even though

I was not motivated only by the desire to produce a time-symmetric theory when I began developing this original approach, the final outcome does share a certain property of time symmetry with some early interpretations of quantum theory which are based on the hypothesis that there must be an equivalence between initial and final conditions.

Given that most of those early time-symmetric interpretations constitute more or less elaborate (and more or less inappropriate) quantum mechanical versions of the original absorber theory discussed in the preceding section, then one may say that absorber theory is their common ancestor. In this respect, it is apparent that those time-symmetric quantum theories also share some of the above discussed weaknesses of the original, classical theory. I believe that if this kind of approach is usually considered to have failed to produce a consistent interpretation of quantum theory, despite the many advantages it offers (which will be discussed below), this is due in part to the fact that absorber theory, itself, is considered a failure. As a result, many generations of physicists were inoculated against time-symmetric approaches in general, even though a few well-informed specialists have recognized that the requirement of time symmetry is essential to a consistent interpretation of quantum theory. But it is also clear that this is not the only reason why the early attempts at formulating a time-symmetric version of quantum theory did not arouse more interest, because, as I came to understand, they also contain hypotheses and constructs that make them inconsistent and inadequate as a representation of quantum reality.

One of the first interpretation of quantum theory that sought to accommodate the requirement of time symmetry was that proposed by John Cramer [46] as an outcome of his work on the problem of advanced waves. As such, it contains hypotheses which are very similar to those of the original absorber theory which I have identified as problematic. But its most important defect, in my opinion, is that, despite the fact that it is proposed as an alternative time-symmetric model, it actually involves a fundamental time asymmetry that is incompatible with this basic requirement. What Cramer proposed, basically, was that a kind of ‘handshake’ process takes place whenever a quantum particle is emitted by a source and then propagates a certain distance before being absorbed by a detector. We may consider, for example, the traditional double slit experiment in which a particle must go from source to detector by passing through the slits. It is known that an accurate estimation of the probability for such a process to occur must take into account the existence of interferences between the individual probability amplitudes

associated with each of the paths through which the particle is allowed to go, whenever both slits are open.

What Cramer's handshake process involves is the emission of a classical wave acting as an 'offer', which is assumed to be sent by the source forward in time and which is allowed to propagate without constraint (it is assumed to go through both slits all at once), followed by the production of another such wave that would constitute its 'confirmation' and which would be sent by the detector backward in time (toward the initial emission event), upon absorption of the offer wave. The most problematic aspect of this description, from my viewpoint, is the fact that the confirmation wave must follow an evolution that is restricted to be compatible with the macroscopic constraints which would have existed if the particle (not the offer wave) had been restricted to follow the unique classical path it is assumed to actually have taken as it propagated forward in time (the confirmation wave only comes back through one of the two open slits)⁴.

It is difficult to see how the advanced wave could be submitted to macroscopic constraints which differ from those that apply to the retarded wave, in the context where the observed macroscopic conditions of the experiment are fixed once and for all. But what is even more incomprehensible with this interpretation is that the evolution of the 'confirmation' wave is actually required to reflect the fact that the particle took a certain path (say the upper slit), while the evolution of the 'offer' wave would not be allowed to reflect the same fact (passage through both slits would initially be allowed). This is how time asymmetry is reintroduced in the model, as a means to allow a unique, *classically* well-defined history to correspond with the process, despite the fact that the statistics of this quantum process can only be explained by assuming that the particle is not restricted to follow a unique path. Of course, even if those problems did not exist, there would still be a difficulty associated with the fact that this approach requires the existence of both classical waves and classical particles (constrained to follow unique trajectories by those classical waves), while it is known that both concepts (which are shared by certain classical hidden-variables theories) are problematic in quantum field theory.

I believe that the source of the problems affecting Cramer's transactional

⁴In fact, Cramer assumes that this handshake is actually repeated several times, for any single quantum process, and is responsible for the transfer of energy and other conserved quantities which take place during the process, but we may ignore this problematic aspect of the handshake process if it simplifies the discussion.

interpretation of quantum theory is to be found in the fact that it assumes that the retarded and advanced waves are actually propagating in the same portion of history, because this is why it needs to be required that the quantum particle submitted to the constraint of those classical waves goes through only one slit, corresponding to this unique history, which, in turn, requires a certain fundamental temporal asymmetry to be introduced into the theory, in violation of the time-symmetric nature of its equations. Also, the fact that, as a particular instance of (quantum-mechanical) absorber theory, Cramer's framework appears to require genuine wave emission and absorption to take place in the course of all quantum processes, may be problematic, because there are situations where quantum measurements are performed without interaction. Those difficulties are more significant than the additional problem that would arise in the context where it is not obvious, from the viewpoint of Cramer's theory, when it is exactly that the handshake would be initiated, while the particle is propagating along its classical path.

Indeed, if the handshake was to be completed when the particle reaches one of the two open slits, then the process would always be that which we expect to occur when one is allowed to observe through which slit the particle goes and under such conditions, the particle would follow a quasiclassical trajectory (interferences would be absent), which is contrary to observation. Thus, there may be a difficulty associated with the apparent arbitrariness of the choice of which macroscopic conditions are necessary to trigger a handshake (do we have to wait for an observer to become aware of the outcome as John Von Neumann once proposed?). But this is in fact the same quantum measurement problem as may affect a more traditional interpretation and therefore we are allowed to assume that any solution to this problem that would be proposed in a more conventional context would also apply to the transactional interpretation. This is an important point, because this difficulty is sometimes proposed as an argument against all time-symmetric approaches to quantum theory, while, when it is properly understood, it no longer stands out as a problem that is specific to time-symmetric models. Of course, it would not be appropriate, either, to assume that Cramer's theory is equivalent to standard quantum theory, as its author suggested, because the validity of the predictions derived from ordinary quantum mechanics does not depend on the existence of advanced waves, while this hypothesis is essential in the context of the transactional interpretation. In fact, when the inadequacy of the boundary conditions that give rise to the destructive interference effects that would allow advanced waves to go unnoticed is rec-

ognized, the theory no longer even agrees with observation, which certainly makes it different from standard quantum mechanics.

What I'm suggesting that we retain from those alternative, semi-classical interpretations is the notion that the squaring of the wave function, which allows to obtain the probability of a process, is made necessary as a consequence of the fact that, somehow, two histories are involved in any quantum process. I believe that this is what explains that it is merely by multiplying the probability amplitudes associated with each of those paired processes that we can obtain (under appropriate conditions) the probability for the entire process to occur. Indeed, the squaring of the wave function (which is necessary to obtain the probability of a process in quantum mechanics) involves taking the complex conjugate of the probability amplitude associated with one history before multiplying it with the probability amplitude associated with another history and it is known that taking the complex conjugate is equivalent to reversing the direction of time for those equations that describe the changes taking place in the quantum state of a system.

Therefore, the core mathematical framework of quantum theory already contains, in embryonic form, the requirement that each process be described as a history that unfolds forward and then backward in time, for some mysterious reason. This otherwise puzzling requirement has been transformed by modern interpretations (such as the consistent-histories interpretation of quantum theory) into a condition, imposed (without any real justification) on certain pairs of minimally coarse-grained histories, that they provide the probability of occurrence of a 'consistent' history. But in the process, it seems that the most important aspect of this requirement, which is the fact that the two histories forming such pairs actually take place in opposite directions of time, was lost and with it, the important insight we should have learned from early time-symmetric interpretations of quantum theory.

At this point, it is important to mention that a more pragmatic approach to achieve symmetry with respect to the direction of time in quantum mechanics had already been proposed by Aharonov, Bergmann and Lebowitz [47] (see also Ref. [48] for a more recent review) long before Cramer introduced his transactional interpretation. Unlike the transactional interpretation, this formulation of quantum mechanics really is mathematically equivalent to the standard theory, but it does not seek to explain the time asymmetry of boundary conditions and merely suggests that two state vectors are required to describe the state of a quantum system. One state vector contains all the information obtained from past measurements (as in

the standard interpretation) and the other contains all the information that will be obtained, concerning the same system, in the future. Between measurements, those two state vectors follow a ‘unitary’ evolution toward the future and toward the past respectively⁵. What this means is that there is no longer a preference for the past over the future in determining the current state of a system (a system can be submitted to both pre- and post selection, although the post selection is only apparent after a future measurement has actually been performed). Of course, there is a natural reluctance to recognize that it might be possible for a state vector to be determined by what ‘happened’ in the future, instead of what happened in the past, but this is merely a consequence of the previously discussed prejudice toward a unidirectional conception of causality, which we inherited from our thermodynamically constrained experience of reality and it does not rest on any rationally formulated argument.

It must be clear that, despite the equivalence between the two-state-vector formalism and standard quantum theory, it has been shown that post selection, or the effect of a future measurement on the past state of a system, is not an optional feature of quantum theory, but arises even in the simplest and most conventional of quantum-mechanical experiments. Indeed, in certain interferometer experiments which bear enormous resemblance to the classical double slit experiment and which will be discussed in section 5.9, the choice of performing either a measurement that determines through which path a photon went on its way to the detector, or a measurement that reveals the existence of quantum interferences attributable to the presence of two possible paths, can be delayed to long after the particle has actually traveled most of the distance to the detector and it does, in effect, influence what the particle did back when it was just leaving the source. The reality of such post selection effects has therefore been experimentally confirmed and contrarily to what is sometimes suggested, it is not possible to assume that no post selection occurs as long as we reject a realist interpretation of quantum phenomena (because it is not possible to reject such an interpretation, as I will explain later). Thus, somehow, the path taken by a photon can be influenced by a measurement that takes place long after the actual process is

⁵I use the term ‘unitary’ to denote the *deterministic* evolution of the wave function or state vector that takes place in the absence of a change (usually performed through an act of measurement) in the observational constraints applied on a quantum system, because using the term ‘deterministic’ would be misleading in the context where I will argue that the evolution of the system itself always takes place randomly.

over⁶. Only a time-symmetric approach to quantum theory that recognizes the existence of a backward-evolving state allows to explain those facts while remaining within the confines of the principle of local causality.

Now, even though some of the originators of the two-state-vector formulation of quantum theory are hesitant to assume the reality of the backward-evolving state that enters the formalism, it is clearly possible to assume that we are indeed dealing with a distinct state that evolves somewhat independently from the forward-propagating state, but which is subjected to the same macroscopic experimental conditions. What I'm proposing is that in order to go beyond early time-symmetric models one must, in effect, recognize that a whole history unfolds backward in time, whose elements are not in causal contact with those of the history that unfolds forward in time. Indeed, I believe that, in order to accommodate the requirement of time symmetry, it is not enough to simply assume that semi-classical waves are propagating backward in time, in the same portion of history, because, as I have already explained, advanced waves are forbidden to exist by the constraint of global entanglement that gives rise to time asymmetry in our universe. The problem here usually is that, even though two kinds of Schrödinger equation appear to exist, which would allow to describe the propagation of the wave function in either the future or the past, only the equation that describes the evolution of the retarded portion of the wave function is retained, given that retarded waves are the only ones which are allowed to evolve without constraint, and this is why it is usually considered appropriate to take into account only the state vector that evolves forward in time in order to obtain the probability of a whole process, even if this process may actually involve a pair of state vectors evolving in opposite chronological orders.

But once it is understood that this limitation is not a fundamental property of the wave function itself, but arises as a consequence of the requirement of diminishing entropy imposed on all past evolution by the global entanglement constraint that applies to the initial Big Bang state, then the two-state-vector formalism becomes not only acceptable (as it does not require the existence of advanced waves), but actually essential to accommodate time symmetry in a quantum-mechanical context. In fact, given that the

⁶It should be clear that I'm not suggesting that post selection would allow information to flow from the future, or that it would allow one to change an observable fact from the past which has already been established. For reasons I have already mentioned, backward causation, as would occur in the context of a consistent, time-symmetric interpretation of quantum theory, is incompatible with both of those conclusions.

direction of time in which any process unfolds is a relatively defined property, the state vector which we may consider to be determined by future measurement conditions (the post-selected state vector) could also be considered, as a matter of convention, to be that which was determined by past conditions, while the state vector which would otherwise be assumed to be determined by past measurement conditions (the ordinary state vector) may be considered as that which actually evolves backward in time, ‘after’ having been determined by future conditions, as long as the other state vector is, in effect, assumed to be that which evolves forward in time. Therefore, we would not be better off by assuming that only past conditions can determine the evolution of the state vector, because this could also be understood to mean that only *future* conditions can determine the same evolution, which would be an even worse conclusion from a conventional perspective.

I may add that an explanation of thermodynamic time asymmetry of the kind I have proposed in section 4.9 does not only render plausible the hypothesis that every quantum process is complemented by a backward-evolving counterpart, but actually seems to require the existence of two histories evolving in opposite chronological order, because, otherwise, it would be difficult to explain what enforces the then unique, classical history, which we would be free to consider as evolving toward the past, to take place with continuously decreasing entropy. But once it is recognized that there necessarily exists at least one history that unfolds from the past toward the future (as one needs to assume in the context of a time-symmetric interpretation), then it becomes possible to explain the thermodynamic arrow of time as being the consequence of the initial condition of low gravitational entropy imposed on the initial Big Bang state by the global entanglement constraint, because the evolution of at least one state vector is then determined by past conditions (while the evolution of the other state vector is determined by future conditions, which allows the same requirement to be fulfilled for that portion of history unfolding past the initial Big Bang singularity, following a quantum bounce). In fact, this is a general requirement that would apply to all processes in the context where global consistency (not to be confused with global entanglement) is required, because, from a quantum-mechanical viewpoint, the consistency of past events with future events can only be fulfilled when those future events are also allowed to influence past events, as I will explain in section 5.12.

Once this is understood, it is easy to see how a relativistically invariant model based on the sum-over-histories approach can be formulated that

embodies the explicit time symmetry of the two-state-vector formalism by assuming that every quantum process involves both a conventional history (evolving without apparent constraint in the future direction of time) and a possibly distinct, time-reversed history evolving independently (from the viewpoint of local causality) toward a state of lower entropy in the past direction of time. This is an issue I will discuss more specifically in section 5.8, but before I can do that, I must first explain why it is that a model involving two unique, but partly unobservable histories, unfolding in opposite directions of time (instead of two wave functions propagating in opposite directions of time), is not merely possible, but actually constitutes an essential requirement of a fully consistent, realist interpretation of quantum theory, despite the fact that what is usually assumed to be required in order to obtain the appropriate statistics is that all possible paths are followed all at once, in one single portion of history, for any given process.

5.7 The constraint of scientific realism

It has often been argued that the counter-intuitive aspect of quantum theory is not a real problem and merely indicates that there is a limit to what one can intuitively understand. It would then be incorrect to assume that the fact that there appears to be something incomprehensible with the current interpretation of the theory is due to the inadequacy of this interpretation. I believe, however, that this argument is invalid. In order to see what is wrong with this long-standing viewpoint, let's first suppose that we humans are, in effect, too dumb to understand quantum theory. The argument would then be that only some artificial super-intelligence from the future (perhaps one that would run on a quantum computer) would eventually be able to overcome those limitations and to properly understand the significance of the empirically derived mathematical framework of quantum theory. Such a super-intelligence would, therefore, succeed at gaining a proper understanding of physical reality in a way that is simply impossible for us to achieve, due to the inherent limitations of our primitive intellect. But what does that mean in concrete terms?

When you carefully think about this question, it becomes obvious that the only thing that could happen is that this super-intelligence would then have developed a better interpretation of quantum theory, because if the current mathematical framework is, in effect, appropriate to describe physical

reality, then the only progress that *could* be achieved would have to arise at the level of interpretation. You do not have to be super-intelligent to understand that and yet this is precisely what we fail to take into account when we suggest that the problem we experience while trying to make sense of quantum theory merely reflects the fact that it is not possible for our brains to understand the theory. I believe that the lack of intelligibility of our current understanding of quantum theory is not a fantastic new property of the universe which we happen to have discovered. It is a failure that originates in the inappropriateness of the current interpretation and if this difficulty may be a consequence of the inadequacy of certain concepts we inherited from our human experience of the world, it is also a problem that can be solved using our human intellect, as long as we do recognize that there is indeed a problem and that it deserves our attention. But those who still doubt the importance of a proper interpretation of quantum theory should take notice of the fact that, without interpretation, it would not even be clear that the theory is about probabilities of measurement outcomes, as this is indeed an aspect that only came to be understood after the mathematics of the theory (regarding the Schrödinger formulation in particular) had already been developed.

Now, it must be clear that quantum theory *is*, in effect, counter-intuitive and that it cannot be reduced to a classical view of the world by using the freedom we may have to interpret experimental facts and the current mathematical framework of the theory. Physical reality cannot be such as it was conceived at the epoch of Isaac Newton. Classical waves (which are not mere manifestations of quantum interference) and classical particles (which would not be subjected to the constraints imposed by the uncertainty principle) are gone and they will never form part of a consistent theory about the fundamental structure of reality ever again. But that does not mean that everything else is possible. What is not allowed of a rational understanding of physical reality is inconsistency. The problem is that all known interpretations of quantum theory do contain inconsistencies. Thus, either they contradict themselves, or else they do not agree with certain facts concerning that portion of reality which can be directly observed. This is usually understood by well-informed authors who recognize that the best that we can do in the present context is to pick as our necessarily inaccurate standpoint the interpretation which may be the least problematic for the kind of problem we are working on.

What I have come to realize is that, while some new conceptual elements

(which have never been considered before) are necessary to formulate a fully consistent, but straightforward interpretation of quantum theory (which actually constitutes a more accurate physical theory), it is also necessary to reject many of the outlandish concepts that came to be associated with a quantum-mechanical description of reality. Thus, I believe that the concept of history, or the concept of reality itself, must be simplified to once again be allowed to agree with the most basic empirical evidence, concerning in particular the uniqueness of facts and the particle nature of physical reality (as a concept more elementary, but also more refined than its classical counterpart). The problem here is that it is often believed that the notion of an elementary particle propagating along a unique trajectory is incompatible with the ‘complexity’ which characterizes the quantum state of a system. But, as best understood by Richard Feynman, given the right formulation of quantum theory, not only is it unnecessary to reject the existence of elementary particles, or even to deny the relevance of the concept of trajectory, but it becomes imperative to recognize that those concepts actually form the substance of reality on the scale at which the most precise experimental data can be obtained.

I think that it is important to emphasize, therefore, that, even though common sense is not always a good guide for judging the validity of a physical theory, as the development of quantum mechanics itself illustrates, it would not be wise to conclude from this that more intuitive models are inappropriate and are necessarily ruled out by the apparent awkwardness of experimental facts, in the sense that our direct experience of reality would need to be considered irrelevant as a guide for elaborating a consistent interpretation of quantum theory. We must keep in mind that classical physics itself once involved counter-intuitive concepts which turned out to be inappropriate because of their awkwardness (like action at a distance), or which only became fully understandable in the context of a more intuitive formulation of quantum theory (like the principle of least action). Thus, I believe that, in the end, quantum reality will not be more difficult to visualize than classical reality, but will rather be more comprehensible, because it will be more consistent from a logical viewpoint. In any case, I believe that I’m justified in adopting a less counter-intuitive approach, given that the persistent problems which we are dealing with here have to do precisely with the apparent impossibility to provide a consistent, but also understandable representation of reality. However, instead of entering into a sterile debate about which of the ontological or the epistemological viewpoint constitutes a

better approach to interpret quantum theory⁷, I will concentrate on explaining what the elements of an empirically accurate approach actually are that allow to reach consistency with the least amount of arbitrary hypotheses (I believe one does not need any).

To begin this discussion, it would be appropriate to point out that the most radical of those deficient approaches which were once proposed in order to make sense of quantum theory is certainly that which is called quantum logic. It was suggested, in effect, that the logic that applies to physical reality may not be the ordinary Boolean logic with which we interpret ordinary facts, but some alternative logic, emerging from the apparently contradictory nature of certain conclusions made on the basis of a strict adherence to the rules which govern quantum reality. But while it is now recognized that such an approach would go too far as a tentative to adapt our mode of thinking to the reality of the quantum world, the fact that, at a certain epoch, quantum logic was considered to constitute a viable candidate for a solution to the problem of interpretation is quite indicative, I believe, of the extent to which we have departed from serving the true objective of science, which is to understand facts by adapting and generalizing our physical laws and concepts to fit new experimental facts, in order precisely to avoid having to change the rules of logic with which we analyze and understand reality.

The best example of such an adaptation is, of course, the shift to Riemannian spacetime that was brought about by relativity theory as a means to retain the validity of the concept of space in view of the equivalence of acceleration and gravitation. Indeed, if we were to reject Einstein's theory of gravitation, the only way we could retain the validity of the concept of physical space would be by altering the rules by which we formulate logical arguments, such as would be necessary to argue that despite all the evidence the Earth is flat. What the whole history of physics tells us is that it is always appropriate to use logical coherence as a means to constrain our representations of reality and as a guide to assess the validity of our assumptions,

⁷The debate concerning interpretation has always centered around the problem of deciding whether the wave function that allows to derive the quantum statistics of a process is a real 'entity' or whether it merely provides the sum of all knowledge about what a (real) system is doing, which I believe is pointless, because even if the wave function is not physical reality itself it does provide the most accurate description of the state of a quantum system at any particular time and this description is a real *aspect* of the system. The approach I will follow may actually be considered to allow a reconciliation of those two apparently incompatible viewpoints.

while the rules of logic themselves are rather like the rules of the game and can only be altered at the expense of invalidating most of everything else we have learned. But the mere fact that quantum logicians were never able to dispense themselves from the need to use ordinary logic in order to reason about their own alternative system is quite indicative of the failure of their approach. I think that this is a particularly good example of the difficulties which the currently favored interpretations of quantum theory are facing as they stretch the notion of consistency, while trying to adapt to some perceived requirement of the mathematical framework of the theory, by going so far as actually allowing for contradictory accounts of factual aspects of the world. I will return to this question later in this section.

Not so long ago, it was suggested that certain difficulties that emerged as a result of the development of quantum field theory may indicate that the concept of an elementary particle is no longer relevant to fundamental theoretical physics. One of those ‘difficulties’ would have to do with the fact that, due to quantum uncertainty, particles can no longer be considered to be localized in space, as would seem to be necessary for the particle concept itself to make sense. Actually, in a relativistic context, it seems that the very fact that a particle is localized may depend on the state of motion of the observer which is assessing this fact, given that a particle’s wavelength varies as a function of its relative velocity (actually its momentum). Another aspect of the quantum-mechanical description of reality which would appear to constitute a serious challenge for the particle concept is quantum entanglement and the demonstration that what one particle does may, under certain conditions, depend on what another particle is doing at the exact same time in a remote location (relative to a given reference system), thereby apparently implying that only the ensemble, consisting of the two particles taken together, has physical significance. Finally, an additional difficulty arises from the fact that, due to the fluctuating nature of the quantum vacuum, the very reality of a particle’s existence may be called into question, because, even in empty space, particles would appear to be present. This problem is particularly severe in the context of a semi-classical approach, in which the effects of acceleration and spacetime curvature on the quantum vacuum are taken into account and the presence of real (observable) particles becomes an observer-dependent property.

While I will not immediately address the issue of quantum entanglement and non-locality, the conclusion I have reached is that, despite the difficulties

mentioned here, the elementary particle concept is still viable in quantum field theory. In the remainder of this section I will provide arguments to the effect that a realist description of physical processes, based on the concept of particle trajectory is still desirable, even in the context where quantum interference, involving multiple distinct position states, must be assumed to constitute an essential aspect of reality. What emerges from this reflection is that it might be incorrect to suggest that particles cannot be localized in any way, because it may well be that particles in a pure momentum state do follow unique, but unobservable trajectories in a certain sense which is merely incompatible with the *classical* concept of trajectory. In such a context, the fact that the ‘wave packet’ which is sometimes associated with the position state of a particle can be more or less localized in space, depending on the state of motion of the observer who measures this position, would not mean that a particle can actually be more or less ‘real’, because such a variation would merely be a reflection of the dependence of macroscopic conditions (here those which constrain the non-classical trajectory of the particle) on the choice of a particular reference system. But a detailed description of the realist picture of quantum processes that allows to articulate those considerations will only be provided in section 5.8. In any case, I believe that the only real problem here is the general confusion that surrounds the question of deciding what it is exactly that remains acceptable about the particle concept in a quantum field theoretical context, because all attempts at completely disposing of this essential concept have failed to provide a sensible, alternative conception of the nature of physical reality at the most elementary level of description.

What I would like to immediately emphasize, though, is that, in light of the developments already introduced in chapter 2, it is possible to conclude that vacuum fluctuations, far from constituting a problem for the elementary particle concept, actually allow to provide a more consistent definition of what a matter particle really is. Indeed, what I previously explained is that positive-energy particles must be considered to arise from an absence of both negative energy and positive or negative charge in the fluctuating vacuum, that is to say, from an absence of virtual particles in the otherwise electrically neutral portion of zero-point vacuum fluctuations that contribute negatively to the maximum measure of vacuum energy density (while negative-energy particles arise from a similar absence of both positive vacuum energy and positive or negative charge). It therefore appears that the distinction between real particles and the virtual particles which are present in the vacuum is not

as significant as one might imagine, given that the presence of real particles is actually equivalent to an absence of virtual particles in the quantum vacuum. But it was also made very clear, in section 4.3 and then in section 4.7, that, despite the fluctuating nature of the vacuum, there is a clear distinction between matter or radiation energy and vacuum energy, which is reflected in the electrical (or non-gravitational) neutrality of vacuum dark matter and in the absence of contribution to gravitational entropy by a uniform distribution of vacuum energy (associated with the cosmological constant). On the basis of those developments, it becomes relatively straightforward to provide a clear and unambiguous definition of when it is that matter is present in a vacuum, that would also apply from the viewpoint of observers accelerating relative to a local inertial reference system, or in the presence of very strong local gravitational fields (such as those which are present in the vicinity of a black hole), and therefore the difficulties identified above would now appear to be rather insignificant.

But, in my opinion, one of the most powerful argument that can be used to support the idea that the elementary particle concept still constitutes a necessary and viable element of a consistent interpretation of quantum theory (when it is allowed to obey the limitations imposed by the uncertainty principle) is the observation that, even in the context where it may seem to be the least appropriate to hypothesize about the usefulness of elementary particles, it nevertheless turns out that this assumption allows to explain, in a surprisingly simple way, certain key aspects of the processes involved. What I'm talking about is the use of virtual particles as the mediators of elementary particle interactions. The fact that it would be very difficult to explain certain properties of those interactions, like their range and their strength, without assuming that the interactions themselves are actually mediated by particles, even if those particles remain unobserved and cannot have classically well-defined energy and momentum states, is indicative of the usefulness and indeed of the necessity of assuming that, from a physical perspective, quantum fields actually consist of particles that propagate between interaction events⁸.

The problem we may have, in relation to this conclusion, is that, if parti-

⁸Feynman himself insisted that the concept of an external field becomes relevant merely in the context where the motion of a particle depends on a probability amplitude to interact with the particles mediating this field that varies only with the particle's position at a certain time, as may arise when a large number of such interactions are taking place over a relatively short period of time.

cles do exist as real physical entities, it would seem that it is not possible to attribute a unique position state to those particles at all times, in the context where it is known that there are interferences between the probability amplitudes associated with the many different trajectories that contribute to determine the transition probability for one single event involving the propagation of a particle in a given momentum eigenstate. This is why so many people prefer to assume that the wave function, despite its immaterial nature, may constitute reality itself; a hypothesis which raises difficulties of its own in the context where it must be recognized that this reality would be submitted, by the act of measurement, to discontinuous changes that may violate the spirit of relativity theory and the principle of local causality, even if no information is communicated at faster-than-light velocities. In any case, it must be clear that the wavelike nature of quantum processes is simply a consequence of the fact that the probability amplitudes that must be used in the calculation of transition probabilities are subject to periodic evolution and there is no sense in saying that a *particle* may sometimes evolve as a wave, because the wavelike property is already well-understood as being a property of processes which always involve particles and the problem really has to do with the apparent impossibility to attribute a definite location to those particles, under general circumstances.

What I will explain, however, is that we have not yet exhausted all possibilities and that a realist interpretation of quantum phenomena that would involve elementary particles can still be formulated that would be compatible with the current mathematical framework of quantum field theory (if we allow for a slightly more elaborate particle concept, while still rejecting the contradictory notion of wave-particle duality). I believe that it is, indeed, possible to assume that a unique history *of some kind* is taking place, even for what regards the physical attribute of a quantum system (say its position in space) that is not under observation. This is a conclusion that would obviously contradict the orthodox interpretation of quantum theory, at least in its original form, given that, according to the conventional doctrine, there is no sense in speaking about the state of some physical attribute when no measurement has been effected to actually determine what this state is at a given time. But if we recognize that the elementary particle concept is essential to a consistent interpretation of quantum theory, then it seems that we have no choice but to conclude that the current interpretation of the theory is incomplete, because it does not provide a clear and unambiguous description of what happens when the position of such a system is not under

direct observation. Of course, certain modern interpretations, such as the consistent-histories interpretation of quantum theory, go some way into providing a more realist picture of quantum phenomena, but they also appear to be incomplete, given, precisely, that they allow reality to be described only under particular circumstances, when a certain more or less arbitrary criterion of ‘consistency’ is met and given, also, that, despite their more appropriate handling of the measurement problem, they still fail to explain the emergence and the persistence of a quasiclassical world, as I will explain in section 5.10.

In the introduction to this report I mentioned that I believe that it is essential to adopt a realist interpretation of quantum phenomena if we are to avoid deviating into a solipsistic and idealistic view of reality, according to which nothing would really exist aside from our own mind (if that could ever be found possible). This criterion is particularly important in the context where the only thing that may be considered undeniable about reality is precisely that it is real. The problem is that the adjective ‘real’ is usually assumed to be the characteristic of something that exists as a fact rather than as a mere possibility and therefore the characterization of *quantum* reality as actually being real would appear to exclude the possibility that this reality may not always consist of *observable* facts. Thus, it is important to emphasize that what I have in mind here is the *scientific* concept of realism, according to which it would be deemed appropriate to seek to describe the actual ways by which certain physical processes can occur, even when it is not possible to determine the specific path which is followed in the course of any one particular process. But in the context of the preceding discussion, it would also appear desirable to apply the criterion of physical reality not to the wave function itself, as is usually proposed, but rather to the elementary particle trajectories that enter the sum-over-histories formulation of quantum theory. The hypothesis would then be that it is appropriate to assume that, even in between position measurements, elementary particles follow real and to a certain extent, unique (but not classically well-defined) trajectories in spacetime, despite the fact that those trajectories must, as a matter of principle, remain mere potentialities.

I understand, of course, that, despite being intuitively appealing, the hypothesis that some unobserved dynamic attribute of a quantum system could exist in a definite, unique, but unknown state, even as the various alternative states available to it interfere quantum mechanically with one another, appears to be ruled out by the fact that all possible histories must, in effect,

be put to contribution in order to derive the right probability for a process to occur (that which is obtained by repeating the experiment a large number of times). Faced with this difficulty, one usually concludes that it is not possible to retain a realist description of quantum phenomena that would involve elementary particle trajectories and that it is more reasonable to simply assume that reality cannot be unique in any way between measurements and that the question of what happens to unobserved attributes is meaningless from a scientific viewpoint, as originally proposed by Bohr and Heisenberg. But, if one recognizes that the uniqueness of history is a condition that cannot be ignored, one may alternatively propose that quantum interferences are not indicative of the fact that multiple trajectories must be taken into consideration simultaneously, but rather arise as a consequence of the existence of non-local, but otherwise classically well-behaved hidden variables that would violate the principle of local causality by allowing to determine the course of a conventional history involving otherwise ordinary objects. Without entering into the details of each proposal, it is clear that they are both unsatisfactory, precisely because they both involve assumptions that contradict one key aspect of physical reality (either the uniqueness of history, as an observational requirement, or the absence of instantly propagated effects, as a theoretical requirement). It must be clear, though, that, despite what is commonly believed, the first proposal is just as problematic as its alternative counterpart, even if it was favored by the originators of quantum mechanics on the basis of the fact that it involves fewer arbitrary assumptions.

It always appeared preferable, in effect, to avoid postulating the existence of classical hidden-variables, given that any model based on such a requirement would necessarily involve complex mechanisms of an unobservable nature, whose validity could never be empirically confirmed. Yet, the argument that it is the non-locality of the hidden-variables models that makes them unacceptable is not very satisfactory. Indeed, if one recognizes that there must necessarily be a reality of some kind, then the only *known* alternative to assuming the existence of hidden variables would be to consider the wave function itself as being this reality, which means that non-locality would also constitute an aspect of the orthodox interpretation, because the wave function is also a non-local entity, which is subject to non-local changes, as would occur in the course of certain measurements. Thus, it would appear that the only alternative to a non-local theory, potentially involving complicated arbitrary constructs whose validity would remain unconfirmed, actually amounts to assume that reality is not real (when it is not subject to direct measure-

ment). This is obviously a *simple* assumption, but I'm not willing to accept that it would be mere scientific progress to consider it as a valid assumption about physical reality. One must come to recognize that such a position is not progress, but simple non-sense of the most scientifically objectionable kind. If a *physical* reality exists, then I believe that what is certainly the most basic property that would need to characterize this reality is that it is, in effect, real. This must be considered an essential consistency requirement and neglecting it would again amount to allow a logical contradiction to stand at the basis of our interpretation of the most fundamental of all physical theories.

Therefore, I suggest that one of the crucial points that cannot be neglected in trying to produce a consistent interpretation of quantum theory is that the unique outcome of measurements is indicative of the uniqueness of the history that takes place in between measurements, even for what regards those dynamic attributes that are not subjected to direct observation. The existence of definite causal relationships between all elements of the universe must be understood to actually require that every element of this physical reality is indeed involved in only one such history in any one particular universe. The right interpretation must, therefore, emerge from a combination of two apparently incompatible requirements, which are provided, on the one hand, by this condition of uniqueness of history and on the other, by the necessity to allow quantum interferences to occur between the many distinct possibilities that may exist for the unobservable aspects of this unique history, even as may affect one single quantum process that need not be repeated many times. It is the description of reality we are considering that must adapt to those two requirements if we are to avoid having to alter the rules of logic to accommodate their simultaneous fulfillment. But I do agree with Copenhagenists that this must not be achieved by postulating the existence of hidden variables, whose effects would propagate at faster-than-light velocities, because, from all that we know, the principle of local causality provides as real a constraint on our description of physical reality as the existence of quantum interferences.

In fact, I have come to understand that the debate between Copenhagenists and classical hidden-variables theorists is not as meaningful as one might assume, because the only hidden-variables models that may allow to retain agreement with observational data are those that postulate that the hidden causes of the unique, classical evolution that takes place in between quantum measurements would remain unobservable as a matter of princi-

ple, even when they evolve deterministically (ignorance of the exact state does not arise from a practical limitation that could eventually be overcome, as in conventional statistical mechanics). Thus, even though such classical hidden-variables models would contradict the orthodox postulate of objective indefiniteness, the fact that the hidden variables could never become part of experimental knowledge means that those models do not require a rejection of the concept of objective chance (associated with fundamental unpredictability). It would therefore appear that it is really just the naive, classical definiteness of the phenomenon which is assumed to govern the behavior of quantum particles that is problematic with those hidden-variables models, given that it necessarily requires the propagation of signals of a conspiratorial nature at faster-than-light velocities to achieve agreement with observational data. The real problem for current (classical) hidden-variables theories would then be that, instead of enhancing the domain of validity of the quantum-mechanical description of reality to the classical world, as an improved realist interpretation of quantum theory should enable to achieve, they just allow to perhaps reproduce the empirically confirmed predictions of the theory through some unnatural and complicated contortion of classical reality that makes them even less appealing than the currently favored traditional approach.

But before I elaborate on what kind of physical reality might agree with the two basic requirements identified above (uniqueness of history and local causality), it is important to mention that the requirement that there exists a unique reality is different from Einstein's proposal that reality should be independent from whether or not a certain parameter is being observed, which assumes more than just a unique reality and which is irreconcilable with the mathematical framework of quantum theory. We must recognize as an established fact that quantum reality is not independent from experimental conditions, even if it might be possible to assume that conjugate physical attributes like position and momentum can simultaneously possess unique (even though partly unobservable) values *in a certain sense*, because, as I already explained, this unique reality must also give rise to quantum interferences among multiple states and it is only the physical attribute that is under direct observation at a given time, or in the course of a certain process, that is free of interferences. Assuming that reality is independent from experimental conditions would require that quantum interferences be absent altogether, which is certainly not compatible with any plausible interpretation of quantum theory in its present form.

If the values taken by conjugate observables cannot be determined at the same time with an arbitrarily high degree of precision, it is precisely because the macroscopic constraints necessary to determine the exact state of those physical attributes cannot be realized together at the same time for the same process, while it is those macroscopic constraints (associated with the existence of records) that determine which physical observable is not subject to quantum interferences (for reasons I will discuss in section 5.12). Thus, even though I believe that it is necessary to assume that a unique reality actually exists, regardless of whether it is being observed or not, I also think that it must be recognized that this reality does not evolve independently from the macroscopic physical conditions which determine what can be known, experimentally, of its actual state. Furthermore, it should be clear that the hypothesis that there exists a unique reality of some sort does not impose on quantum particles (say negatively-charged, positive-energy electrons propagating forward in time) that they be distinct individually, even when they possess the same *static* attributes. What we must ask ourselves, therefore, is what the unobservable reality actually is if it does not conform to a classical representation in terms of simple, identifiable objects. Quantum theory, from the viewpoint of its current interpretation, is not so much an answer to the problem of the fundamental nature of reality, as a constraint that must be obeyed by any realist description of physical phenomena.

At this point, it is necessary to mention that I do know that from the viewpoint of someone who has been introduced to quantum mechanics in the conventional way, the requirements discussed above may appear irreconcilable, as it seems that the formalism of the theory itself cannot be dissociated from the Copenhagen interpretation, while the traditional definition of a quantum state would appear to be totally incompatible with a realist interpretation that would involve a unique history. It is only when one begins studying relativistic quantum field theory, that one is introduced to Feynman's approach and the sum-over-histories formalism (despite the fact that conventional quantum mechanics can also be formulated using path integrals), at which point one has already been conditioned to believe that it is not possible to visualize quantum processes as involving unique histories of some sort, while, in fact, this is precisely what the sum-over-histories approach suggests and from a certain viewpoint even requires⁹. In this partic-

⁹It is important to note that even a conventional formulation of quantum mechanics, like Heisenberg's matrix mechanics, can be interpreted as involving a summation over a

ular sense, I was lucky, because I first learned of quantum theory by reading about the problem of interpretation and Feynman's original approach, while I became familiar with the conventional formalism of quantum mechanics only later on, which means that rather than being critical of the reality of Feynman's histories, I was rather critical of the conventional interpretation. I believe that this uncommon course is what allowed me to see more clearly how it can be that each independent elementary particle process consists of a unique (even though partly unobservable) history, despite the fact that there always arise interference effects between the multiple histories which are allowed by the macroscopic, experimental conditions which are imposed on the process. What I would like to explain, therefore, is why it is necessary to assume that the multiple *unique* histories depicted in Feynman's diagrams correspond more than is usually recognized to the actual reality behind all quantum phenomena.

I believe that it is merely the fact that no truly acceptable realist interpretation of quantum theory has ever been proposed that motivates the widespread belief that the multiple histories depicted in Feynman diagrams do not correspond to anything actually occurring (must be considered purely fictitious) and merely constitute useful computational apparatus, despite the obvious similarity between the processes so described and the actual reality we experience. It has become very clear to me that what this formalism provides is nothing but a description of what is actually going on, which we are not able to directly observe, when some dynamic physical attribute of a quantum system is evolving in between measurements. Even ignoring the arguments provided so far in this section, concerning the relevance of the concept of elementary particle in quantum field theory, I think that one must recognize that it is very unlikely that the individual paths entering a sum-over-histories formulation of quantum theory could happen to be intuitively significant simply as a coincidence, without being related to what actually goes on in between measurements of the observable concerned. Perhaps that instead of insisting that our experience of reality is not a reliable guide for judging the value of certain hypotheses concerning unobservable aspects of this very same reality, we should instead try to figure out how the phenomena that cannot be directly observed can be described in a way that would

series of intermediate, unobserved, 'virtual' processes and it is significant that some of the originators of quantum theory were, in effect, open to such an interpretation (perhaps because they were not told by others how they should interpret their own theory), even though they did not see how it could be made truly viable.

better agree with what we do know about physical reality.

It is remarkable in this regard that, while Feynman himself believed that quantum reality involves particles and only particles, he also said that there is no way to explain or to understand what happens to those particles, even during the most simple of quantum processes, because it is not possible to assume that a particle in a given momentum state goes one way or another in space, so that it may be preferable to give up trying to create a model of what is actually happening. I believe that this shows how deeply the philosophy behind the Copenhagen interpretation of quantum theory has become ingrained in our conception of reality, because if one person might have been allowed to understand what is the reality behind all quantum phenomena, it should certainly have been Feynman and it is clear that his failure is in part attributable to the fact that, despite his remarkable insights, as all physicists of his generation he adhered to the notion that a realist representation of quantum phenomena is not possible. But if those difficulties have been allowed to persist to this day, it is merely because we still do not understand the profound meaning of quantum strangeness and remain ignorant of the fact that quantum phenomena *can* be visualized.

What remains to explain, therefore, is how it can be that one and only one of the histories which can be depicted using Feynman diagrams corresponds to what really happens in the course of a specific quantum process¹⁰, despite the fact that it is not possible to attribute to a quantum particle the properties of a classical object, in the sense that one cannot simultaneously determine both its momentum and its position with an arbitrarily high degree of precision. For that purpose, it is necessary to point out that there is something highly problematic with the conventional viewpoint provided by Bohr's complementarity principle. What Bohr suggested, in effect, is not just that the conditions necessary for the measurement of a certain dynamic attribute is incompatible with those necessary for the measurement of its conjugate counterpart, but really that the concepts of momentum and posi-

¹⁰I must mention that I'm aware that a method called 'unitarity' is often used as a shortcut for the determination of quantum probabilities that constitutes a modification of Feynman's original approach, but this alternative technique does not require assuming that the original sum-over-histories formulation of quantum theory is not fundamentally the most accurate and it remains that the summation over all possible histories is more representative of what really goes on at a fundamental level, even if, from a practical viewpoint, it may be even less appropriate than the alternative approach for performing certain calculations under particular circumstances.

tion, for example, constitute mutually exclusive representations of a quantum system, so that it would not even be logically appropriate to conceive of a particle with a given momentum as being a localized entity. If one was to hold on to such a viewpoint, then, clearly, a realist description of phenomena based on the sum-over-histories formulation of quantum theory would become impossible to achieve.

But, in fact, there is absolutely no reason to assume that the indefiniteness of the state of some unobserved attribute of a quantum system cannot be the consequence of a mere incompatibility between the macroscopic conditions necessary for the measurement of one dynamic attribute and those necessary for the measurement of its conjugate counterpart. When one understands the true nature of the constraints which allow decoherence to take place and to rapidly eliminate quantum interferences for the physical attribute that is subjected to measurement (an issue I will address only in section 5.12), it appears quite plausible that quantum indefiniteness actually arises as a consequence of this practical (but fundamental and inescapable) limitation. Therefore, it is not *a priori* impossible for a quantum particle which is known to be in a pure momentum state to follow a unique, but *observationally* undetermined trajectory in space and only the existence of quantum interferences involving multiple distinct trajectories would appear to contradict this conclusion.

There is certainly something true in Heisenberg's statement to the effect that "the progress achieved (through the elaboration of quantum theory) was obtained at the price of having to abandon the possibility of visualizing natural phenomena in a way that is immediately and directly comprehensible to our conventional way of thinking". However, I would insist that what is inappropriate is not the requirement that it should be possible to visualize physical reality, but the requirement that this reality should, in effect, be similar in every way to what it appeared to be before experiments began revealing the existence of interferences between the probability amplitudes (a purely quantum mechanical concept) associated with alternative potential histories. In order to progress toward this legitimate objective of visualizing quantum reality, we may once again consider the classical double slit experiment. What can be learned using this simple, but very general experimental arrangement is that, despite the fact that we are always dealing with discrete, localized particles, interferences, similar to those which can be observed when what is propagating is a classical wave, must be assumed to occur whenever a particle is allowed to propagate between a source and a detector through

more that one possible path without giving rise to the formation of a permanent record that would indicate through which trajectory the particle actually went. Even though such interferences become apparent only in the statistical distribution of measurement results, which is known to depend on the differences between the lengths of the possible paths along which a particle can propagate before its position is detected, the interferences must be considered to apply even for a single process involving the propagation of one unique particle (because the observed interference patterns can be produced even when the photons are sent from source to detector only one at a time). The problem, then, is to figure out how it is possible for a localized particle in a given momentum eigenstate to give rise to those interferences involving many distinct potential paths if, as a particle, it must necessarily propagate in space by going through a definite, yet unobservable trajectory.

Stated in such a way, this aspect of the problem of interpretation appears at once very clear and quite unsolvable. But it took me a very considerable amount of time to simply realize that this is, in effect, how the problem must be stated, as this is not how most people see things. Indeed, it is *not* usually assumed that the particle, as a particle, must necessarily go through a single trajectory or even through any trajectory at all, as this would immediately appear to give rise to an unavoidable contradiction, because, ‘obviously’, a particle cannot go through one trajectory and produce interference effects which involve multiple distinct trajectories. Anyone arguing that this is not necessarily the case would merely be a nostalgic of classical reality that does not accept the ‘undeniable’ strangeness of reality unveiled by the observation of quantum phenomena. Such an approach to the problem of interpretation would necessarily have to deviate into classical hidden-variables and non-local causality. But in fact, that is not the case. Not only is it possible to visualize what is going on when one acknowledges the validity of those premises, but this is the only way to arrive at an interpretation of all quantum phenomena that does not involve any arbitrary and undesirable assumptions that would either conflict with the observed uniqueness of experimental facts, or else contradict one another (as when one speaks of a ‘probability wave’ going through both slits all at once, which then ‘becomes’ a particle when its position is detected), therefore implicitly or explicitly requiring an alteration of the conventional rules of logic.

It is important to understand that, while it is usually believed that logical contradictions may arise when one insists on requiring a realist interpretation of quantum theory, those contradictions are merely a consequence of holding

on to a conventional, or naive conception of reality, according to which it might be possible to obtain simultaneous experimental knowledge about the state of all physical attributes of a quantum system. Indeed, it is usually believed that one cannot assume that all dynamic attributes of a system could be in a unique state at all times without assuming that a precise knowledge of the state of those dynamic attributes would be available (which would violate the uncertainty principle). But, once we recognize that only the second assumption is inadequate and could give rise to factual contradictions, while an absence of knowledge concerning the state of some dynamic attribute that is not subjected to measurement may allow one to assume, without contradiction, that this attribute still evolves along a unique path in any given portion of history (in a certain sense which will be clarified later), then it becomes possible for a realist interpretation to be formulated that is not logically inconsistent (even though such a proposal would normally appear to contain a contradiction).

What I will eventually explain is that the fact that a purely phenomenological model of reality (such as that which constitutes the core of the orthodox interpretation of quantum theory) may appear to be better suited than a realist model for explaining certain observations is merely a consequence of the fact that a realist model cannot be applied to quantum phenomena as they are traditionally described, but only becomes appropriate in the context of a time-symmetric description of those phenomena. Following Einstein, I believe that one must be ready to take an intuitive leap and to derive, based on available experimental data, general postulates that may not always be immediately confirmed through direct observation, but which allow to better model the reality underlying those empirical facts. For what regards the problem of the interpretation of quantum theory, this intuitive leap would actually consist in assuming that the particles involved in the description of elementary quantum processes are, in effect, real and that they are taking part in one unique history of some kind. Once this is recognized to be a legitimate and necessary requirement, the difficulty would then consist in understanding how such a realist description of reality could be made compatible with both the observational constraint imposed by the existence of quantum interferences and the theoretical constraint of a time-symmetric conception of causality.

I think that one cannot be satisfied with assuming that what explains the existence of quantum interferences is the ‘fact’ that a particle doesn’t follow a unique path and actually propagates, from emission to detection,

by simultaneously following, at once, *all* possible trajectories. I believe that the notion that all the available paths are actually followed together in the course of any single quantum process occurring in a given universe actually constitutes one of those strange aspects of quantum reality (as it is usually conceived) which are not merely unexpected, yet unavoidable, but which remain strange because they actually conflict with certain factual aspects of reality. What is quite amazing is that, even though such a notion is only slightly different from the usually rejected viewpoint according to which a particle may go partly through one slit and partly through the other (in a double slit experiment), it is often considered to provide an appropriate depiction of quantum reality. But if one recognizes that such a representation is indeed incompatible with a realist interpretation of quantum phenomena that would not reject the empirical evidence for the uniqueness of historical facts, it remains that one must take into account, in the determination of transition probabilities, any possible trajectory which is *allowed* by the macroscopic conditions which are in effect while those transitions are taking place. In order to accommodate this fact, what is sometimes assumed (as I briefly mentioned in section 5.6) is that a single unique process may actually always involve two interfering histories which, for some reason, would share the same observational conditions. But it remains to explain what justifies this assumption (which would still appear to conflict with the uniqueness of historical fact) and why it can be expected to give rise to the kind of classically well-defined reality we do experience.

It is certainly true that one of the criteria that may allow one to judge the validity of a conception of reality involving unobservable theoretical constructs is its usefulness for producing accurate predictions of experimental phenomena, but this is precisely why the currently favored interpretation must be rejected. Indeed, I believe that if the notion that all histories occur all at once in the same universe is incompatible with the experimentally derived uniqueness of historical facts (in the context where the tentative solution to the quantum measurement problem that is provided by a ‘many-worlds’ approach is recognized to be ineffective, as I will argue in section 5.10), then it must be rejected in favor of a conception of reality that does not require this uniqueness to be a mere illusion. The problem, however, has always been that it would appear that the only realist alternative to such an interpretation would require assuming that the wave function itself is the reality, because, in the context where quantum interference is possible for unobserved attributes, this mathematical object (the state vector more gen-

erally) does not merely provide a probability distribution for the position of a particle which is in a momentum eigenstate, but may involve superpositions of position states with complex-number weighting coefficients, which means that position may sometimes appear to constitute an inappropriate element of physical reality (of course the same is true for momentum under distinct experimental conditions). But, while this is not necessarily inadequate from a mathematical viewpoint, it remains unsatisfactory from a physical viewpoint, especially in the context where this wave function can be subjected to discontinuous changes that would violate the spirit of the principle of local causality whenever the potentialities involved are actualized, as I previously mentioned.

I believe that it is merely the fact that we fail to correctly visualize what is going on in between position measurements (for instance) that makes it look like physical reality cannot involve a unique history of some kind and needs to be replaced by some strange picture which deviates from a conventional representation to the point where reality itself looks unreal, in the sense that the proposed picture is not only incompatible with observable aspects of reality, but also with the logical consistency which is known to apply under more general circumstances (which would allow one to reject the possibility that a particle could be in one place and also in a different one, all at the same time, in one single portion of history). What holds the key to a better understanding of quantum reality is the acknowledgment that what can be known about a quantum system does not allow one to tell everything about how it actually evolves, even though, as a matter of principle, there does not exist a more accurate description of the processes involved. Such a standpoint is the only alternative that is available when one considers it inappropriate to assume that dynamic attributes simply do not exist when they are not those about which direct experimental knowledge is available.

Although the approach I favor may, at first, seem problematic, it is actually much simpler to apply than its logical alternative, because the idea that a dynamic attribute does not exist when it is not subject to observation cannot be adapted to the case where such an attribute is only known to an intermediary level of precision, because, clearly, either an attribute exists or it doesn't, while it is undeniable that the state of any attribute can be determined with more or less precision by the appropriate measurement, as long as an inversely proportional uncertainty applies to its conjugate counterpart. If, at least, it was possible to speak of certain quantum systems as definitely being observed, while other systems would not (under particular

circumstances), then it might perhaps make sense to assume that what is measured is real and what is not measured doesn't exist, but in fact, there is always something that is known with arbitrarily high precision about a physical system, as long as the system remains causally related to the rest of the universe (this is what is implied by the linearity of Hilbert space) and it is merely the conjugate dynamic attribute of this known attribute (however unnatural it is) which is completely undetermined, so that if one was to choose to follow the orthodox approach, then, based on those considerations, one would be forced to somehow ascribe both reality and absence of reality to the same physical system (even though not to the same physical attribute), which again constitutes a logical contradiction.

Thus, despite what one is usually encouraged to believe, it seems necessary to assume (particularly if one wants to avoid having to consider the possibility of a reality created through observation) that two systems prepared in the same quantum state may evolve differently at the level of the dynamic physical attribute whose state is not determined by the macroscopic conditions of an experiment. I believe that this is what explains that a subsequent measurement of this originally undetermined attribute may produce outcomes that differ from one system to the other (from one experiment to another) and if this is correct, it would mean that it is inappropriate to assume that it is the act of measurement itself that introduces randomness into our description of quantum phenomena. The more consistent approach I will propose, therefore, allows physical systems which are described by the same wave function to actually be different at a certain unfathomable level, even if the wave function still provides the most complete description of a quantum system and of how it will evolve.

From that perspective, it becomes apparent that something very problematic comes into play with the conventional interpretation whenever post selection is involved in the determination of which physical attribute of a system is actually measured (as would occur in the context of the delayed choice experiments discussed in section 5.6). Indeed, if one assumes that only measured attributes are real, then it would mean that what is real at the present moment depends on what choice will be made in the future regarding which attributes are to be measured. This is so embarrassing that it is usually considered to support the view that quantum theory is not about reality at all, but about the outcome of measurements, while, in fact, what the reality of post selection illustrates is rather the awkwardness of the conventional interpretation of quantum theory, in the context of which the reality of a

physical attribute itself is dependent on what measurements are performed, either now or in the remote future. Once the necessity of a realist approach (according to which all relevant physical attributes are assumed to be real regardless of experimental conditions) is recognized, then all that one must avoid is taking the easy way out and postulate that there exist hidden variables of a classical kind that would require violations of the principle of local causality as a consequence of trying to explain in too simplistic (but actually quite complicated) a way how the multiple possible histories of unobserved dynamic attributes are allowed to interfere, even in the course of one single quantum process.

In order to achieve a realist description of quantum phenomena that does not contradict other essential aspects of reality, it is necessary to first understand that the most significant difference between the sum-over-histories formulation of quantum theory and the statistical mechanics of classical systems (think about the phenomenon of Brownian motion in particular) has to do with the existence of the quantum phase that allows interferences to arise between the different possible histories involved and which is attributable to the use of probability amplitudes, instead of classical probabilities, as elements of the summation process. From that viewpoint, what needs to be explained is how it is possible for a particle to follow a path along which all of its unobserved dynamic attributes have unique values at all times, despite the fact that the many trajectories which can be followed by any one such attribute would seem to interfere with one another, as if no definite trajectory was ever followed. At this point, it may still appear justified to simply reject as implausible the hypothesis that there must, in effect, exist such a unique path. Once the requirement of a time-symmetric description will be taken into consideration, however, it will become clear that it is just as inappropriate to refuse to admit the existence of those unique trajectories, as it would be to refuse to recognize the existence of elementary particles themselves. John Von Neumann was certainly right when he claimed to have demonstrated that the ordinary reality of everyday objects cannot apply to quantum particles if those objects are to obey the principle of local causality. But, as I will explain, that does not necessarily mean that we need to reject the notion that particles always follow a unique trajectory of some kind (in the space of their unobserved attribute), which still constitutes a valid hypothesis as long as we allow for this unknown trajectory to conform to the requirements of a time-symmetric conception of causality.

If the sum-over-histories formulation really constitutes a fundamentally

different formulation of quantum theory that cannot be derived from earlier formulations by a simple mathematical transformation, as is usually understood, then one cannot reject the possibility that it is only by considering the reality it describes for what it is that we can begin to understand quantum theory. From that perspective, it is certainly incorrect to argue, as many authors do, that quantum theory is only about the probability of measurement results and does not tell us anything about what goes on in between measurements. If the most adequate and general of quantum-mechanical formalisms does involve a certain description of what happens in between observations, then it would seem that it is merely our failure to understand why it is exactly that this description is relevant from a physical viewpoint that motivates our rejection of this realist picture of phenomena and that justifies the commonplace belief that the formalism is not indicative of anything more profound. In any case, one must keep in mind that the prevalent opinion that what the sum-over-histories formalism indicates is that all paths are followed all at once in the course of any single process is not an unavoidable conclusion and that it cannot be claimed that no other choice exists for a realist description of quantum phenomena. What I will explain is that it is still possible, in effect, to assume that a quantum particle must merely be *allowed* to take any of the available paths, but that it does not actually go through all paths in the course of one single process. It is not true that we are confined to contradictory assessments of reality and that it is necessary to assume that quantum theory is about particles and yet that it is not about unique particle *histories*.

What I would like to emphasize is that it is not the hypothesis that there exists a unique and variable (but unobservable) history which is incompatible with experimental facts, but rather the usually preferred hypothesis that similarly prepared systems always evolve in identical ways in between measurements. Indeed, it is clearly the measurement results which are unique and variable, while it is merely our current assumptions regarding what remains unobservable which may turn out to be inappropriate. It must be clear, though, that I'm not claiming that the mathematical framework of quantum theory is incomplete, because I do recognize that it is impossible to provide a more accurate description of the state of a system than is allowed by the uncertainty principle, so that, even if it is real, the unique history of an unobserved dynamic attribute remains a mere potentiality for any specific process. Experimental knowledge of both the exact momentum and the exact position of a particle is not allowed by the basic structure of quantum

theory and I believe that this is a conclusion that cannot be overturned. In the language of the consistent-histories interpretation of quantum mechanics, one would say that the simultaneous determination of a particle's momentum and position can only take place on decoherent 'branches' of history, which, from my viewpoint, actually means that it cannot occur at all, because this would require distinct macroscopic constraints to exist together simultaneously (for the same system) and if one wants to preserve the character of uniqueness of physical reality, then, obviously, one cannot argue that one set of mutually exclusive macroscopic constraints exist at the same time as a different one.

As a means to accommodate the uniqueness of measurement results in light of the existence of quantum interferences between the multiple possible histories of unobserved dynamic attributes, what is usually proposed is that all histories actually occur all at once in different 'branches' of the same universe, but that it is precisely the decoherence effect that allows *observed* reality to appear unique, given that it requires the interferences that may exist between the different states of a dynamic attribute to vanish very rapidly upon a measurement of this physical attribute. However, as I will explain in section 5.12, it appears that decoherence can only achieve the goal of giving rise to a quasiclassical world if we require the existence of a unique history of some kind. Once the dust has settled, it appears that no valid argument actually remains that would support the validity of the hypothesis that all histories are followed at the same time, in the same universe, as different co-existing and interfering branches. Thus, by assimilating what I believe to be the only appropriate interpretation of quantum phenomena, we will go from a situation where it is necessary to assume that, either there is no reality at all between measurements, or else that all histories are followed all at once, to a situation where it is no longer possible or necessary to embrace such logically inconsistent viewpoints and where we are allowed to once again conceive of a universe as involving one single and unique history of some kind, which, in effect, constitutes the most essential element of a physically meaningful definition of what a universe actually is.

What emerges from those considerations is that it is the very notion that decoherence is responsible for eliminating the interferences between the many co-existing 'branches' of history that makes quantum entanglement problematic, given that it requires the existence of non-local hidden variables, to enforce the selection of one branch over another following measurement, due to the fact that this selection is, in effect, a global phenomenon. I know that

many people do not agree with that, because they assume that the multiple branches of history are causally independent from one another, as if they actually consisted of different universes. But the problem, once again, is that there is a logical contradiction here, because we cannot assume that we are dealing with truly independent branches, while those branches would nevertheless be assumed to exist in the same universe (so that they can interfere with one another). A lot of crazy things have been said concerning why those two assumptions may not be incompatible with one another, but in the end, one must recognize that the simple truth is that there is a contradiction and if the branches interfere prior to a measurement, then there must exist non-local hidden variables propagating effects faster than the relativistic speed limit, to enforce the global consistency of measurement outcomes at multiple remote locations, in the presence of quantum entanglement, when it is assumed that all possible histories are indeed followed together in the absence of measurement. Thus, it is not absolutely true that the phenomenon of quantum non-locality cannot be used to constrain our concept of physical reality in a way that would require it to better agree with certain properties of this reality which are known to apply under more general conditions, as Einstein once sought to achieve.

It is telling, therefore, that it is quantum entanglement which is usually assumed to forbid a more conventional, realist description of quantum phenomena. Indeed, the violation of Bell's inequality by the results of multiple different experiments which have been performed on pairs of entangled elementary particles proves that a naive concept of reality, according to which all dynamic physical attributes are in a unique *classical* state at all times, could not be considered valid unless this reality explicitly involves non-local influences. In fact, what was shown by the experiments in which a violation of Bell's inequality occurs is that non-local correlations do arise at the most fundamental level of description of physical reality. But this does not necessarily mean that non-local hidden variables must exist that would propagate effects faster than the relativistic speed limit, because this property may instead be a simple reflection of the fact that the basic structure of reality is richer than we usually assume, in the sense that it could be governed by a more general concept of causality that is not limited by the constraint of thermodynamic time asymmetry. Given that quantum entanglement is made manifest through quantum interference, the non-locality that is discussed here is not different from that I have already identified as emerging whenever one assumes that the wave function itself constitutes physical re-

ality. I believe that what this actually means is that it is not the hypothesis that there exists a unique history which is problematic, but the notion that quantum non-locality must necessarily involve a violation of the causal structure of spacetime imposed by relativity theory.

I have already emphasized, in the discussion about time-symmetric causality from section 5.3, that backward-in-time causation is not forbidden by relativity theory. But it should be clear that backward causation, even when it is restricted to operate in accordance with the principle of local causality, may actually give rise to non-local correlations. The important point here is that the existence of such correlations would not allow faster-than-light communication, given that the backward propagated influences are submitted to the constraint of diminishing entropy in the past that is imposed by the constraint of global entanglement and in such a context information is only allowed to flow from the past toward the future and never in the opposite direction, while a flow of information toward the past would be required for faster-than-light communication to occur. Amazingly, this is precisely the property that is observed to be obeyed by the non-local correlations which have been experimentally demonstrated to occur in the course of certain quantum phenomena, as a result of entanglement. I believe that this is not just a coincidence, but that it actually confirms what I have said concerning the time-symmetric nature of causality and the crucial role played by this property in a quantum-mechanical context.

If this is the true origin of quantum non-locality, then it would mean that what is actually ruled out is merely the existence of non-local hidden variables that would violate the principle of local causality by propagating effects at faster-than-light velocities (which would allow faster-than-light communication and therefore also the flow of information from the future toward the past), while the non-local correlations that follow from backward-in-time causation would actually be a fact which we were traditionally allowed to ignore only because it does not allow signals or information to be communicated instantaneously (due precisely to the origin of those non-local correlations) and therefore can only be revealed through subtle correlations of otherwise random outcomes of measurements, performed on carefully entangled quantum systems. What should be clear, in any case, is that the observed absence of backward-in-time signaling need not be a consequence of the inadequacy of a realist time-symmetric interpretation of quantum theory, as it can also be a consequence of the effectiveness of the constraint which was identified in section 4.9 and that gives rise to the thermodynamic arrow of time under more

general circumstances. Only if that was not the case, would the backward-in-time causation that may be involved in giving rise to quantum non-locality be allowed to violate the principle of local causality that is enforced by relativity theory. It is not appropriate to conclude that the experiments which have confirmed that certain quantum phenomena involve non-local correlations have proven that those phenomena are irreconcilable with any commonsense interpretation of the theory. What must be abandoned is not scientific realism, but the traditional interpretation of quantum theory which forces us to reject the principle of local causality and to return to a conception of reality that would involve instantaneous action at a distance.

It is important to note, in this regard, that it is the locality assumption that would allow one to conclude, based on the results of certain recently performed experiments described in Ref. [49] which involve multiple entangled photons, that there may coexist many mutually incompatible accounts about what constitutes a known, or observationally confirmed fact. Those experimental results, which involve the violation of certain inequalities similar to, but distinct from the conventional Bell inequality, were initially assumed to support the claim that factual truth is a relative notion (and therefore that reality may not be objective), a conclusion which would appear to confirm the relevance of the relational interpretation of quantum theory. But once we recognize that quantum non-locality is not optional and that it was actually shown, by even more straightforward methods, to itself constitute an unavoidable aspect of reality, then the inappropriateness of the radical conclusions which were drawn, based on the results of the above discussed experiments (regarding the lack of objectivity of observationally derived facts), becomes all the more obvious, even aside from the fact that they would (once again) have given rise to logical contradictions. Thus, it should be clear that the assumption that the experimental results obtained in one part of such an experimental setup cannot be correlated non-locally with those obtained in a remote part of the same setup is incorrect and it is only when we are not willing to take this aspect into consideration that those experiments seem to imply that reality is not objective and that the truth of certain experimentally established facts, which all happened in the same universe, may be an observer-dependent aspect of reality. I believe that this only shows how important it is to recognize that non-local correlations do arise in the quantum realm, even if effects are always constrained to propagate at velocities no larger than that which is imposed by the relativistic speed limit associated with the light-cone structure of spacetime, either forward or backward

in time.

What I have tried to make clear in this section is that it is highly preferable to adopt a realist interpretation of quantum phenomena, because all alternative proposals involve logical contradictions at one point or another and those difficulties are always attributable precisely to a rejection of scientific realism. What is unsatisfactory, however, is the absence of a realist interpretation that would agree with the multiple specific constraints imposed by the mathematical structure of quantum theory, like non-locality or quantum interference (more generally). I believe that if the orthodox interpretation of quantum theory is still preferred by most researchers in the field, despite the fact that it requires rejecting scientific realism, it is because something essential is missing from all known realist interpretations that could make one of them acceptable. The problem to which I will now turn, therefore, is that of explaining in an intuitively satisfactory, but logically consistent way, without rejecting as mere illusion the uniqueness of historical facts, why it is that the probability amplitudes associated with the many trajectories available to a quantum particle interfere with one another when its position state is not under direct observation, as if the particle actually followed several different trajectories all at once in the course of one single process. It is here that it will finally be shown that, despite what is usually believed, this is not an impossible task.

5.8 Time-symmetric quantum theory

It is quite amazing that one single requirement allows to satisfy, all at once, both the condition of scientific realism in face of quantum interference or state superposition and the principle of local causality in face of quantum entanglement. This requirement is that of time-symmetric causality. There should be no doubt, indeed, that the only way one can avoid having to conclude that there exist non-local influences propagating faster than the relativistic speed limit, in the context of a realist description of quantum phenomena, is by assuming that certain effects actually propagate backward in time. But it is usually believed that such backward causation would be even more problematic than the existence of non-local hidden variables. I must admit that I don't really understand what motivates that opinion, which to me appears even more arbitrary than the rejection of negative energy states. What could be worse, indeed, than an outright rejection of relativity theory

and the principle of local causality and what could be more difficult a task than rebuilding quantum physics from the ground up, while trying to provide a consistent classical hidden-variables theory that would allow to match all empirical constraints, by postulating explicitly non-local influences? But what is even more significant is that, as I have explained in sections 5.3, 5.4, and 5.5, the alternative of a time-symmetric conception of causality, is actually well-founded from a purely theoretical or epistemological viewpoint and would constitute a highly desirable development in the context where it is recognized that there can be no fundamental distinction between the past and the future at the most fundamental level of description.

What must be understood is that backward-in-time causation is not necessarily problematic, even if the finality it involves may appear unnatural from the viewpoint of our conventional, unidirectional experience of time. First of all, in a universe where entropy cannot grow in the past, backward-in-time causation would not allow us to tell the future in advance. But, as I already explained, it is also clear that backward causation does not allow one to change a known fact from the past. Classical causality, or the pairing of the distinction between causes and effects with the thermodynamic distinction between past and future, only comes into play at the macroscopic level where time asymmetry emerges from the constraint imposed by the presence of negative-energy matter on the initial Big Bang state at which global entanglement must take place. In other words, our experience of classical, unidirectional causality is not necessarily incompatible with backward causation, as long as the effects which are propagated backward in time do not give rise to the kind of backward-in-time signaling that would require entropy to grow in the past.

Now, I previously mentioned that what quantum entanglement appears to allow is precisely the kind of non-local correlations that would arise from such backward-in-time propagation of effects, which is required to occur with ever decreasing entropy in the past and which, for that reason, is not allowed to give rise to faster-than-light communication, as would be the case if classical hidden variables were responsible for quantum non-locality. A consistent interpretation of quantum theory would be one that naturally agrees with this limitation in all situations, despite the fact that it would allow to explain non-local correlations. This must be considered an absolute requirement of any realist approach in the context where no violation of the classical (unidirectional) principle of causality has ever been observed to take place in the course of any measurement on entangled systems.

If this is correct, then we need to ask how it is exactly that such backward causation is allowed to take place, in the context where the only particles we know about that do propagate backward in time are antimatter particles, while such particles are not always involved in the experiments which have revealed the existence of non-local correlations. What I have come to understand is that, in fact, such time-symmetry is precisely what the mathematical structure of quantum theory naturally requires, as my discussion of the two-state-vector formalism from section 5.6 emphasized. Indeed, as I previously explained, a mathematically equivalent formulation of quantum theory is possible that involves two state vectors, one of which provides the state of a system as determined by past measurements, and the other the state of the same system as will be determined by future measurements. In between measurements, those two state vectors evolve in a conventional ‘unitary’ manner, in the future following a past measurement, and in the past preceding a future measurement. Of course, this is not a realist representation of quantum phenomena, as we are still dealing with wave functions, but at least, it shows that a formulation can be provided that allows to reproduce all the predictions of quantum theory (sometimes more naturally than even the standard theory) while satisfying the requirement of a time-symmetric description of quantum reality (whatever this reality turns out to be).

One clear advantage of such an approach is that it allows the time-symmetry that is implicit in the original theory to be preserved even when non-local correlations exist and the order in time of two measurements performed on a pair of entangled particles is dependent on the state of motion of an observer. Indeed, when the chronological order of two measurements is an observer-dependent property (which would occur whenever the events are separated by a space-like interval), a process of state vector reduction which may appear to be triggered by a measurement performed on one entangled particle, from the viewpoint of a certain observer, would appear to be triggered by the measurement performed on its entangled counterpart, for a different observer. But it would be problematic to have to choose one or the other of two such measurements as being the cause of the outcome of the other measurement if it was not also possible to assume that it is this other measurement that is the cause of the outcome of the first one, because, in such a case at least, there is no objective criterion that would allow one to tell which event is the cause and which is the effect. Yet, from the viewpoint of the conventional approach, it would appear that it is necessarily the event that happens first that is the cause of the other event,

even if this first event actually happens later from the viewpoint of a different observer. In the context of a time-symmetric formulation of quantum theory, however, the fact that future measurements are allowed to influence the present state of a system means that a certain reciprocity is allowed between the measurement that influences and the measurement that is influenced (both measurements exert an influence on the outcome of their remote counterpart). In other words, it is no longer necessary to assume that there exists an absolute distinction between a cause and its effect, from a purely quantum-mechanical viewpoint, and this actually allows to avoid the contradiction that would otherwise emerge when we are dealing with measurements performed at space-like separated events on entangled systems.

What one must retain is that if it was not for the existence of backward-in-time causation, then a clear distinction would need to exist between the causes of state vector reduction and their effects, even when we are dealing with entangled particles. But given that, in such a case, this distinction may be an observer-dependent property, then it would appear that the spirit of relativity theory would be violated, even if it would be impossible to say exactly what distinguishes the cause from its effect, because, from a traditional viewpoint, this distinction would be required to exist. The fact that, in all known situations where non-local correlations may arise, backward-in-time signaling is not allowed to occur, clearly shows that unidirectional causality is not involved in the determination of the outcome of the second of two measurements performed on a pair of entangled particles, because if it was involved, then there would be no reason *not* to expect backward-in-time signaling to occur, at least in some reference systems. From such a viewpoint, it would appear that the prevalent belief that causality must always operate forward in time is motivated by expectations similar in nature to those which *originally* motivated the formulation of the Lorentz transformation (the contraction of physical objects in motion relative to absolute space), because imposing a unidirectional conception of causality, in the context of quantum non-locality, amounts to postulate a property of reality which, even if it did pertain to the physical world, would be required to have absolutely no distinguishable effect on it.

Now, even though the two-state-vector formulation of quantum mechanics represents a step forward, the fact that it still does not provide a realist picture of quantum phenomena that would fully accommodate the particle concept and the requirement of a unique history of some kind, means that it cannot be the final answer to the problem of interpretation. Clearly, some-

thing essential is still missing and it is only after much questioning and while trying to figure out how the two-state-vector formalism could be generalized to agree with the sum-over-histories formulation of quantum theory that I was able to obtain a truly consistent, realist picture of quantum phenomena. I have become convinced, in effect, that the bold intuitive leap which I previously suggested one must be ready to take to achieve a more realist interpretation of quantum theory actually consists in recognizing that what we are dealing with here is a set of two unique histories (involving unique particle trajectories) unfolding in opposite directions of time without directly interacting with one another in any way.

In such a context, what matters is not really the direction of propagation in time of the particles involved in those processes, but an overall direction of time that only differs in a relationally defined way, such that, if the two histories were to be otherwise identical, they would still differ in that the directions of propagation in time of all the particles involved would be opposite for those two histories. But in fact, it is not possible to differentiate in any absolute way initial causes from final ‘causes’ and it is only the difference between the directions in which the two histories unfold in time that has physical meaning and this relationship must be preserved even when the processes actually occurring in the course of those two histories differ in ways not forbidden by the macroscopic experimental conditions imposed on those processes.

The important point here is that the path followed by a quantum system in the space of its unobserved dynamic attributes must, in effect, be allowed to differ for the retarded and the advanced portions of a process (the ordinary process and its time-reverse counterpart), even though the observed attributes of both portions of the process must share the exact same history. What really happens, therefore, to a photon on its way to a detector in the double slit experiment (see Figure 5.1) is not that it passes through both slits all at once, but that it has the *possibility* to pass through any one of the two open slits in *both* the retarded and the advanced portions of the same process (when the actual trajectory remains observationally undetermined), which therefore requires that both paths be taken into account in the determination of transition probabilities for any given process, even though a photon only ever goes through one particular slit in the retarded portion of history and then again through one particular (but possibly different) slit in the advanced portion of history. It is simple to verify that those assumptions allow to reproduce the predictions of the standard theory in any specific and

possibly more complex situation (I will explain below why this should indeed be expected).

It is merely the fact that we do not observe the actual trajectory followed by the photon that makes it necessary to consider both possibilities, all at once, for any single process, given that under such conditions this trajectory can be different for the retarded and the advanced portions of the process. But this does not necessarily mean that the trajectory is actually different for the two histories, only that it *can* be and, as I will explain below, this is sufficient a motive for requiring that both trajectories be considered to contribute to the estimation of transition probabilities. Any one history still involves a particle following a unique, unobservable trajectory, only, each process involves both a retarded history and an advanced history (a pair of histories taking place in opposite directions of time) which are only required to share the same macroscopic experimental constraints. Those histories are therefore allowed to differ in all aspects which are not constrained to a particular subset of possibilities by the observable ‘macroscopic’ conditions (the paths can differ as long as no permanent record of those differences ever becomes available) and this is why the many different possible paths available to a quantum system interfere with one another and must therefore be taken into account in the determination of the probability for the complete process (comprising those two histories) to occur.

Remarkably, if it was not for the fact that probability amplitudes, unlike conventional probabilities, involve periodic variations, which allows them to interfere constructively or destructively, then it would be impossible to deduce the existence of the advanced portion of a process (which may actually be any one of the two histories), because it is the periodic or wavelike aspect of probability amplitudes which allows the retarded and advanced portions of history to interfere, when the dynamic attributes involved are not subjected to direct experimental determination. The greater consistency of the viewpoint proposed here is apparent in the fact that it is no longer necessary to assume that, when the path followed by a particle is not observed, the object actually behaves as if it was a different entity (a classical wave), because the interferences which are made conspicuous in the statistical distribution of measurement results can be explained without requiring one to assume that the particle behaves differently when its position is not observed. What changes, when a different dynamic attribute is submitted to observation, is merely the macroscopic conditions imposed on a process, while the system involved still follows a unique, but unknown, and possibly different trajectory

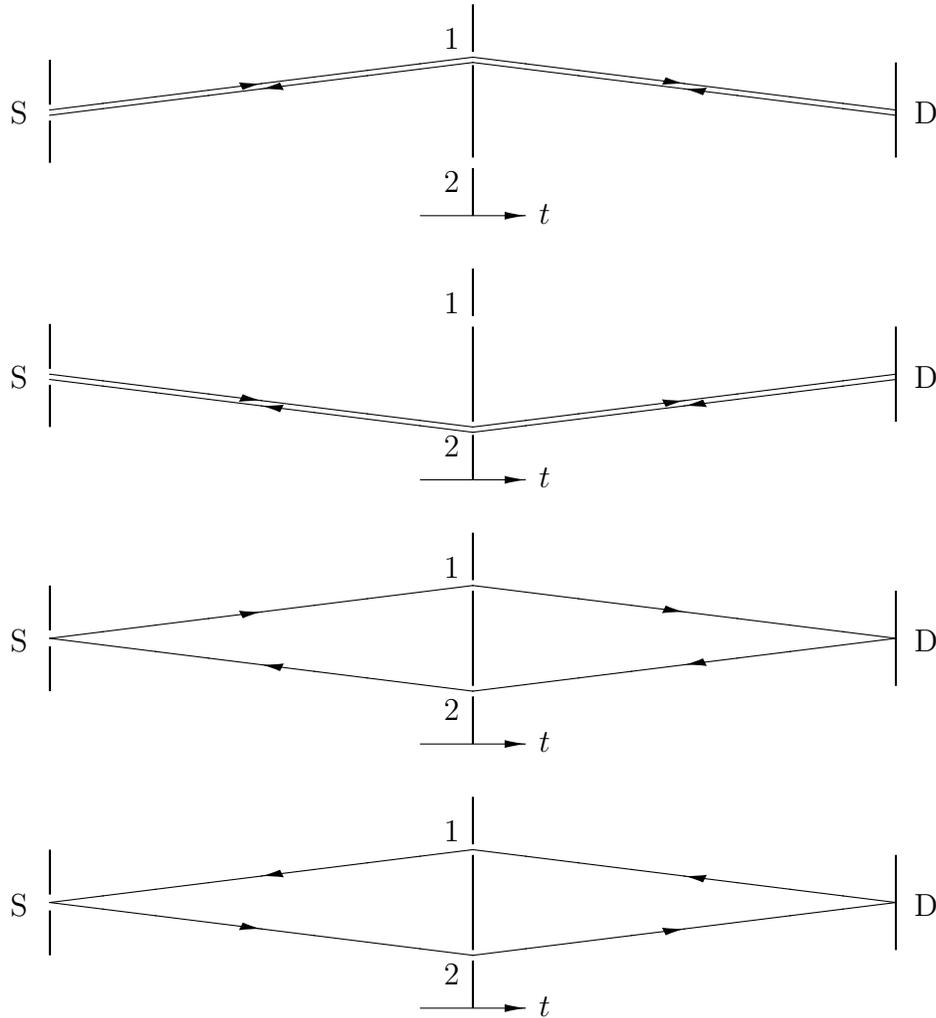


Figure 5.1: The four possible, combined, retarded and advanced histories of a double slit or simple interferometer experiment, with one source S and one position detector D , when the actual trajectory of the quantum particle remains unknown. The direction of the arrows corresponds to the flow of time. When the actual trajectory of the particle is subject to experimental determination, only the first two combined histories remain possible and the two trajectories no longer interfere, as the retarded and the advanced histories must be the same for any complete process.

in the retarded and advanced portions of history, which unfold in the space of the unobserved attribute.

It is only when a particle is constrained by the experimental conditions to follow a certain definite path (when a record of the actual slit through which the particle went is available) that interferences are absent, for the dynamic attribute involved, because, in such a case, the particle must follow the same path during both the retarded and the advanced portions of history. It would, therefore, be incorrect to maintain that it is not possible to visualize what occurs to a photon as it propagates from source to detector in a double slit experiment, when its trajectory is not observed. It is not nonsense to speak of the passage of the photon through one particular slit, even when this trajectory remains experimentally undetermined, as long as one recognizes that the actual trajectory can be different for the retarded and advanced portions of the process. From this viewpoint, what looks rather absurd is the conventional assumption that an elementary particle whose trajectory is undetermined follows, at once, all possible paths. When it is properly understood, quantum theory is no longer as unsettling as it used to be (this comment will become even more apposite when other essential aspects of this approach are discussed, which allow to justify its inevitability).

From the viewpoint of the interpretation of quantum theory proposed here, there would no longer arise logical contradictions in the description of the state of a system when a certain dynamic attribute of the system is in a state of superposition (which is always the case for at least one physical attribute). We may consider, for example, an electron whose spin has been measured to be up along the horizontal axis. Under such conditions, the spin of this electron along the vertical axis must be considered undetermined. But this cannot be understood to mean that the spin of the electron is either up along the horizontal axis and up along the vertical axis, or else up along the horizontal axis and down along the vertical axis, as one might consider appropriate from a classical perspective. Whenever one tries to experimentally confirm the apparently indisputable validity of this legitimate hypothesis, the results one obtains show that it cannot be true. It may therefore appear that whenever an electron is in a definite state of spin relative to the horizontal axis, its spin state along the vertical axis, if it is real, must be such that it cannot be described without violating the conventional rules of logic, because it would seem to be allowed to point along two mutually exclusive directions all at once, which, from a realist viewpoint, does, in effect, constitute a contradiction.

But once it is understood that two independent histories are involved in any single process, then it becomes clear that what the results of the discussed experiments mean is not that the vertical spin of the electron is in no state at all (which would require rejecting the possibility of a realist description of quantum phenomena), or that it is, at once, in all possible states (which would require rejecting the conventional rules of logic), but merely that while its vertical spin state in the retarded portion of history can be either up or down, the same vertical spin state can also be either up or down in the advanced portion of history, which means that four different combinations of states are allowed, thereby contradicting the hypothesis that this vertical state can only be either up or down and nothing else, for any single process, or at any single time (which actually consists of two different times that must simply correspond with one another for the retarded and advanced portions of history, as I will later explain).

One is therefore allowed to assume that the spin of an electron along any axis is always in a unique, but possibly unobservable state in any one portion of history, even though it is not in a unique state for any *process* (when a process is adequately considered to involve both a retarded and an advanced portion), as experiments confirm. Thus, if those experiments with electrons, as well as more decisive observations of the same kind, do show that quantum strangeness is unavoidable, it would be incorrect to assume that what they demonstrate is that a realist interpretation of quantum theory is impossible and that there cannot be an unique reality of some kind behind the observed phenomena. Indeed, the contradictions which are encountered in the context of a more conventional, realist interpretation are only made apparent in the statistical distribution of measurement results and always concern physical attributes which actually remain unobserved, while it is precisely at this level that the alternative interpretation proposed here differs from the conventional approach. But given that this realist, time-symmetric interpretation of quantum theory allows non-local correlations to arise, even when no effect is propagated at faster-than-light velocities (as I will show in the following section), then it also appears inappropriate to argue, as is often done, that only a rejection of scientific realism (the idea that there must exist an ‘objective’ reality between measurements) may allow to avoid the conclusion that quantum non-locality arises from instantaneous action at a distance. Quantum non-locality is not an illusion, but action at a distance can be avoided, even in a realist interpretation.

It is the fact that, traditionally, time-symmetric interpretations of quan-

tum theory involved classical wavelike phenomena that made them undesirable as realist interpretations. But once the dual nature of the state vector is understood to be a consequence of the existence of two actual histories in which particles propagate, once through any of the available paths and then again through any of those same available paths, but in the opposite direction of time, then the time-symmetric nature of quantum reality becomes a more significant asset, given that it allows to reproduce the statistics of quantum measurement results and to explain interference effects involving distinct paths, while naturally providing a picture of quantum reality that satisfies the requirements of scientific realism. Of course, the reality unveiled here is not classical, because it involves probability amplitudes instead of classical probabilities and it requires the existence of an unobserved counterpart to every process (because we really experience only one of the two portions of history at any single time). But then, what we are dealing with *is* quantum reality and not classical reality and only consistency provides an unavoidable criterion for judging the validity of any experimentally accurate representation of reality. If quantum strangeness itself cannot be avoided, then there must certainly remain some unexpected element in any empirically established model. In fact, it appears that it is the remaining ‘incomprehensible’ aspects of quantum reality that make the theory truly consistent in a way that would be impossible classically and, as such, they are certainly not undesirable.

As I explained in the preceding section, consistency merely dictates that physical reality must, in effect, be real and therefore unique in *some* particular way, but it does not *a priori* constrain this reality to conform to some preconceived criteria of appropriateness we may believe should apply, that would be based on an experience of physical reality which is restricted to a subset of experimental conditions, namely those where quantum interference and entanglement are usually unapparent. What’s more, I’m not suggesting that two processes are taking place in parallel that could differ from one another in an observable way, in violation of the uniqueness of historical facts, but merely and precisely that there is a counterpart to history which, even if possibly distinct from its time-reverse version at the level of intricate details, would nevertheless remain identical from the viewpoint of its observable macroscopic features, even though it would still be required to exist in order to explain certain observable features of reality (the interferences). Therefore, what constitutes a decisive advantage of the time-symmetric interpretation of quantum theory proposed here (over the usually favored approach accord-

ing to which all paths are followed together, all at once, in one single portion of history) is that it naturally agrees with the observation that all results of quantum measurement are, in effect, unique, so that it is no longer necessary to try to provide an independent explanation (such as the hypothetical splitting process of a many-worlds interpretation) for the fact that all potentialities are not actualized all at once, as would appear to be required if all histories actually occurred all at once, as is usually assumed. This means that there is no longer a problem associated with the objectification of measurement results.

It seems that the error that is made, in the context of most current interpretations of quantum theory, is that we fail to recognize that, if we were to take into account the existence of the advanced portion of every quantum process, it would simply no longer be necessary to assume that all paths from either the retarded or the advanced portion of history are somehow being followed all at once, because the simple fact that the advanced portion of the process can be different from the retarded portion, while still obeying the macroscopic constraints of the experiment, is sufficient to guarantee that it is only when all possible paths are taken into consideration that the right predictions, concerning the probability of occurrence of the whole process, will be obtained. Those considerations are also valid in the case where we are dealing with the outcome of one single event (like the passage of a unique photon from source to detector in a double slit experiment), even if time-symmetric processes do not always involve all interfering paths all at once (but merely two of them), because, in the context where probability amplitudes are involved, it is possible for the probability of one single process (composed of a retarded and an advanced portion) to contribute negatively to the final probability of a process and as I will explain below, this actually allows all the different alternatives to contribute to the probability of one single process.

Another advantage of such a realist, time-symmetric interpretation (involving both forward and backward propagated effects) is that it allows to enforce the global consistency of factual aspects of the world in a way that is particularly significant in the case of entangled systems. Indeed, if one is to assume that the retarded and advanced portions of history share the same observable, macroscopic conditions (a requirement whose validity will be justified in section 5.12), then the result of a measurement performed on one of two entangled particles must be compatible with both the experimental conditions of this measurement and those of any measurement that

may eventually be performed on the other particle, because in any reference system (from the viewpoint of any observer) there is as much causal influence from the first measurement on the second, as there is from that second measurement on the first (especially when those two events are separated by space-like intervals). What is apparent here, therefore, is that a quantum-mechanical description of reality involves some form of causal circularity of the kind that would arise if time travel was a possibility. But, as I mentioned in section 5.4, the only problem that may arise in the context where such closed causal chains would be considered a possibility does *not* have to do with the fact that, if global consistency is always preserved, this would seem to contradict our expectations regarding free-will (a difficulty which is significant merely from the viewpoint of our conventional, unidirectional experience of time), but with explaining how it is, in effect, that global consistency (the idea that all facts must agree with one another under all circumstances) can be preserved, regardless of what happens. What remains to understand is how it is that this requirement is enforced at the level of time-symmetric quantum-mechanical processes.

It should be clear, first of all, that quantum theory does appear to be the appropriate framework for implementing global consistency, as it already allows to appropriately handle the closed causal chains occurring as a result of the existence of antiparticles as negative-energy particles propagating backward in time. Thus, the usual approach to estimating the probability for a process to occur, which amounts to sum up the probability amplitudes for all possible ways by which a process can occur and then to take the square of this complex number, appears to allow global consistency to be satisfied, only, it is not completely clear why, in effect, such an annoying procedure allows to produce consistency, in the context where backward causation would be assumed to constitute an unavoidable aspect of a quantum-mechanical description of reality. To understand what is going on, it is necessary to first recognize that a complete quantum process (one to which can be attributed a certain probability) actually consists in the combination of a retarded history, unfolding from a given past state toward a given future state through one particular unobservable path forward in time, followed by an advanced history, unfolding from the same future state toward the same past state through another particular and still unobservable path backward in time, or vice versa (as it may be the advanced history that is ‘followed’ by the retarded history backward in time).

Thus, it is essential that the two possible segments of history, which are

unfolding in opposite directions of time, be combined to actually give rise to one complete time-symmetric process, to which can indeed be assigned a definite classical probability (instead of a mere probability amplitude). It would then be by adding the probabilities for all such combined, time-symmetric processes which are compatible with the observable past and future experimental conditions that we would obtain the final correlation probability. Now, even though such a procedure can be shown to produce transition probabilities equivalent to those of the conventional approach under similar circumstances, the problem is that it is not always possible to obtain meaningful results from such a procedure, unless one limits the scope of the questions that can be asked, concerning the history of a system and its environment, by adopting a suitable coarse-graining. It is only under such conditions (when certain details are left aside concerning the processes which are described) that classically meaningful probabilities can be obtained for various alternative histories.

In the context of the conventional, modern interpretation of quantum theory (the consistent-histories interpretation), what this would be assumed to mean is that, when described with a maximum level of details, certain histories are simply nonsense and cannot be considered to actually occur as real physical phenomena. This would be the case, for example, of the history of a photon as it goes from source to detector in the conventional double slit experiment, when the particular path taken by the particle is not subjected to direct observation, because it seems that one cannot obtain a classically meaningful probability for a unique history of such a kind. But I believe that this self-imposed and somewhat arbitrary restriction, concerning what can be considered real of reality itself, is not appropriate and arises merely because we do not understand the profound significance of those apparently inconsistent probabilities, which only emerges when they are considered in the context of a realist and fully time-symmetric interpretation of quantum theory.

It must be clear that what I find problematic about the formalism of consistent histories is the restriction that is usually imposed regarding what can be *meaningfully* described of quantum reality, not the logic of the conclusion, made in the context of the conventional interpretation of quantum theory, concerning what can be *classically* described of quantum reality (which histories can be assigned classically meaningful probabilities) and under which circumstances. What I'm trying to explain is that the criterion of decoherence, which is imposed on families of coarse-grained histories in the context

of the consistent-histories interpretation of quantum theory, is not really a criterion for assessing what is consistent of reality, but merely a criterion for assessing what is *classically* well-defined of this reality. I believe that it is incorrect to argue that common sense logic (conventional logic) is increasingly less adequate for describing reality, when we consider increasingly smaller scales, even though it is certainly true that the probability that various alternative histories interfere with one another rises as those histories are being described with an increasing amount of detail (using a finer coarse-graining). Clearly, either conventional logic applies or it doesn't and one cannot try to justify how nonsense could be made acceptable by relying on the confused notion of complementarity (the apparent freedom to describe the same reality with mutually incompatible concepts).

The problem with this conventional interpretation of quantum strangeness is that it would cease to provide a logically consistent picture of reality¹¹ precisely on the scale where the unconventional phenomena we may want to understand are occurring (the quantum scale). But this difficulty arises merely when we fail to understand that conventional logic applies *not* to the observed phenomena themselves, but to the unobservable, unique reality, which consists in each of the two portions of history that unfold in opposite directions of time for every process, on any scale. The fact that conventional logic still applies on the classical scale, even from a more traditional viewpoint, can therefore be understood to result, not from the fact that reality is only consistent on such a scale, but from the fact that the two time-reversed portions of history must always be the same on such a scale (for reasons I will explain in section 5.12).

Anyhow, what is usually considered undesirable of the probabilities that may sometimes be obtained for a combined pair of histories is that they can assume negative values, or normalized values larger than one. I believe that one can only begin to understand why the existence of negative probabilities in the intermediary stages of the estimation of a final transition probability is not catastrophic when one recognizes that it is precisely the circularity of all quantum causal chains (that follows from the existence of an advanced portion to every quantum process) that enforces the consistency of the present with a given future (while the retarded portion enforces the consistency of

¹¹It must be clear that my use of the term 'consistency' does not have the meaning it has in the context of the consistent interpretation of quantum theory, where it refers merely to the classical definiteness of a history.

this future with the known present, as is usually understood). In the context where one must take into account the existence of quantum interferences, it appears necessary, in effect, to impose on the quantum phase that it returns to a value that is as close to its initial value as possible, after a complete, time-symmetric process has occurred (once forward and once backward in time), if this unobservable parameter is to have any physical significance. Of course, this initial value can be any arbitrarily-chosen one, as only *changes* to the phase and the amplitude of the wave function, occurring as a result of the evolution that takes place during the retarded and advanced portions of a process, are significant. In other words, if the phase was originally π radiant it cannot end up being 2π radiant (if the amplitude of the wave function remains unchanged) after a complete time-symmetric process has taken place, otherwise a contradiction would have occurred, as those two phases are the perfect opposite of one another and therefore correspond to two maximally distinct unobservable initial conditions (of the phase itself) which can only belong to two mutually exclusive instances of physical reality.

In the present context, probabilities larger than one merely constitute another facet of the same problem, because, as Feynman pointed out [50], a greater than one probability for a given process to occur is equivalent to a negative probability that the same process will *not* occur. What I'm suggesting, then, is that, whenever the probability for a process to occur in one specific way is negative, one must assume that, if the process would occur in this specific way, it would diminish the chances that the observable macroscopic conditions which would have actually given rise to it existed in the first place, thereby making the sum of probabilities for all the possible ways the process could occur smaller than it would otherwise be, given that it would make the initial conditions themselves less likely to *have* occurred (because the probability that the process would occur in such a way would decrease the likelihood that the process may occur in any possible way, instead of increasing it as is usually the case). Thus, when a pair of minimally coarse-grained histories (composed of both a retarded and an advanced process) has a negative probability of occurrence, this can be interpreted as diminishing the chances that the process involved may occur by following any possible path (even those for which there is no destructive interference). Likewise, when the probability for an individual pair of minimally coarse-grained histories to occur is larger than one, this can be interpreted as decreasing the chances that *alternative* initial conditions existed, which is another way to say that it would actually increase the chances that the actual initial conditions that

gave rise to this history did indeed occur.

Thus, one can speak meaningfully not just about a reduction or an increase in the probability that some future outcome is realized when some past conditions are observed, but also about a reduction or an increase in the probability that a given set of initial conditions has actually been observed to occur, whenever a given set of future conditions will be satisfied. Those additional contributions to the conventional measures of transition probability are dependent only on the degree of compatibility between the unobservable, initial quantum phase and the quantum phase that is obtained as a result of the phase change that occurs in the course of the whole time-symmetric process (the combination of a retarded and an advanced history). From this viewpoint, therefore, a process is allowed to influence the very probability that certain boundary conditions necessary for its occurrence may be found to exist, not just in the future, but in the past as well. In the context of a time-symmetric interpretation of quantum theory, it should actually be expected that such effects would arise, given that there is necessarily as much influence of the future on the past, as there is of the past on the future, which forbids the initial macroscopic conditions to be determined independently from what happens in the unknown future.

What transpires, therefore, is that when the retarded and advanced portions of the history of a given unobserved physical attribute are such that they require changes to the quantum phase that would not allow it to return to its initial value (as would occur when the probability amplitudes associated with the two possible paths in a double slit experiment interfere destructively), then one must assume that the probability that the very initial conditions of the process (the emission of a photon with such an energy by this particular detector at this particular time) could themselves be observed to have occurred is reduced in proportion to the magnitude of this contradiction. But if those initial conditions cannot be expected to have happened, then it means that the pair of minimally coarse-grained histories with which is associated a negative probability would merely contribute to reduce the probability that the process is actually observed to take place following *any* of the available paths, forward and backward in time, and this is why no single instance of a complete time-symmetric history (involving only two possibly different paths) contributes to the probability of a process independently from the other possible time-symmetric histories (involving all the other possible paths), despite the fact that only one such history actually happens.

As a result, even though the normalized probability that a given observable history actually occurs must, in effect, be a number between zero and one, any *unobserved* portion of history that contributes to determine this final probability could have a negative probability of occurrence, or a probability larger than one and still be describable as consistently as any other portion of history. The interpretation I'm proposing here, therefore, does not require one to reject as meaningless the histories with which are associated negative probabilities, as those time-symmetric processes can be interpreted in a realist way and do not differ fundamentally from other time-symmetric histories (occurring once forward and then backward in time), given that they do contribute, in a meaningful way, to establish the final, positive transition probabilities for an observable (sufficiently coarse-grained) history to occur.

It is merely the fact that negative probabilities can only arise when quantum interferences are present, while, in general, interferences are only apparent when the actual path followed by a quantum system is not subjected to direct observation, that explains that we appear to be justified to assume that negative probabilities cannot arise and must be rejected as physically meaningless, because it is true that the validity of theoretical estimates, regarding the probability for such individual portions of history to occur, cannot be confirmed through direct observation, as a matter of principle. Once again, it is merely the fact that our experience of reality is limited to the portion of it that is directly accessible to our senses that explains that we have never experienced negative probabilities and that we view them with suspicion, as if the histories they characterize could not be real. What I have explained is that this self-imposed limitation, concerning the scope of a realist description of quantum phenomena, is not necessary and once this is understood, then *all* the histories that contribute to establish the statistics of quantum processes can be given the status of physical reality, as required.

Thus, the occurrence of negative probabilities in the context of the realist, time-symmetric interpretation of quantum theory I'm proposing should not be considered a problem all by itself (that would justify rejecting as unreal the histories which give rise to those unconventional measures of probability), because I have shown above how it can be assigned a clear meaning in such a context (this is the originality of my approach), as long the negative values involved do not show up in the final results of the estimation of a transition probability for an *observed* history. In fact, from a purely formal perspective, the proposed approach may be considered even more adequate than the

traditional method, given that it always involves only the summation of real probabilities (one real, but possibly negative number for each time-symmetric process) instead of mere probability amplitudes (complex numbers with no independent physical meaning).

In any case, it is now apparent that the most important weakness of early time-symmetric quantum models (such as Cramer's transactional theory, discussed in section 5.6) is that they required assuming that the advanced waves which are part of a complete 'handshake' process propagate backward in time in the same portion of history as that in which the retarded waves propagate forward in time, instead of occurring as part of an independent segment of history, which would forbid any interaction with the processes taking place in the retarded segment¹². It is important, therefore, to understand that, even though the retarded and advanced portions of history share macroscopic experimental conditions and even though their durations also correspond to a certain extent (for a given process), they do not take place simultaneously (even in the opposite order) in the same segment of history (how this can actually be made reasonable will be discussed in section 5.12). It is precisely the fact that we are dealing with two different portions of history (not really occurring at the same epoch) that allows the principle of local causality to be satisfied, despite the fact that the model proposed allows non-local correlations to arise, because the particles which are propagating in one of those two portions of history do not, in effect, interact with those which are propagating, at the corresponding moment, in the time-reversed portion of history. As a result, this alternative approach allows to do away with advanced waves as real, *classical* waves and this means that, contrarily to the early time-symmetric models, the interpretation of quantum theory proposed here is not a particular instance of classical hidden-variables theory.

This is certainly a suitable characteristic of the proposed model, because, as I previously mentioned, it is now understood that in order to reproduce the results of certain experiments in which quantum entanglement is involved (the EPR-type experiments, which will be discussed in the following section), classical hidden variables would need to violate the principle of local causality

¹²For those reasons, the time-symmetric interpretation of quantum theory proposed here cannot form the basis of a solution to the problem of advanced waves, because, in the present case, we are not dealing with advanced propagation as it could be observed to occur in the same portion of history and this shows, again, that the problem of the absence of advanced waves must be considered independently from the problem of the interpretation of quantum theory.

through complex and highly unnatural mechanisms. I believe that a similar unnatural coordination of influences, now affecting experimental conditions, would be required if we were to assume, instead, that quantum non-locality is an illusion attributable to what has been called ‘absolute determinism’, or the idea that every choice of measurement is determined in advance as a consequence of deterministic evolution. It is clear to me, indeed, that from a physical viewpoint, this latter proposal merely constitutes the same classical hidden-variables theory in disguise, because, in the absence of non-local hidden variables, the puzzling predetermination which it requires would remain unexplained and this means that the hypothesis would simply be inadequate.

What adds to the difficulties facing all such interpretations is that it was experimentally demonstrated, not so long ago [51], that the classical hidden-variables hypothesis is in fact incompatible with the results of certain measurements that can be performed on a single quantum object, for which entanglement is irrelevant. Basically, what those experiments are designed to achieve is a measurement of five pairs of attributes of a photon that is in a state of superposition of three position states. When the experiments are performed, it emerges that the statistical distribution of measurement results is incompatible with what is allowed in the case where classical hidden variables (of the naive realist kind) determine the outcome of those measurements, because the choice of which pairs of attributes are to be measured affects the outcome of the measurements. What those results were immediately assumed to imply is that what is not measured of a quantum system cannot be considered to exist independently. But it must be clear that, in this particular case, just as in the cases where quantum entanglement is involved, what has really been demonstrated is not that there cannot exist a unique reality of some kind in between measurements, but that this unique reality cannot be classical in nature (it cannot involve a retarded history and nothing else, for the unobserved physical attributes).

I have already explained, however, why reality cannot be uniquely characterized, in the classical sense, in between measurements and those developments clearly show that it is not necessary to reject the hypothesis of a unique reality for the retarded and advanced portions of a process (independently), regardless of whether a physical attribute is being measured or not. In the case at hand, it seems that what is happening is that the different possible measurements are affecting the constraints which are exerted on the unobserved retarded and advanced portions of history and are thus

allowed to give rise to different patterns of interference for some related physical attributes, as also occurs in the case of entangled systems (more about this in the following section). But the conclusion that there is no reality, independent of what is revealed by measurements, is not made unavoidable by those experiments, even if it is certainly true that this reality cannot be classical and must be conceived of in accordance with the requirements of time-symmetric causality.

It is also the fact that reality is not classical, even though it is unique in a certain sense, that allows to explain the otherwise puzzling thought experiments proposed by Yakir Aharonov, Jeff Tollaksen, Sandu Popescu, and their colleagues which are described in Ref. [52]. Those experiments involve sending three electrons on two possible paths in an interferometer and then effecting some post selection (see section 5.6) on some of the electrons to influence their past states backward in time and in the process give rise to quantum correlations between the states of the electrons involved. What is remarkable, here, is that according to quantum theory, even if you send three electrons at a time in the interferometer, no two electrons will ever appear to have gone through the same arm of the interferometer during any single trial, as if it was possible for three electrons to simultaneously go through two possible paths without any two electrons ever going through the same path. But the paradox associated with such a thought experiment only arises when we fail to understand the fact that the trajectory of the electrons is only unique in the sense discussed above.

What would be proved by those experiments (if they were actually performed) is merely that, when none of the particles are directly observed to follow one path instead of another, then no pair of particles can, in effect, be determined to follow the same path classically, that is, for both the retarded and the advanced portions of the process. But this does not mean that a given pair of particles may not be following the same unobservable paths during either the retarded or the advanced portions of the process, as long as those trajectories actually remain unobserved, which is precisely the outcome of the condition imposed on the final state in the experiments discussed here. No three particles can go through two different paths without two going through the same path, only, particles can go through no specific path during the complete time-symmetric process and obviously, if no particle goes through a single path from a classical viewpoint, then no *pair* of particles can go through a single path either (from the same viewpoint), even if, from a realist, time-symmetric viewpoint, there are always at least two

particles following the same unique, but unobservable path, either forward or backward in time.

Finally, it is also important to mention that, even though the energy signs of the particles present in the advanced portion of history considered here are well-defined relative to the energy signs of the particles present in the retarded portion of history, in the sense that any positive-energy particle that is observed to be propagating forward in time would be related by the observable macroscopic conditions to a negative-energy particle propagating backward in time (those assumptions will be justified in section 5.12), this does not mean that all the particles present in the retarded portion of history would have positive energy signs, while those present in the advanced portion of history would have negative energy signs. In fact, in each of the two corresponding segments of history there may be both positive- and negative-energy particles propagating in any direction of time and all that we can assess is that the positive-energy particles which are *observed* to propagate forward in time in one of the two portions of history must have negative energy (and positive action) as they propagate backward in time in the corresponding time-reversed portion of history. It should be clear, therefore, that there is no correspondence between the particles present in the advanced portion of history that is required to exist by the time-symmetric formulation of quantum theory discussed here and the unobserved negative-action particles, whose properties were described in chapter 2 and which actually propagate in the same segment of history as ordinary positive-action particles, even though they also have their own counterparts in the advanced portion of history, as any other matter component.

5.9 Quantum entanglement and non-locality

Before I turn to the quantum measurement problem and share the most significant insights I have gained while working on the problem of the interpretation of quantum theory, I would like to return to the important question of the viability of a realist description of quantum phenomena in the context of the existence of non-local correlations. It is possible, in effect, to apply the interpretation which was developed in the preceding section to provide a realist, yet locally causal description of the processes taking place in the course of an experiment of the Einstein-Podolsky-Rosen type, involving pairs of entangled photons. What I will explain is that the experimentally confirmed

violation of Bell's inequality does not make unavoidable the conclusion that instantaneous action at a distance must be an integral aspect of any realist interpretation, in the sense that we are still allowed to assume that no effect can propagate faster than the relativistic speed limit in the course of the retarded and advanced portions of history, when those histories are conceived as taking place independently at two different epochs. The fact that I'm allowed to actually explain the existence of non-local correlations in such a way is significant, because, contrarily to what is often believed, quantum non-locality is not only unexplained from a *classical* perspective, it is not explainable at all in the context of the conventional interpretation of quantum theory.

To help visualize the phenomenon of quantum non-locality, we may consider, for example, a simple interferometer experiment in which the source, instead of emitting one photon in one direction (for which two possible trajectories would be allowed), would emit a pair of entangled photons which would be allowed to travel in opposite directions, each along one or another of two possible trajectories in which they would meet a mirror (one mirror for each path of a given photon) that would direct them toward a detector (one detector for each photon) that would allow to determine either the presence of interferences between the two paths available to a given photon (by simply detecting the arrival of the photon), or the exact path a photon took on its way to the detector (through a measurement of the photon's angle of impact). What's particular with such an interferometer experiment is that, when we choose to measure the angle of impact of one of the two photons we also, inevitably, determine which path the other photon took in its own otherwise independent part of the experiment, so that, even if we try to measure interferences between the two paths available to this other photon, we do not observe any. This correspondence is made unavoidable by the fact that, when the angle of impact of the first photon is determined, the angle of impact of the second photon must be the exact opposite of that of the first photon in order that momentum be conserved in the initial state. Thus, it is only when we choose not to determine the exact path taken by *any* of the two photons that we are, in effect, allowed to observe the presence of interferences between the two paths available to each of them.

From this perspective, it is apparent that, when we are dealing with entangled systems (when the phases associated with the propagation of two otherwise independent systems have become entangled as a result of local contact), the choice of whether to measure an angle of impact or the pres-

ence of interferences between multiple paths is not made locally, but constitutes a global property of the experiment, because, whenever an angle of impact is measured for one of the entangled photons, it is no longer possible to measure interferences between the multiple paths available to the other photon. All that must be understood is that what enforces the global character of this choice of measurement is the existence of an advanced portion of history. Indeed, whenever we choose to determine the exact path taken by one of the two photons, the result of any measurement performed on the other photon must reflect that choice, because the effect of the measurement performed in the future on the first photon propagates backward in time (through the advanced portion of the first process) to the initial entangled state (at the source) and then forward in time (through the retarded portion of the second process) to affect the measurement result performed on the second photon, even when this measurement is separated from the measurement performed on the first photon by a space-like interval. It must be clear that this backward causation does not determine what the result of one particular measurement is whenever a measurement is performed on the other particle (this is the outcome of conservation laws), it only determines if interferences can actually be observed at any of the two detectors¹³.

Indeed, you may recall my earlier discussion regarding the fact that, in the context of the proposed realist time-symmetric interpretation of quantum theory it becomes possible for an unobservable quantum process to influence the very possibility that the observable (initial and final) macroscopic conditions to which it is subjected could have existed. What I have explained is that those conditions should allow the quantum phase to get back to a value as close to its initial unobservable value as possible after a complete time-symmetric process has occurred (once forward and once backward in time), if the conditions necessary for the process to happen are to themselves be allowed to have occurred in the first place. When the outcome of a complete

¹³This is why information cannot be sent at faster-than-light velocities using the property of quantum entanglement, because the outcomes of individual measurements are still random from a local viewpoint and to verify the presence of interferences one would need to repeat the experience a large number of times, while using different configurations of the interferometer, but this would allow one to tell through which path a given photon went, by noting the time at which it arrived. Therefore, in practice, it is not possible to modulate a signal that would produce immediately recognizable effects remotely, by locally varying the type of measurement performed, even though more subtle non-local correlations can be observed which still originate from the advanced portion of those individual processes.

time-symmetric process, along one possible trajectory, results in a final phase that is interfering destructively with the initial phase, then negative probabilities arise (with an amplitude determined by the phase shift involved) which contribute to decrease the final measurable probability that the observed process actually occurs by following any of the possible trajectories available to it, because the probability that the (initial and final) conditions necessary for the process to occur in such a way are satisfied is then itself reduced.

But it is exactly in such a way that non-local correlations are made to happen in the experiments discussed above, because, under such conditions, interference effects observed at one location (on one particle) depend on the experimental conditions observed at a different location (on a different particle), simply because the phase changes which arise in the course of those processes are occurring as a result of the boundary conditions applying on the complete time-symmetric process and not just on some portion of it associated only with one or another particle. What those experiments (during which a violation of Bell's inequality is observed to occur) have revealed is that it is possible to demonstrate the existence of such non-local correlations, which cannot be attributed merely to the conditions imposed by conservation principles on the total momentum (or polarization state) of the two entangled photons, in the context where they are created by pair in the initial state (which merely requires that one photon goes through the upper path when the other goes through the lower path, even when those trajectories are not observed and remain classically undetermined).

In any case, as soon as the angle of impact of a photon is measured at one or the other detector in the experiment described above, then it is no longer possible to measure interferences between the two possible trajectories for any of the two photons, because the retarded and the advanced trajectories are then exactly the same in all possible cases and for each photon, which means that there is no phase change for a complete time-symmetric process. If it seemed impossible, from a conventional viewpoint, to assume that the photons propagated along a unique trajectory prior to such a measurement, it is because the measurement is what determines whether interferences will be observed or not, for both particles, and when interferences are indeed present the trajectories of the two photons are no longer well-defined from a classical viewpoint, which has always been interpreted to mean that there is nothing we can say of reality under such conditions. But what emerges, from the more consistent perspective adopted here, is that it does appear possible

to assume that each photon follows a unique, causally independent trajectory as it propagates toward its detector, from the moment when it is emitted by the source and right up to the moment when a measurement is performed on it, contrarily to what would appear to be allowed according to the orthodox interpretation of quantum theory, only, it turns out that in certain cases (when interferences are present) this unique trajectory can be different for the retarded and the advanced portions of the propagation process, so that it cannot be argued that the photons are always in a *classically* well-defined position state as they propagate toward their detectors.

What happens, therefore, is that the presence or the absence of quantum interferences between the multiple trajectories available to one of the two entangled photons as it propagates toward its detector is determined by the choice of which measurement is performed on the second photon in the future, as a consequence of the existence of the backward-in-time-propagating influences attributable to the advanced portion of the history of this second photon. Thus, the trajectory of the first photon must already be either classically well-defined or quantum-mechanically superposed right from the moment when the particle emerges from the initial entangled state, at the source, in order that those conditions actually agree with the observational constraint set by any measurement that could be performed on the second photon in the future. Any measurement performed on the first photon itself must, therefore, agree with the constraints set by the choice of which measurement is performed on the remotely-located second photon. But, given that there is also an advanced portion of history that is experienced by the *first* photon, then the classical or superposed nature of the trajectory followed by its entangled counterpart, as it propagates in the future toward its own detector, is also required to agree with the choice of measurement that is to be performed on the first photon itself in the future.

As a result, even if the experimental conditions that determine which attribute will be measured by the detectors are changed once the photons have already been emitted, the initial retarded and advanced states will already be such as to reflect that future change and will evolve in accordance with those altered conditions, because the initial retarded states of the two entangled photons are influenced by the choice of measurements performed at each detector in the future (through the advanced portion of the processes). Thus, the whole experimental setup with which the photons will interact in both the past and the future determines what is allowed to happen to both of them, even as they just leave the source in the forward-propagating version

of history, before they interact with a detector. It should be clear, therefore, that the measurement that is performed on any one particular photon cannot *alone* and in every circumstance determine what happens to both photons, as one must assume from a conventional viewpoint, because unidirectional causality is not involved in conveying the influence that propagates along the advanced portions of the processes, from each future measurement back into the initial entangled state.

There is no additional complexity involved here (non-local correlations can be established without any information carrying signal being sent backward in time, as required). Each measurement unveils the state of a photon at the moment when this measurement takes place, but the choice of which measurement is performed influences the state of the photon as it reaches the source in the advanced portion of history, just as when a past state influences a future state, only now backward in time and without entropy increase. But given that, in the above discussed experiments, this past state is an entangled state which results in the two photons sharing a common phase, then it follows that the past state of the first photon is also causally influenced by the choice of measurement performed on the second photon in the future, in a way that is not that different from the usual manner by which effects are propagated forward in time, except that information cannot be carried by the effects so produced, given that entropy cannot rise as they propagate in the past direction of time. This is all made unavoidable in the context where the retarded portion of history experienced by the two photons must share the same observational constraints as apply to the advanced portion of history which is experienced by the same two photons (actually two identically prepared photons existing in a corresponding, but different portion of history), for reasons I will discuss in section 5.12.

What this shows is that, instead of insisting that the wave function may not be real, or that it merely represents the state of knowledge of one particular observer, which must be actualized on contact with information from another, previously independent observer (as one postulates in the context of an interpretation of quantum theory such as ‘QBism’), we should instead recognize that the wave function does provide our best account of the exact quantum state of a system at any time, but that there are two such states (associated with two actual histories), one of which is evolving forward in time and the other of which is evolving backward in time. In such a context, the fact that the wave function may sometimes appear to be a subjective property, dependent on whether information concerning the conditions of a

future measurement to be performed on a system is available or not, can be seen to be a mere consequence of the fact that we cannot know in advance what the backward-in-time-evolving state is, before we obtain information about the outcome of that future measurement, even though it already affects the present state of the system. The adequacy of the latter viewpoint appears to have been confirmed by the fact that the process of actualization of quantum potentialities is now understood to be a consequence of concrete changes that take place in the environment with which a system becomes correlated under very specific conditions (those responsible for triggering decoherence), thereby contradicting the hypothesis that it might be a subjective phenomenon.

In any case, what must be understood concerning the interferometer experiment discussed above is that the unique trajectories of the two entangled photons are only required to be made identical in the retarded and advanced portions of history, when it is the angle of impact of the photon that is measured at one or the other (or both) of the two detectors and not the interferences. Indeed, even if what happens at the source is influenced through backward causation by what occurs at the detectors, if the two detectors are set to determine merely the presence of quantum interferences, the measurement performed by the second detector cannot determine the trajectory that was not determined by the first one and neither is the measurement performed by the first detector allowed to determine the trajectory that was not determined by the second one, given that neither of those two measurements allow to determine through which slit one of the photons went on the way to its detector and this must be reflected in the interfering nature of the trajectories of *both* photons, as they propagate between the source and the detectors in the retarded and advanced portions of history.

Thus, if one of the two photons in an EPR experiment of the kind described above is found to have traveled along one particular path, as a result of the choice of measurement performed on it, then both photons will propagate along one particular path during both the retarded portion of history (away from the source) and the advanced portion of history (toward the source). It is only when none of the two photons in an entangled pair is experimentally determined to have traveled along one or another of the two possible paths (as a result of the choice of measurement to be performed on both particles) that it can no longer be assumed that the trajectory of both photons is classically well-defined and in such a case interferences between the multiple possible trajectories would indeed be observed for both photons.

This does not mean, however, that the two photons would not have unique and corresponding trajectories in both the retarded and advanced portions of the process under any conditions (when one photon goes through the upper path, the other photon always goes through the lower path), merely that, when those trajectories remain unobservable, they can be different, for the same photon, in the two portions of history. For any type of EPR experiment, it is only the attribute of a particle that is correlated to that of the other particle by the conservation principles applying on the initial state (which in the example discussed here would be the one-parameter angle of impact associated with the momentum state of the photon) that must necessarily be classically determined for one particle, when it is so determined for the other particle.

As I mentioned above, what justifies the widespread belief that the unobserved attribute of the photons in an EPR experiment were in no state at all before at least one of the two measurements was performed is simply the fact that, when neither detectors are set to measure the correlated attribute, interferences between multiple intermediary states can arise which cannot be classically described. But once we recognize the influence exerted by the advanced portion of the processes involved, then it once again becomes possible to consider that the photons follow unique trajectories at all times in both the retarded and the advanced portions of any one particular history, even when it is not the correlated attribute (the angle of impact) which is measured at any of the two detectors. This is certainly appropriate, given that it is not possible, in general, to tell which of two measurements (on one or the other photon) determines the time at which the intermediary states could be considered to no longer interfere and to actually become real¹⁴. The only requirement, therefore, is that there is always a correspondence between the trajectories followed by the two photons in both portions of history (such as if one photon goes through the upper path, then the other must go through

¹⁴This is easier to understand in the context of EPR-type experiments involving pairs of linearly polarized photons, in which case it is merely the *difference* between the angles of polarization which are measured by the two detectors that determines if there are interferences or not. The fact that, in the experiment described above, one of the detectors may appear to be privileged in determining the presence or the absence of interferences at both detectors, when only one of the detectors measures the state of the correlated physical attribute, should not be considered to undermine the validity of the hypothesis that the classical or superposed nature of the trajectories is, in general, determined by the configuration of both detectors.

the lower path), as we are, in effect, dealing with correlated states. But unless it is the state of the correlated attribute that is actually measured at one of the two detectors, it is not possible, as a matter of principle, to tell what those corresponding trajectories really are in any particular case. What's interesting is that, once the validity of this viewpoint is recognized, it follows that the idea that the concept of a localized particle may no longer be valid in the presence of quantum entanglement and that it should be replaced by a holistic concept of reality at a fundamental level is no longer justified and actually loses most of its appeal.

At this point, it is necessary to mention that I'm perfectly aware of the fact that Murray Gell-Mann (among others) once argued that the idea that EPR-type experiments imply a certain form of non-locality is merely a distortion of reality, because (so he argued) what occurs when the angle of impact is measured for one of the two photons is merely that we find ourselves in one particular 'branch' of history, where both photons happen to follow a definite trajectory. What is problematic here, however, is not merely the fact that such an explanation would depend on the validity of the contradictory notion that a photon goes, at once, through all available paths (in many 'branches' of the same history, until a 'splitting' process takes place and all potentialities are actualized, all at once, but presumably no longer interfere with one another), the real difficulty has to do with the fact that, from such a viewpoint, unnatural coincidences would still be observed that would remain unexplained, because the choice of which measurement is to be performed on one of the two photons affects the outcome of measurements performed on the other photon in a given 'branch' of the universe's history and it is not possible to explain how even such a coordination of measurement results would occur in a certain branch of history where the global outcome would, in effect, be observed. The truth is that this rejection of quantum non-locality is equivalent to the absolute deterministic viewpoint discussed in the preceding section, given that it requires one to assume that it is possible for pre-existing correlations to exist, which are not attributable to any locally propagated effect. It is quite ironical, therefore, that it was suggested that this viewpoint constitutes an alternative to classical action at a distance, because, as I previously explained, in the context of a realist interpretation, absolute determinism is actually a form of classical hidden-variables theory and therefore must involve faster-than-light signal propagation.

It is now possible to be more specific regarding why it is that we would not be justified to assume that local causality is violated when a state vector

is reduced, in a more general context, as when a photon is emitted by a source, whose propagation is described by an expanding spherical wave function, after which its presence is detected in one particular location, thereby affecting the wave function over the entire volume. I believe that, if the conclusion that the principle of local causality is violated under such conditions is not unavoidable, even when we acknowledge the fact that the wave function always provides the most accurate description of the state of a quantum system, it is because this phenomenon can be described using the same realist, time-symmetric approach which I used to explain the origin of quantum non-locality as it arises in the case of entangled systems and which involves two histories (independent from the viewpoint of local causality) unfolding in opposite directions of time.

What happens is that the spreading wave function allows to accurately describe the results of any measurement that would reveal the existence of interferences between the multiple paths through which the photon might have traveled as its position remained unobserved and this requires that the wave function does indeed provide the most accurate account of the situation, before a position measurement is performed, but only as long as such a measurement is *not*, in effect, performed, because, under such conditions, the retarded and advanced portions of the propagation process might take place along different paths. However, if a position measurement is effected at some point (before an alternative measurement is performed on the same photon that would allow to reveal an interference between two different positions separated from the source by the same distance), then the photon can nevertheless be considered to have been constrained to follow a unique well-defined trajectory, all the way back to the emission process (when no condition seemed to apply on the wave function), as a result of this position measurement, given that in such a case the advanced portion of the process is allowed to enforce that condition of classical definiteness (imposed on the photon's trajectory by the future measurement), through the effect it exerts on the source, backward in time.

Now, given that what I'm proposing is a time-symmetric interpretation of quantum theory, it is important to mention that, in such a context, there must exist a time-reverse analog to ordinary quantum entanglement, which can be shown to actually give rise to non-local correlations arising from post selection (the phenomenon discussed in section 5.6 by which the choice of a measurement to be performed in the future is allowed to influence the evolution of a quantum system backward in time, as originally described in

the context of the two-state-vector formulation of quantum theory). Thus, even in the apparent absence of ordinary quantum entanglement established through local contact in the past, it should be possible to observe the existence of non-local correlations of the same type as arise in a more conventional context, when it is a future state that is entangled in a certain way as a result of post selection. From my viewpoint, the fact that those correlations do not allow faster-than-light communication can, once again, be explained as being a consequence of the fact that the constraint of global entanglement discussed in section 4.9 requires entropy to decrease in the past for all processes (occurring in a given universe), which means that no causal signal can propagate toward the future and then backward in time to a distant location, as a result of post selection, even if causality does operate both forward and backward in time at a more fundamental level.

But, despite what one might be tempted to believe, the possibility that such future entanglement may arise does not mean that every measurement result obtained at the present moment must be correlated to every other measurement result obtained at remote locations, as a result of post selection, because, contrarily to the situation that existed in the far past, most elementary particles present in the remote future are not in contact with one another at any point and therefore, even if those effects do exist, they should not be as commonplace as the effects arising from ordinary, past entanglement. In fact, the very significance of the cosmological constraint of global entanglement is that every particle of matter or radiation in the universe must have been entangled with at least one other particle, which was itself entangled with another particle, and so on, in the maximum-density state of the Big Bang. But no such a condition exists for the future (especially in the presence of negative-energy matter, for reasons I have explained in section 4.9) and given that we would be justified to expect that no low-gravitational-entropy Big Crunch will ever occur, then it would seem appropriate to conclude that future entanglement is not as essential a requirement for the universe, as the existence of the globally entangled state from which all matter emerged in the past.

It is now possible to reflect on the traditional positions, regarding the significance of EPR-type experiments. If we consider, first of all, the orthodox view and Bohr's position, it amounted to assume that only the wave function can be considered real, given that, when it is not the angle of impact that is detected for one or another of the two entangled photons, then it seems to be impossible to say anything meaningful about the trajectories of both

photons. The problem with this viewpoint is that, if the wave function is considered to be a real entity, then instantaneous action at a distance appears to be required, which is why the orthodox interpretation retreated into its idealistic position, according to which it simply doesn't make sense to speak about a reality behind observed phenomena (which may allow to avoid the conclusion that this reality is non-local). The position held by Einstein and the advocates of a realist approach was that this rejection of scientific realism is not acceptable and that quantum theory must simply be wrong or incomplete, given that it appears to require instantaneous action at a distance, when the consequence of a measurement on one of the two photons spreads to its entangled counterpart. Basically, then, one position required assuming that there is no reality, while the other required assuming that there is no entanglement. I believe that both positions were inappropriate in some way, but also accurate in a distinct way¹⁵. Clearly quantum theory and quantum entanglement are there to stay, but what I have tried to explain is that scientific realism is not optional either and can be accommodated without requiring instantaneous action at a distance, when time-symmetric causality is recognized to be an essential aspect of this physical reality and the appropriate (gravitational-entropy-reducing) constraint applies to the advanced portion of every quantum process.

5.10 The quantum measurement problem

It is usually recognized that the two main conceptual difficulties with which we are faced when trying to formulate a consistent interpretation of quantum theory are the existence of non-local correlations and the absence of an objective criterion for judging when it is that the multiple interfering potentialities characterizing the state of some unobserved dynamic attribute of a quantum system are actualized to a unique definite value, as happens when a

¹⁵It has been argued by some of the originators of the consistent-histories interpretation of quantum theory that Einstein was misguided in trying to uphold a certain requirement of scientific realism, to which the *conventional* interpretation of the theory does not conform, because it must be the theory that determines what is true of reality, even when it appears to require contradictory descriptions of it to be valid together at the same time. But I believe that it is rather this position which is misguided and this precisely because it constitutes an attempt at limiting what can be consistently described of reality in order to satisfy the perceived requirements of what is merely an inadequate or incomplete interpretation of the theory.

measurement of this attribute is performed. In the preceding sections I have offered a viable solution to the problem of quantum non-locality, in the context of a realist interpretation of quantum theory based on the requirement of time-symmetric causality. But while progress was achieved in the last few decades in identifying the conditions necessary for the decoherence process to occur, it remains that we haven't yet been able to determine exactly what is responsible for the persistence of quasiclassicality that is observed to characterize the evolution of quantum systems when they become entangled with their environment, following a measurement.

The currently favored approach for a solution to the quantum measurement problem plays a role that is much the same as the late nineteenth-century approaches which were adopted in an attempt to solve the problem of the origin of thermodynamic time asymmetry through the use of statistical methods. Indeed, at some point, Boltzmann thought that he had solved the problem of the origin of the arrow of time, because he had achieved significant progress in identifying its true origin. But as we now understand, it appears that he had not really provided a satisfactory explanation and that the remaining difficulties had not even been clearly identified. Today, it is widely believed that the problem of quantum measurement has been solved by the recent advances achieved in identifying the conditions necessary for the phenomenon of decoherence to occur, while in fact this is not entirely correct, precisely because the consequences of thermodynamic time asymmetry on the evolution of quantum systems haven't yet been properly assimilated. This is the problem I will attempt to circumscribe in this section and to which I will be able to provide a satisfactory solution in section 5.12. This will allow me to confirm, once again, that a realist approach, according to which there must exist a unique reality of some kind, independently from whether or not a system is being observed, is not incompatible with the empirical evidence that singles out quantum measurement as the necessary condition for the factual definiteness of reality.

Traditionally, the quantum measurement problem had to do with the difficulty we were experiencing in trying to identify the exact nature of the conditions that give rise to the actualization of quantum potentialities. In fact, the linearity of the equation that describes the evolution of the state vector made it difficult to understand how it could be that quantum interferences are, in effect, allowed to vanish for the dynamic attribute of a system under observation, so that they can give rise to a definite set of outcomes, to which meaningful probabilities can be ascribed. Thus, there appeared to

be a conflict between observations, which indicate that quantum potentialities are actualized to definite non-interfering outcomes following what we call a measurement, and the theory itself, which seems to require quantum superposition of states to persist indefinitely. From a conventional perspective, it would appear that when each possible outcome of a quantum process becomes correlated with one possible state of a measuring device, if the quantum system was in a state of superposition of the observable concerned at the time when this correlation was established, then the whole measuring device should also be found in a state of superposition following the moment at which the measurement took place.

One of the earliest attempts at solving this measurement problem became the actual justification for the conventional interpretation of quantum theory, according to which interferences arise because all possible histories occur all at once in the same universe as different ‘branches’ of history. What was proposed by Hugh Everett III is that there is no actualization process, but that the superposed macroscopic states of a measuring apparatus, which result from its correlation with the multiple interfering states of a quantum system, all exist simultaneously in parallel versions of history, while, for some reason, a ‘splitting’ occurs following measurement, which is responsible for the fact that the multiple branches of history no longer interfere with one another. The difficulty with this proposal, however, does not have to do only with the fact that it would involve logical contradictions in the context of a realist interpretation of quantum phenomena (a particle could be in one location as well as in another, in the very same portion of history), it also has to do with the fact that, if all branches are followed together, then there is no *a priori* reason why there could not be branches where a measuring device is in a state of superposition of macroscopic observables (in the context where decoherence alone would not be sufficient a requirement to allow one to expect that the absence of interferences that follows a measurement should persist indefinitely). But it is also contradictory to suggest that no actualization process takes place, while it is recognized that a splitting of branches is required to eliminate interferences, following a measurement.

Nevertheless, the idea endured and was later revived when it was discovered that, under specific conditions, the phenomenon of decoherence must give rise to a diagonalization of the reduced density operator in the basis of the attribute under measurement (the probabilities of occurrence for superposed results must rapidly decrease to negligible values upon measurement of a given observable), which would appear to legitimate the splitting branches

hypothesis. But if there should be no doubt that the discovery of decoherence itself was a step in the right direction, this does not mean that the hypothesis that there may exist many continuously ‘splitting’ branches of history in the very same universe has been confirmed. Given that decoherence does not require the existence of those multiple branches of history, it appears that the only advantage that Everett’s original proposal *might* provide does not really have to do with solving the quantum measurement problem, but with allowing one to avoid having to explain the uniqueness of measurement results in the context where the formalism of the current theory does not allow one single state to be attributed preferred status as *the* actual outcome of a measurement on a previously interfering physical attribute.

It is usually recognized, in effect, that the only legitimate purpose of the multiple-branches hypothesis would be to allow one to avoid having to postulate the existence of distinct dynamical laws that would apply only during processes that can be qualified as measurements, given that, if it was possible to assume that all histories do occur all at once, following a measurement, then it would no longer be necessary to explain why it is that one unique measurement result appears to be actualized among the many different potentialities. Thus, it is argued that when all the superposed states of some physical attribute are assumed to be actualized together in different splitting ‘branches’ of the same history, there no longer needs to be a cause (of unknown origin) that would give rise to the one particular outcome that is actually observed following measurement.

But given that I have already argued, based on more general considerations, that it is not really necessary, in order to explain the existence of quantum interferences, to assume that all possible histories are followed all at once in the course of each and every quantum process, then it would appear that it is preferable to recognize that the unique reality we do observe during measurements is a reflection of the uniqueness of the non-classical (time-symmetric) reality that exists in between measurements, instead of trying to argue that there must be a multiplicity of measurement results, which we do not observe, that would correspond with a multiplicity of histories, which we cannot observe either, but that would need to exist in between measurements. Thus, what must be clear is that the uniqueness of measurement results is not less, but rather more problematic when one assumes that all trajectories are followed all at once while a physical attribute is not subject to measurement, which is a hypothesis that is actually necessary only in the context of the tentative solution that is provided to the quantum mea-

surement problem by this many-worlds interpretation itself. But what's even more significant is that, as I will explain below, it appears that decoherence is not sufficient, all by itself, to predict that a physical attribute which has once been measured remains in a definite quasiclassical state and does not give rise to interferences involving macroscopic attributes of a measuring device. Therefore, it seems that we should still expect that in some of those hypothetical branches, macroscopic state superpositions would develop at some point.

What I find most difficult to accept regarding the many-worlds interpretation, however, is the fact that we are required to believe that the unique character of reality that we do observe on a classical scale is just an illusion, while we are also expected to assume that the hypothesis of a multiplicity of coexisting branches of history, which has never been directly confirmed by any observation, is valid under all circumstances. In other words, we are required to assume that what we see is not the true reality, while what is a mere hypothesis, that cannot yet be observationally confirmed, must be considered true, even though it is clearly incompatible with what we do know about reality. It must be clear that it is not possible to assume that the existence of interferences between the multiple paths available to a quantum system simply means that in between measurements a system goes through one path in one universe and through another path in another universe, because if that was the case then we should not, in fact, observe interferences in any one particular universe, given that it is an essential consistency requirement to assume that universes are completely independent from the viewpoint of causality (hence it would be a mistake to think that the expression 'many worlds' is synonymous with 'many universes').

But, it is also difficult to explain how it could be that an observer present in one of the multiple branches of history would not be allowed to perceive what happens in the other branches, while those branches would be allowed to interfere with one another, at least prior to measurement, thereby implying that they actually exist all at once in the very same universe. Here, again, a lot of silly things have been said to try to justify how it can be that those two requirements do not contradict one another, but in the end, one must recognize that this constitutes a basic inconsistency of the many-worlds interpretation of quantum theory that invalidates it as a solution to the problem of quantum measurement. It is merely because this objection is so simple that it has avoided the attention of the most knowledgeable experts, who usually prefer to concentrate their efforts on more complex and more

challenging issues.

Those criticisms, however, must not be understood to mean that the hypothesis of a multiplicity of universes existing *independently* from one another is wrong, because, in fact, there may be good reasons to recognize the validity of this clearly distinct hypothesis (which is not dependent on the validity of the many-worlds interpretation of quantum theory), in the context where the weak anthropic principle appears to constitute the only possible explanation for certain otherwise unlikely properties of our universe. Thus, while it may not be possible to reject the hypothesis that an infinity of causally independent universes exist in parallel, it must be clear that the idea that many interfering branches exist in the *same* portion of a universe's history is a distinct hypothesis, which is certainly not as unavoidable. But if a multitude of realities are to be allowed to interfere with one another so as to explain quantum state superposition, then they must definitely be present all at once in the same portion of the universe's history and therefore cannot constitute different universes, as is often suggested.

Personally, I always felt that the whole idea that there may exist multiple parallel branches of history in the same universe, but that it is only when those alternative branches should become observable (following the measurement of a physical attribute initially in a state of quantum superposition) that they actually 'split' and become totally independent (as a result of decoherence), therefore precluding a confirmation of their existence, has all the characteristics of a conspiracy theory and this only reinforces my conviction that the many-worlds interpretation is not good science. A truly appropriate solution to the problem of quantum measurement would then need to be based on the hypothesis that reality is unique, in the particular way proposed above, even in between measurements, which is the only way one could avoid having to appeal to the problematic splitting-branches hypothesis in order to explain the uniqueness of measurement results.

I'm aware, though, that it has been argued by Heinz Dieter Zeh [53] that the multiple-branches hypothesis may be unavoidable if one does not want to have to modify quantum theory, because this hypothesis allows all possibilities to be actualized all at once as different branches, which is the only way one can avoid having to assume that a unique state of such an unobserved attribute existed before decoherence took place, that would merely have been revealed by the measurement. Indeed, it is known that, for various reasons, a quantum measurement cannot be considered to simply consist in acquiring knowledge about the unique, pre-existing state of an unobserved attribute.

But I have explained, in section 5.8, that, if it is impossible to assume that a unique reality existed before a measurement was performed on some unobserved attribute of a quantum system, it is only because we usually assume a unique reality to be unique in the classical sense, while in fact the unique reality that would characterize a quantum process (in the absence of measurement on a certain dynamic attribute) is of a time-symmetric nature and involves both a unique retarded state and a unique and possibly different advanced state at all times, which guarantees the consistency of past evolution with *any* future measurement and which requires all possible intermediary states to contribute to the final probability amplitude, so that the future measurement does not allow to reveal a unique *classical* path through which the system would have propagated¹⁶.

It is therefore simply the fact that, under such circumstances, the future measurement also exerts an influence on the past state preceding it (as when a system is submitted to post selection) that forbids one from assuming that a unique classical state existed in the past, independent from what measurement is performed in the future. But it must be clear, once again, that this does not prevent a unique state from having actually existed at all times in the retarded and advanced portions of history and therefore the conclusion that the unique character of measurement results can only be explained by postulating that all histories are followed all at once (in the same universe) cannot be considered valid. In any case, if reality was not of the unique time-symmetric type and the decoherent branches hypothesis was assumed to alone provide a solution to the quantum measurement problem, then an alternative explanation of quantum non-locality would have to be found, as it cannot be provided by Everett's interpretation and this is an additional difficulty for the conventional approach. Thus, I think that I have explained with enough clarity why it is that the frequently stated conclusion that it is just as difficult to provide decisive arguments in favor of the many-worlds

¹⁶Contrarily to what Zeh suggests in another publication [54], the fact that there would exist a unique, but unknown state prior to measurement of an unobserved attribute does not violate the condition imposed by the Von Neumann equation (the quantum-mechanical generalization of Liouville's equation) that ensemble entropy should not decrease during measurement, because this conclusion would only be valid based on the hypothesis that the unknown, but definite state would actually be a classical (hidden-variables) state, while, in a time-symmetric context, when information about a measured attribute is obtained, information concerning its conjugate counterpart is lost, which allows information to be conserved.

interpretation, as it is to provide arguments that would invalidate the idea, is not well founded, because the hypothetical, multiple branches of history are not necessary, or even adequate to explain quantum strangeness, while they also do not appear to be required to solve the quantum measurement problem (especially in the context where it is recognized that the splitting process would not, all by itself, allow one to avoid the difficulty associated with the non-local aspect of state-vector reduction).

Now, some theoreticians are worried about the fact that decoherence would seem to be insufficient to solve the problem of the actualization of quantum potentialities, when we are considering the system under observation to be the universe as a whole (as becomes necessary in a quantum cosmological context). The problem they see is that, in such a case, there would be no outside environment degrees of freedom to effect decoherence, while this is known to be a requirement under ordinary conditions. What constitutes a more serious difficulty, however, is that, from the viewpoint of the currently favored approach (the consistent-histories interpretation of quantum theory), decoherence is insufficient to explain the persistence of quasiclassicality, not just in the cosmological case, but even under more general circumstances and on a much smaller scale, as was first pointed out by Fay Dowker and Adrian Kent [55]. But this is not just a consequence of the fact that (ignoring my own contribution) we do not yet have a valid explanation for the irreversibility that characterizes the processes which give rise to decoherence, it rather appears to be a basic insufficiency of the current approach, which does not allow to predict that classical behavior would persist following decoherence, even when irreversibility is assumed to characterize the evolution of the environment degrees of freedom without explanation, due to some boundary condition of low entropy that presumably apply to the initial state of the universe at the Big Bang (I will return to this question below). An appropriate solution to the quantum measurement problem, therefore, must allow to predict the emergence of quasiclassicality, not just on the cosmological scale, but also on the much smaller scale of measuring devices, where the conventional approach is insufficient as well.

In any case, it is my intention to demonstrate that it is not necessary, in order to explain the nature of the outcomes of quantum measurement, to postulate the existence of distinct (perhaps fundamentally irreversible) evolution laws that would apply only during a process that could, in effect, be characterized as a measurement. Thus, while I do agree with the most knowledgeable authors that quantum theory, as it is currently interpreted,

fails to explain the persistence of quasiclassicality that is observed to follow any measurement, I do not believe that what is required in order to address this difficulty is a modification of the basic mathematical framework of the theory that would need to apply whenever measurements are performed, as was once proposed. We cannot reject a requirement like that of time symmetry, whose value is indisputable, to seek a solution in terms of fundamentally irreversible physical laws, when there is no evidence that such a choice is absolutely essential for a solution to the problem of the emergence and the persistence of quasiclassicality. I still believe that it is at the level of interpretation that the appropriate solution will emerge that will allow us to solve the remaining difficulties surrounding quantum measurement. As I will explain in section 5.12, what the current theory needs is not so much a modification of its structure, as an extension of its meaning.

It must be recognized, however, that the distinctive feature of all processes that can be characterized as giving rise to a measurement is indeed irreversibility. A quantum measurement is nothing but the entanglement of a particular state of some attribute of a quantum system with some distinguishable macroscopic property of its environment whose future evolution is irreversibly influenced by this particular event. The fact that no quantum interference is ever observed for irreversibly evolving systems indicates that the non-superposed nature of measurement results is related to the irreversible character of the measurement process. Thus, decoherence itself (literally the loss of phase coherence) can only occur when a microscopic (quantum mechanically evolving) system becomes entangled with some irreversibly evolving (entropy increasing) processes taking place in its environment (usually involving dissipation), so that the phase relations that could have given rise to interferences become delocalized and are assumed to no longer be accessible to observation, as is already well understood. In fact, the ultimate manifestation of irreversibility appears to be decoherence itself.

It should not be unexpected, therefore, that all measurements involve the formation of a record, given that for a record of some past event to form, entropy must be growing in the future. Indeed, the formation of a record merely consists in the production of multiple persistent and somewhat independent effects in the future, which all emerge as the outcomes of one single identifiable cause in the past and this is undoubtedly a process that is asymmetric with respect to the direction of time. What this means is that there is something very tangible occurring when a quantum measurement is performed and therefore, if it is true that our knowledge of a quantum

system changes when quantum potentialities are actualized, it would not be appropriate to assume that the changes which are taking place in the course of a measurement are merely subjective, because, following measurement, the observed attribute is no longer unique merely in a time-symmetric quantum way, but acquires the same unique value in both the retarded and the advanced portions of history.

It is not difficult, in effect, to show that irreversibility is essential for a measurement to occur, while the mere complexity, or the large number of independent degrees of freedom of a macroscopic system with which a quantum system may become entangled, alone, is not sufficient a condition for triggering a measurement. Indeed, it is apparent in the formalism of quantum field theory that there is a near infinite amount of structure that must be taken into account in estimating the probability amplitude of any process, as is apparent in the fact that additional fermion loops and radiative correction terms arise at every level of approximation on shorter scales. If we were to consider this small-scale complexity to provide the conditions for quantum measurement to take place, then it should be the case that the world would be quasiclassical down to a much smaller scale, given that all the complexity that is present at higher energies (and which can only be ignored as a result of the validity of the renormalization procedure) would allow measurement to take place long before a quantum system even has the chance to become entangled with a macroscopic system¹⁷. When no irreversible change that could potentially carry information about a former state is allowed to take place, the predictions of quantum theory do not merely apply as much for the future as for the past, they do not apply at all, because there can be no measurement, that is to say, no irreversible process of amplification of alternative microscopic states. From those considerations one can only conclude that the defining characteristic of the processes that allow quantum measurement to happen is not merely their complexity, but really their irreversibility.

This asymmetry must not be confused with that which also characterizes

¹⁷I have provided strong arguments, in section 4.7, to the effect that in the absence of matter there can be no persistent microscopic structure in the distribution of vacuum energy and this means that no record of what takes place in the vacuum on smaller scales can exist, unless we directly reveal the existence of those processes by entangling them with irreversibly evolving degrees of freedom which leave persistent traces of their occurrence and this allows to confirm that the complexity of virtual processes cannot be considered sufficient as a condition for quantum measurement to take place.

the otherwise time-reversible ‘unitary’ evolution that takes place in between measurements and which is made conspicuous by the fact that the predictions of quantum theory are only valid for future evolution. Indeed, the impossibility to accurately ‘predict’ the past arises as a consequence of the fact that only a subset of states can be actualized in the past, due to the constraint of diminishing entropy that exists for this direction of time and which also applies to classical evolution. It is the fact that no such a constraint applies on future evolution that allows predictions of future transition probabilities to be valid, while predictions of transition to past states do not apply in general. In section 4.9 I have explained that this constraint arises from the requirement that there exist relations of causality between all particles present in the expanding universe, which in the presence of negative-energy matter implies that the initial state at the Big Bang was characterized by a condition of minimum gravitational entropy from which all later irreversibility follows. But while the time asymmetry that characterizes all measurement processes has the same origin, it is a distinct phenomenon, that usually operates on a much shorter time scale and that does, in effect, give rise to a reduction of the state vector. Yet, it is appropriate to remark that it is the global entanglement constraint unveiled in chapter 4 that actually explains the fact that decoherence is allowed to occur, which is necessary (even though not entirely sufficient) to explain the persistence of the quasiclassical nature of history that follows quantum measurements. In fact, this is the only explanation of time asymmetry that allows to deduce (rather than merely assume) that decoherence always occurs in one and the same direction of time for all measurement processes (as is required for the logical consistency of history according to Roland Omnès [56] (p. 237)), as decoherence itself does not *a priori* favor one direction of time over the other.

But when a measurement is performed on a dynamic physical attribute of a quantum system in an interfering state, what must also happen is a variation of the macroscopic constraints which apply on the wave function that describes the evolution of the system. When this is allowed to occur and the state vector is reduced, an irreversible change is introduced in the evolution of the system itself and the outcome of this evolution is, in general, unpredictable (even if the wave function itself always evolves deterministically in between measurements). But in the context of a realist interpretation of quantum theory, this cannot be understood to mean that it is the evolution that takes place in the course of a measurement which is alone responsible for giving rise to the unpredictability of quantum phenomena, as is sometimes

proposed. If the state of some interfering dynamic attribute of a quantum system is unique, in the time-symmetric sense, before a measurement is performed, as I previously argued one must recognize, then it certainly cannot be assumed that the randomness of its evolution is merely a consequence of the events that take place during the subsequent measurement and it becomes necessary to admit that it is the unobserved paths followed forward and backward in time by the system as it approaches or emerges from the event at which a measurement is performed which are randomly determined and that this is what explains the unpredictability of the outcome of this measurement. It must be clear, in any case, that the randomness of quantum processes, like their uniqueness, is not an illusion that emerges from the fact that an observer may be unable to perceive the evolution that supposedly takes place all at once in multiple branches of history that would exist together in one single universe, as is sometimes suggested in the context of a many-worlds interpretation. Randomness is a fact of the reality we experience that becomes perfectly acceptable in the context of a time-symmetric formulation of quantum theory, where the deterministically-evolving wave function is not reality itself and there exists a unique history of some kind, even in between measurements.

Thus, if randomness appears to take place only during measurements, it is simply because it is only as a result of processes which can be characterized as measurements that the uniqueness of reality (in the time-symmetric sense) is made apparent, while it is only at the level of individual histories that reality may be observed to vary unpredictably (given that the wave function itself evolves deterministically). But quantum evolution must be understood to always be random, even though in the absence of measurement, or when the macroscopic constraints applying on a system remain unchanged, this unpredictability is not apparent, because it has no observable consequences. Once again, therefore, it seems that it is incorrect to assume that a fundamental distinction must exist between the ‘unitary’ evolution that takes place in between measurements and the evolution that characterizes a process during which quantum potentialities are actualized and this means that it should be possible to explain the quasiclassical nature of the evolution that follows a quantum measurement while remaining within the confines of the current mathematical framework of quantum theory. The difference between observed and unobserved evolution is real, but only because the conditions that exist when there is an absence of knowledge provide a quantum system with more freedom regarding what it is allowed to do as it randomly evolves

in the retarded and advanced portions of history.

It also transpires that the standard account, regarding the distinction between those situations in which a measurement takes place and those in which the usual ‘unitary’ evolution law applies, is somewhat misleading, because, in fact, a quantum system is always in a state where at least one dynamic attribute (as unnatural as it may be) is in a classically well-defined state, even though this means that the conjugate attribute is completely undetermined. This is a very important fact that is often overlooked and which actually holds the key to a solution to the remaining issues that prevent the formulation of a satisfactory explanation of the persistence of quasiclassicality. When a measurement is performed, all that really happens is that the state of a system changes in such a way that an attribute (say position) which was in a state of quantum superposition the moment before, becomes classically well-defined the moment after, while its conjugate attribute (say momentum), which was classically well-defined initially, actually becomes quantum-mechanically superposed. In such a context, it would certainly be inappropriate to argue that a *fundamental* change occurs in the course of a process that can be qualified as a measurement, even though it is clear that some constraint, not present before the process took place, does, in effect, become significant for the future evolution of the attribute of the system which is subjected to measurement (I will have more to say concerning this issue in section 5.12).

Now, the modern formulation of quantum theory, the one which can most naturally accommodate the decoherence process, is usually considered to be that of consistent histories, which was developed in three steps by Robert Griffiths [57], Roland Omnès [58], and Murray Gell-Mann and James Hartle [59]. From this formalism emerges an interpretation according to which it is merely the fact that one may choose to ignore certain aspects of reality, and submit them to a summation process, that allows one to obtain meaningful probabilities (which are positive and which add up to one) for the possible histories of a quantum system which has become entangled with the summed-over portion of reality. It would then merely be the fact that one may choose to ignore what goes on in the environment with which a system has become entangled that would allow one to find the system to be in a mixed quantum state, instead of a pure quantum state for which interferences would be observed, following measurement. More specifically, what the formalism of consistent histories provides is a criterion for judging when it is that sufficiently coarse-grained histories are obtained (by ignoring certain

details of the historical description of reality), which do not interfere with one another and which can therefore be attributed meaningful probabilities. Interestingly, the manner by which this is achieved is by considering pairs of coarse-grained histories (consisting of sets of alternative fine-grained histories whose ignored details are allowed to differ in any possible way) which are subjected to decoherence and between which there are virtually no interferences. When those conditions are satisfied, a meaningful probability for the process so described to occur can be obtained by applying the usual rule, which consists in multiplying the probability amplitude for a history with the complex conjugate of the amplitude for the same coarse-grained history. But no interpretation is given for why it might be necessary to consider pairs of coarse-grained histories rather than single histories, even though this appears to be required from a mathematical viewpoint.

I believe that the formalism of consistent histories must be considered an essential element of a fully satisfactory interpretation of quantum theory, even if merely because it constitutes the basis of the only solution to the quantum measurement problem that would also apply on a cosmological scale, where no *external* environment degrees of freedom exist which, according to a more conventional theory, would be required to give rise to a measurement. It is incorrect, therefore, to argue that the state of the universe cannot decohere because no environment exists outside the universe, because if decoherence is an outcome of temporal irreversibility, then there is enough opportunity for decoherence to occur *on a much smaller scale*. Indeed, what the formalism of consistent histories allows is a more appropriate definition of quantum measurement as taking place continuously over the entire duration of a process, rather than at one particular event. This becomes possible as long as the *local* environment degrees of freedom which are left out of the description of the process evolve irreversibly, thereby allowing decoherence to arise. One of the advantages of such a viewpoint is that it is easier to see how it can be that the simple possibility for an event to happen allows a measurement to be performed, even if this event does not happen (as in the case of interaction-free measurements), because when something is, in effect, allowed to happen we simply are in a situation where one specific set of macroscopic experimental constraints exists throughout the duration of a process, which would not exist otherwise, while different constraints mean a different measurement, not an absence of measurement.

But while the consistent-histories approach is certainly well-founded all by itself, given that it allows one to avoid having to refer to classical ob-

servers and classical measuring devices that would not be describable using the formalism of quantum theory, it appears to be insufficient to predict the emergence of a classical world (a maximum quasiclassical domain). It is as if decoherence alone was not enough constraining a condition to guarantee an absence of quantum interferences between all the coarse-grained histories to which it may give rise, while no criterion currently exists to select as physically relevant only those future histories which actually describe a quasiclassical evolution. As was the case with the original many-worlds interpretation of quantum theory, it is not possible to avoid the conclusion that, in the course of certain otherwise ‘consistent’ histories, a macroscopic measuring device may end up in a superposition of states, after becoming entangled with a quantum system.

There are, then, two problems affecting the consistent histories interpretation of quantum theory. The first problem one must face has to do with the previously discussed lack of motive for justifying the application of the criterion of ‘consistency’ that is attributed to families of coarse-grained histories and according to which certain histories would simply be meaningless, given that classically meaningful probabilities cannot be assigned to them. I have already mentioned that it appears preferable to allow our conception of reality to adapt to the fact that classical probability theory does not always apply, instead of trying to limit what may be consistently described of this reality through some arbitrary criterion that only serves to accommodate the limitations and the inadequacies of an interpretation that cannot fully satisfy the requirement of a realist, time-symmetric description of reality. Thus, I believe that it is important not to commit the error of enforcing consistency at the price of rejecting a realist interpretation of facts, which would simply contribute to perpetuate the difficulties which are known to affect the original Copenhagen interpretation.

What should be recognized as nonsense is not the hypothesis that a photon follows a unique but unobservable trajectory of some kind in between measurements, but the decree that we should not even try to describe reality in situations where we do not yet know how to make sense of it. This reflection is especially relevant given that, in a quantum mechanical context, we are always dealing with probabilistic inferences, so that even histories which we may expect to be ‘consistent’ might in some rare circumstances turn out to be ‘nonsense’, which is certainly indicative of the arbitrariness of the restrictions imposed by the consistent-histories interpretation of quantum theory on our concept of reality. Therefore, to achieve further progress

regarding the issue of quantum measurement, one must first realize that in face of the experimental evidence from which quantum theory emerged, the desire to restrict the application of the criterion of logical consistency to aspects of reality which behave in conformity with classical expectations is just as irrational as the desire to uphold determinism, that is to say, the predictability of future evolution.

The additional issue we need to consider, however, is more pragmatic. It has to do with the fact that, in the absence of a stronger and more specific constraint, there would be histories which could be characterized as ‘consistent’ by the formalism of the theory, but which would not remain quasiclassical as time goes, following decoherence. This is the problem discussed in Ref. [55] and which I have mentioned earlier in this section. As Dowker and Kent explain, predictions only become possible, within the formalism of consistent histories, once a set of histories, the physically relevant set, which is based on a specific choice of dynamic attributes and a particular choice of coarse-graining, has been selected, whose elements can then be attributed meaningful probabilities. But in a quantum-mechanical context, there appears to be total freedom over the choice of which dynamic attributes are used to specify the exact state of our physical systems and what elements of reality can be ignored and summed-over, and this is where the problem originates, because when no criterion exists to limit those choices, most ‘consistent’ histories do not remain quasiclassical in the future, even if they were so characterized in the past. Thus, the criterion of ‘consistency’ appears to be insufficient to predict the persistence of the quasiclassical nature of history. In fact, it seems that the condition would not even allow one to assume that the past itself must have been classical up to the present moment, despite the fact that the existence of mutually consistent records of a unique past appears to indicate that the whole observable universe evolved classically (without large-scale quantum interferences) as far back in time as one can tell. What remains problematic with the current approach, therefore, is the absence of a criterion, within the interpretation itself, for choosing the appropriate, physically relevant set which would allow to describe the quasiclassical world we do experience.

What was originally proposed by Murray Gell-Mann and James Hartle is that, if we perceive a quasiclassical world it is because, as observers, we have evolved to take advantage of only those formulations of history according to which the world does, in effect, remain quasiclassical. The problem is that it appears that in the absence of a criterion for justifying the selection of the

appropriate, physically relevant set of histories, the above mentioned results imply that the most likely explanation for the fact that one experiences a quasiclassical world would require one to reject all evidence of past quasiclassicality and all expectations of future quasiclassicality as being mere illusions and to satisfy oneself with having ‘explained’ why it is that, at the present moment, one goes through a classical experience that is such as to give one the impression of living in a world that remained quasiclassical on a global scale during most of its history, even though that would not be the case. But I have already explained why such solipsistic explanations, which require one to assume that one’s current state of awareness is all that truly exists (or that evolves classically), are not acceptable in general, from the viewpoint of scientific realism, and if there is one situation where this criticism would definitely need to apply it is certainly here.

It seems to me that if such an approach is still often considered to constitute a valid explanation of the quasiclassical character of reality, it is merely because we cannot see how the remaining issues facing the current state-of-the-art interpretation of quantum theory could be resolved, so that we have come to believe that the solution may be that there is no problem after all, as long as we consider the world in the ‘appropriate’ way. But, if we really want to explain something, then, clearly, we must identify the constraint that allows to select the physically relevant set of histories in which quasiclassicality is experienced by *all* observers, because the only alternative would be to retreat into a paranoid vision of reality, where all that exists (in the classical sense) is the *impression* of a persistent, large-scale, quasiclassical reality, despite the fact that there would be absolutely no reason for why such a deceptive state of consciousness should be experienced (which is the real problem). I believe that what those difficulties illustrate is the incorrectness of the basic assumption that no logically consistent interpretation exists for the interfering fine-grained histories which actually constitute the most fundamental elements of the consistent-histories formulation of quantum theory.

It is significant, in this context, that certain specialists have proposed a weaker and more general form of consistency conditions [60] that merely amounts to impose that the probabilities of coarse-grained histories be positive, while still satisfying the usual probability sum rules. Those generalized ‘consistency’ conditions result in a formalism that is time-reversal invariant (which from my viewpoint is certainly a desirable property) and which selects sets of histories called *linearly positive histories* that include consistent histories as a subset of possibilities. Once this is recognized to be a viable

approach, however, one may be tempted to go one step further and simply allow negative probabilities as well, by considering the most complete sets of histories that would include all sets of linearly positive histories as a subset. If such an even more complete generalization was never considered viable it is obviously due to the fact that negative probabilities cannot be classically interpreted (in such a context) and therefore appear meaningless and undesirable. Yet, Robert Griffiths, suggested that it might be desirable to try to provide an interpretation of the probabilities which are known to arise when we consider histories that do not satisfy the ‘consistency’ criterion. Dowker and Kent themselves insist that there would be no logical contradiction in using an ‘inconsistent’ set of histories if a criterion existed that would allow one to select from it the physically relevant set and it was found that it allows a logically consistent description of historical facts on a sufficiently ‘large’ scale.

The problem is that, in the current context, the ‘consistency’ criterion appears to be necessary for selecting sets of coarse-grained histories that do not interfere with one another, as required by observations, while no satisfactory interpretation exists for negative probabilities. But in the context where we still need to identify the constraint that allows one to choose the physically relevant set, it cannot be ruled out that it might be this condition which enables to generate a historical description of reality that naturally satisfies both the criterion of ‘consistency’ and that of persistent quasiclassicality. I have already suggested that, in the more appropriate context of a realist, time-symmetric interpretation of quantum theory, logical consistency (in the general sense) would rather need to be satisfied by the unique retarded and advanced portions of history. But I also explained that, from such a perspective, an adequate interpretation of negative probabilities can be formulated that would confine them to unobservable aspects of physical processes. Thus, if a criterion can be found for the selection of a set of histories that is not *a priori* ‘consistent’, but that would nevertheless allow the quasiclassical character of reality to naturally emerge on the appropriate scale, then we may finally obtain a satisfactory extension of the current formalism that would allow to solve the quantum measurement problem.

In fact, we may have another motive for recognizing that an additional constraint is necessary to explain the quasiclassical nature of reality that is observed on a sufficiently irreversible scale. It was pointed out by Roger Penrose and apparently also by John Bell and Bernard d’Espagnat that the current explanation for the reduction of the state vector through decoher-

ence is dependent on the hypothesis that it is impossible, in effect, to reveal the existence of quantum interferences involving the detailed configuration of the degrees of freedom of that part of the environment which has become entangled with a quantum system. But there is presently no valid reason to assume that such an unlikely procedure could not be carried out at some point in the future (even without deliberate intervention) and this means that the current explanation for the disappearance of quantum interferences following measurement is only valid based on the assumption that the practical limitations that may prevent the observation of interferences between macroscopic states will *never* be overturned.

Given that the existence of practical limitations to unveil superpositions of macroscopic states through a manipulation of the delocalized environment degrees of freedom has been shown by Roland Omnès to be necessary for the validity of the factual definiteness of reality and the applicability of the conventional rules of logic, it is certainly significant that Omnès himself has argued that one cannot definitely rule out the possibility that such an unlikely evolution could happen, but that given that it would mean that the world would no longer be ‘consistent’, then he prefers to simply assume that the low probabilities involved imply that the decoherence process is definitive in principle. In the context of a conventional many-worlds interpretation, we would certainly be justified to assume that this condition needs to be fulfilled, as if it was not the case, then we should actually observe the multiple branches of history to interfere among themselves, even on the macroscopic level of measuring apparatuses, which does not only constitute an additional difficulty for this particular interpretation of quantum theory, but which also illustrates the necessity of providing a satisfactory explanation for the absolute irreversibility of the decoherence process.

Of course, we do observe an absence of interferences between alternative coarse-grained histories past a certain level of irreversibility of the ignored (summed-over) portions of a process¹⁸ and this may appear to confirm the

¹⁸It was once suggested that quantum interferences between alternative states are actually always allowed to occur, regardless of the size of the system under observation or the degree of irreversibility of its evolution, but that, if the existence of such interferences can be ignored, it is simply because they would be too difficult to reveal in the case of macroscopic systems. But it is usually recognized that this is not a valid proposal, because, in fact, nothing would be easier to distinguish than interferences between two different states of a pointer on a measuring device, given that this would necessarily be apparent in the statistical distribution of subsequent measurement results.

validity of the assumption that the practical limitations discussed here cannot be overcome. But we must recognize that we have, at present, no reason, from a theoretical viewpoint, to assume that such an unlikely reversal of fortune could not happen at some point in the future, because, even if there is only an infinitesimal chance that it does, given an infinite amount of time it should eventually happen and in such a case the consequences would be felt right now (this is made unavoidable in the context where the time-symmetric nature of quantum evolution allows future measurements to exert an effect on past evolution). Even if such a phenomenon was to occur only once on a large scale, it would be possible to observe its consequences, because the usual assumption to the effect that there is no state superposition following measurement would then no longer allow our prediction of transition probabilities to match observations, therefore indicating that the conventional hypothesis is incorrect. The fact that we usually do not observe such a disagreement means that the assumption that the decoherence process is in general truly irreversible is appropriate, even if it is not, at present, entirely justified. A satisfactory solution to the problem of quantum measurement should therefore allow one to gain confidence that, once decoherence has occurred, there is no chance that it may somehow be overturned at *any* time in the future, which would allow to justify attributing the status of established facts to measurement results.

In any case, the often encountered statement to the effect that quantum theory has never been proven wrong, which would seem to invalidate the claim that the currently favored interpretation is incomplete, can no longer be considered accurate, given that, in the context of the developments discussed above, it seems that what the theory predicts is an absence of quasiclassicality in both the future and the past and this is clearly in conflict with what we do observe (for the past) and with what we have very good reasons to expect to observe (for the future). Therefore, a solution to the quantum measurement problem, the central problem of the interpretation of quantum theory, cannot merely consist in assuming that elementary particles acquire reality as a consequence of interaction with another part of reality (presumably a measuring device), as was originally proposed by some of the founders of quantum mechanics and as is still considered appropriate by advocates of the relational interpretation of quantum theory.

What I have tried to explain in this and the preceding sections of this chapter is that it is not necessary and not appropriate, or even possible to assume that no unique reality of some kind exists, for a quantum system, in

between interactions with a measuring device. The difficulty to explain the emergence of quasiclassicality cannot be considered to mean that the theory only allows to describe how quantum systems interact with the rest of the world, as if this was a requirement of a relational description of reality. In fact, as I will soon explain, it rather appears that a satisfactory solution to the quantum measurement problem actually requires considering that a well-defined and, in some way, unique, but unobservable reality does exist between measurements. Particles do not become real through interactions, and the uniqueness of reality, which is observed during measurements, is not an effect that propagates as a result of further interaction, because, even if that was considered to be true, the emergence of quasiclassicality would remain unexplained. It is not our intuition that such an explanation of quantum strangeness must be wrong that is at fault, but rather the orthodox interpretation of the theory and the insistence that we should not attempt to describe reality when it is not observed.

What emerges from those considerations is that, as undesirable as it may once have appeared, there seems to be something unavoidable with John Von Neumann's conclusion that something essential (although not necessarily fundamental) must differentiate a quantum system from the measuring apparatus and observer who effect a measurement on this system. Unless we are to allow for grossly inaccurate predictions, it is necessary to explain what justifies this distinction. But even though this difference can be recognized to have something to do with time irreversibility and even though it must come into effect following decoherence, its exact nature remains unidentified from the viewpoint of all known interpretations. What explains that Von Neumann's conclusion was never taken seriously is certainly his early proposal that the dividing line between superposed system and observing system may be determined by the level at which consciousness occurs, which could perhaps explain why it is that human observers never experience quantum interferences. Indeed, any reference to such qualitative aspects of physical reality as a degree of consciousness, or a level of cerebral development as possible causes of state vector reduction is properly viewed with extreme suspicion by any physicist with a minimum level of cerebral development, while, in fact, such a reference is not necessary for the validity of Von Neumann's conclusion. Once again, a perfectly valid deduction was ignored as a consequence of being associated with questionable assumptions which are not essential to its validity. But, if this is the truth, then it remains to identify the nature of this distinguishing property and to explain why it has the

decisive consequences it is observed to have, in the context where the basic mathematical framework of quantum theory is assumed to be valid under all circumstances. This is the task I will try to accomplish once I have clarified the role played by time in the most fundamental of quantum-mechanical frameworks.

5.11 The emergence of time in quantum cosmology

When searching for an adequate solution to the quantum measurement problem and a plausible explanation for the emergence of our quasiclassical world, what one must first decide is whether quantum theory needs to be replaced by a better theory, or whether the current framework is appropriate to deal with those apparently insoluble difficulties. What I have been led to conclude is that quantum theory is indeed incomplete and that it must be supplemented with new conceptual elements if it is to be made fully consistent with what we already know of physical reality that currently appears to conflict with its predictions. But, as I already mentioned, this does not mean that the current mathematical framework of quantum theory (in its most appropriate form) must be rejected, or that the progress which was already achieved while we were trying to develop a better interpretation of the theory has become useless. It is, in effect, by building on earlier developments towards a time-symmetric formulation of quantum theory that I will be able to address the remaining difficulties affecting the consistent-histories interpretation and to finally explain the quasiclassical nature of reality. For that purpose, however, it is necessary to first examine the extent to which time itself can still be assumed to constitute a meaningful concept in quantum cosmology and to explain how it is allowed to emerge from a fundamental theory in which it may only be present in embryonic form. This has been made unavoidable by certain developments that took place in the field of quantum gravitation, which appear to imply that the notion of a universal time variable may no longer be relevant to a fundamental description of reality, whether on the Planck scale or on the cosmological scale.

Even though it was originally suggested that time may be irrelevant to cosmology only when the first tentative quantum-mechanical descriptions of the universe as a whole were introduced, the perceived difficulty is actually

also present in classical cosmology. Indeed, it appears desirable, from both a practical and a theoretical viewpoint, to formulate relativity as a dynamical theory that would describe the evolution in time of the curvature of three-dimensional space, given that such an approach can be more easily extended to a background-independent quantum-mechanical theory. But in a general-relativistic context, when it is recognized that all the meaningful physical attributes of the universe must be defined in a purely relational way, without reference to any absolutely defined, external parameter, it transpires that any slicing of spacetime into three-dimensional space-like hypersurfaces and a time dimension (any particular choice of foliation) is equivalent to any other. A general-relativistic description of the dynamics of the universe as a whole, therefore, does not allow to identify one particular dimension from among the four dimensions of spacetime as being that of time, given that the gravitational field equations remain valid regardless of the choice of a particular signature for the metric of spacetime. An additional difficulty also arises, due to the fact that the universe, as a particular instance of isolated system, must have an invariant total energy¹⁹, which would appear to imply that no meaningful change can take place on the cosmological scale, as if time was, in effect, irrelevant.

It seems that a similar conclusion would have to be drawn about the status of time in canonical quantum cosmology, where the same arbitrariness in the choice of a particular foliation and the same absence of change to the energy content of the universe would now apply to the many different histories of extended three-dimensional space-like hypersurfaces, which must be allowed to interfere with one another quantum mechanically. This is reflected in the fact that the most straightforward interpretation of the Wheeler-DeWitt equation (the equation that would allow to determine the wave function of the universe) requires assuming that it is similar in form to the *stationary* Schrödinger equation, while time is notoriously absent from such an equation. It is sometimes suggested that what those difficulties demonstrate is that the hypothesis that time exists as a unique dimension, distinct from the other

¹⁹The reader may recall that I have provided arguments in section 4.5 to the effect that the energy of the universe (just like its momentum and its angular momentum) must actually be null (even when space is assumed to be flat on the largest scale) if no characterization of the physical properties of the universe is to refer to external, or metaphysical elements of reality, because if the universe as a whole had a positive or negative energy it would become possible to identify a particular direction in time as being of absolute (non-relational) significance.

three dimensions of space, is incorrect.

It should be clear, however, that the absence of change on a global scale, which is assumed to be a consequence of the fact that the universe has a fixed value of energy, does not mean that time is not a meaningful concept for relating the changes taking place in one part of the universe with those occurring in another part of it, as long as we are actually dealing with different portions of the same universe, because it is not required of local subsystems that they have invariant energies as a consistency requirement and therefore change can certainly be observed to take place on an intermediary scale. In other words, even if we were to assume that time is irrelevant on a global scale, this could not be understood to mean that it has no clear significance as a means to relate local measures of changes. What's important to recognize is precisely that, from a cosmological viewpoint, time, as a dimension distinct from space, has meaning only as a relationally defined physical quantity that allows multiple local measures of change to be compared, thereby enabling all observers to provide a unique description of the various processes taking place in the universe (or within their associated causal horizon). Thus, it is an exaggeration to suggest that time does not constitute a meaningful concept in quantum cosmology²⁰. But to show that a conventional notion of time is not irrelevant to our description of reality on the cosmological scale, one must first explain how it is possible, in effect, for time to differentiate itself from the other three dimensions of spacetime, despite the fact that all four dimensions are kept on an equal footing and are required to be equivalent, from a fundamental viewpoint, by relativity theory. It is regarding this particular aspect of the problem of time in quantum cosmology that I would like to offer some original insight.

Two points must be taken into account in order to explain the existence of a uniquely significant slicing of four-dimensional spacetime into three-dimensional space-like hypersurfaces that would consistently select one single dimension as being that of time. First, it needs to be recognized that there must exist unique relationships of causality between all local events comprising an extended four-dimensional universe. Second, it must be recognized

²⁰The idea that a null value of energy, for the universe as a whole, would be indicative that time does not exist may be no more reasonable than the idea that a universe with null momentum (relative to the global inertial reference system determined by the average state of motion of all matter in the universe) would be indicative that space does not exist, which is so obviously inadequate a hypothesis that no one has ever suggested it could apply.

that, at the fundamental quantum gravitational level, it is possible for the principle of local causality to be enforced due to the existence of an embryonic element of time directionality in the causal structure of spin foams. Once this is recognized, then it becomes possible for a metric of spacetime with a unique signature to emerge that singles out one particular direction of four-dimensional spacetime as being that which is associated with the dimension of time across an entire space-like hypersurface (throughout the universe, on a given slice of spacetime). This is because, as I have explained in section 4.9, the homogeneity of the initial matter distribution at the Big Bang (which is responsible for the existence of a thermodynamic arrow of time) arises precisely as a consequence of requiring a constraint of global entanglement to apply uniformly over that entire slice of spacetime and this constraint is actually a condition for the existence of causal relationships between all elements of the universe which are present in this initial state.

It is important to understand that what distinguishes time from the other dimensions of spacetime, in a relativistic context, is merely the choice of a particular signature for the metric of spacetime which is arbitrarily imposed on solutions of the gravitational field equations in order that they satisfy observational constraints. But what this distinction provides is merely a separation of spacetime into past and future light cones along one particular dimension, which is really a requirement of local causality. Thus, if the signature of the metric was different and causality still operated uniformly, but along another dimension of spacetime, we would simply call this dimension time, while the other three dimensions would then all be analogous to conventional space. In fact, given that general relativity allows for local variations of the light-cone structure, one may say that what is produced as a result of spacetime curvature, or due to the presence of local gravitational fields attributable to the presence of matter inhomogeneities, are merely smooth local alterations of the direction in which causality operates.

Now, all that is required by the global entanglement constraint is that at least one space-like hypersurface exists over which the matter density is sufficiently uniform, down to the quantum gravitational scale, that no macroscopic event horizon is present. But, given that global entanglement is a condition that is imposed in order that causal relationships be allowed to exist between all parts of the universe, what is implied by this absence of macroscopic event horizon in the initial state at the Big Bang is that the direction in which causality operates actually is the same over all space, because the embryonic, quantum gravitational element of causal order must

itself be found to operate in the same direction of spacetime in all locations, over at least one such hypersurface and right down to the quantum gravitational scale (the Planck scale), thereby consistently imparting on spacetime a unique signature that is shared throughout the universe. I believe that this is what explains that the direction in which time is flowing is still mostly the same over all of space today (except in the presence of strong local gravitational fields and macroscopic event horizons), as necessary for the existence of a universal time variable.

To put things a little differently, one could say that, if there were significant local differences in the alignment of light cones on the quantum gravitational scale, in the initial Big Bang state, this would be equivalent (for what regards causality) to the presence of macroscopic event horizons and the presence of event horizons on all but the shortest scale is precisely what is forbidden by the global entanglement constraint, in the presence of negative-energy matter. Therefore, if global entanglement is necessary for the existence of the universe as an ensemble of causally interrelated parts, then there must exist one space-like hypersurface over which the light cones and time itself are oriented in the same direction of spacetime in every location (this is easier to visualize when space is assumed to be two-dimensional). It should be clear, however, that it is not merely the existence of causal or cosmic horizons that imposes a condition of global entanglement, because global entanglement is an independent consistency requirement for the existence of causal relationships between all elements of the universe, which actually allows the emergence of causal horizons as unidirectional phenomena (given that it allows the emergence of a thermodynamic arrow of time). It is not logically inappropriate, therefore, to argue that it is when global entanglement is imposed that causal order must be found to apply in the same direction of spacetime, uniformly, throughout the universe, as long as one recognizes that what is involved here is time-symmetric causality.

This is a significant result, because when a constraint of global entanglement is imposed on the initial Big Bang state in the presence of negative-energy matter particles, a strong limit is found to apply to early fluctuations in the density of matter, which means that local variations of the light-cone structure that determines how proper time intervals vary over extended space-like hypersurfaces are virtually absent, so that time flows uniformly over all space, as would be the case, by default, in a Newtonian context. The above argument would therefore appear to provide the basis for a satisfactory solution to one of the last major unsolved problem still facing the

most appropriate of current tentative quantum theories of gravitation, which is the question of how it is possible for a universal time variable to emerge from the timeless equations of the theory. Thus, it would no longer be necessary to appeal to anthropic arguments to explain, not only the observed time asymmetry and the unidirectional nature of causality, but really the very existence of a universal time variable.

Even though, from a classical perspective, relativity theory does not *a priori* require that there is a preference for one particular dimension of four-dimensional space over any other, the condition that there should exist causal relationships between all parts of that undifferentiated four-dimensional reality (between all the events taking place in it) implies that one direction in four-dimensional space is singled out, *uniformly*, as being that along which effects are propagated in the emerging spacetime and this is what gives rise to time as the continuous and uniformly flowing variable we are accustomed to experience on a macroscopic scale. The validity of the hypothesis that there does emerge such a singular dimension of time out of four-dimensional spacetime is what legitimizes a formulation of quantum cosmology as having to do with the dynamics of extended three-dimensional space-like hypersurfaces, whose histories can be described as unique trajectories in superspace (the configuration space of those three-dimensional objects). What is remarkable is that the viability of such a description is, in fact, a necessary condition for the elaboration of a consistent explanation of the quasiclassical nature of reality that emerges under conditions where irreversibility is a characteristic of the processes involved, as I will explain in the following section.

The problem that there was, originally, with the proposal that quantum cosmology has to do with the dynamics of extended three-dimensional space-like hypersurfaces is that the introduction of a fundamental element of causality in quantum gravitation requires a decomposition into positive- and negative-energy solutions, as in conventional, relativistic quantum field theory, and it was not clear how this could be achieved in the context of such a model. But even though this difficulty appears to have been overcome, I still believe that significant progress could be achieved in developing the current covariant framework of spin-foam quantum gravity into a fully satisfactory theory by taking into account the possibility for negative energy states to propagate both forward and backward in time, which constitutes a necessary step in allowing a proper integration of the requirements imposed by the generalized, classical theory of gravitation I have introduced in chapter 2. In any case, if local causality is, in effect, a decisive constraint on the

quantum gravitational scale, then time itself necessarily constitutes a meaningful parameter in quantum cosmology, even on a global scale, because the separation of four-dimensional spacetime into three dimensions of space and one uniformly-pointing dimension of time appears to be the defining character of a world that obeys the principle of local causality in the presence of negative-energy matter.

One must be careful, however, when considering a quantum-mechanical theory that purports to describe the whole universe, because, from a realist viewpoint, it would not be appropriate to describe the universe by using a wave function evolving deterministically over its entire history. Indeed, by doing so, we would commit the same error we make in the classical theory of relating all past and future three-dimensional space-like hypersurfaces in a predetermined way to some arbitrarily-chosen present state, which makes it look like everything about history is resumed in one single stationary state. In a more realistic situation, the whole history would not be determined from knowledge of one particular global state and following each local measurement the state of the universe and its wave function would need to be actualized, which would reveal the random nature of the history that actually takes place and the absence of *predetermined* relationships between the multiple extended three-dimensional spaces forming a history, which in turn illustrates the relevance of time in characterizing the actual relationships. Even in the context where a unique future is assumed to exist in the same way a unique past does, there is no rational motive to argue that time, as a measure of change, becomes an irrelevant notion, because such a conclusion would only be valid if we ignored the random aspect of quantum-mechanical processes (which is particularly unavoidable in the context of the existence of closed causal chains) and if we neglected the constraint imposed by the requirement that all parts of the universe be causally related, which singles out the state of maximum matter density of the Big Bang as a state of minimum gravitational entropy from which all future evolution is taking place irreversibly, as I explained in section 4.9.

Anyhow, it must be clear that, despite what is sometimes suggested, it is not true that time, or even space, do not exist at all in modern quantum gravity. Indeed, a certain embryonic notion of space is clearly present in the structure of spin networks, which allows classical space to emerge naturally when a sufficiently large number of fundamental, discrete elements of structure are combined according to purely quantum rules. Furthermore, even in such a context, we are still dealing with four-dimensional boundary

conditions and this is certainly indicative of the relevance of time, even if this parameter may not explicitly appear in the equations which allow to determine the correlation probabilities associated with those four-dimensional boundary conditions. Actually, the mere fact that, even in a quantum gravitational context, we are still speaking about ‘local’ changes occurring in the configuration of spin networks means that an additional degree of freedom *must*, as a fundamental requirement, be allowed to emerge, which relates those local changes to one another. The problem that there was, originally, is simply that, in the absence of a constraint of global entanglement, no universal time variable was allowed to emerge, because no unique direction appeared to exist that would be associated with this degree of freedom and along which events could be sequentially ordered into some kind of *universal causal chain*. Once it is recognized that causal relationships must exist among all elements of the universe, however, then the most appropriate of the current fundamental theories do allow a certain notion of history to emerge given that, in the presence of negative-energy matter, this condition allows one particular dimension of four-dimensional spacetime to be singled out uniformly, throughout one extended spin-network configuration, as being that along which causal order is established and for this reason alone, those extended space-like configurations may be considered to constitute the dynamic elements of a quantum theory of cosmology.

However, in my opinion, what would definitely invalidate a truly timeless quantum theory of gravitation is precisely the fact that such a theory would be incompatible with the existence of a fundamental time-direction degree of freedom (such as revealed in particular by violations of time-reversal symmetry T), while I have shown, in chapters 2 and 3, that such a property is essential to a consistent description of physical reality, in a semi-classical context. Indeed, once it is recognized that, in quantum field theory, the propagation of elementary particles can take place along any of two opposite directions of time, independently from the constraints imposed by thermodynamic irreversibility, then a conflict emerges with the timeless viewpoint, given that if there is no time, then obviously there cannot be a fundamental direction in time, because any relationship of time directionality must necessarily involve a sequence of events causally related to one another following a definite and unique order, distinct from their spatial order, even when the classical spacetime structure in which those events are embedded is assumed to emerge from the combination of discrete elements.

Given the nature of the arguments which are usually proposed to support

the conclusion that time is irrelevant in quantum cosmology and therefore that it may not even exist, it would seem that solipsism is once again to blame for misleading even some of the most brilliant thinkers into this theoretical dead-end. Indeed, what a rejection of time would require us to assume is that there can be change and that all changes can be related to one another by the use of a reference system we call time, but that this is not enough to justify the conclusion that this reference system is the reflection of something real. Thus, while we are allowed to recognize the emergence of a certain variable, distinct from spatial position, which is useful for comparing various local measures of change involving one or another physical attribute, and while the assumption that such a variable exists is undeniably useful and allows to simplify our description of reality, the fact that it is not possible to directly measure any changes relative to that additional variable itself and the fact that this variable may no longer be globally significant under the most extreme conditions (in the presence of very strong local gravitational fields) would mean that it cannot be considered a real physical property, even under more ordinary circumstances. All arguments against the existence of time as a meaningful concept in quantum cosmology involve such an element of solipsism. Time does not exist because it cannot be subjected to direct observation, or be the object of some measurement that would confirm that it is real. But that is just a perfect example of the kind of irrational conclusion one can draw based on such considerations, because what can be more obvious in fact, from our experience of physical reality, than the existence of change and the reality of time?

Now, it has been argued that it might be possible for time to emerge as a mere thermodynamic phenomenon, despite the fact that it would not really exist from a fundamental viewpoint. What I'm talking about is the concept of 'thermal time', according to which the passage of time would actually be an illusion attributable to the fact that the irreversible time of our conscious experience appears to always be associated with heat dissipation, which would appear to single out one particular physical variable as that relative to which energy remains unchanged, while in fact there would be nothing fundamental with such a variable. But the problem with this proposal is that there is, in fact, plenty of evidence for the relevance of a more conventional notion of time at the level of elementary particles, where irreversibility is not a defining characteristic. Of course, the fact that there would be no preferred direction of time in the absence of heat dissipation is not completely irrelevant to the problem of the existence of a classical spacetime continuum (given that dis-

sipation appears to be necessary to explain the decoherent nature of space and time), but it is not that significant either, because we are not merely trying to decide whether unidirectional time is a valid concept, but with deciding if the whole concept of time is, in effect, relevant to a description of physical reality. However, if thermodynamics was the ultimate explanation for the existence of time, it would not be necessary to wait until we begin to explore reality on the quantum gravitational scale to witness an absence of time, because many phenomena are known to exist, on a much larger scale, that do not involve any irreversibility and yet they are still describable using space and time coordinates²¹.

It is important to point out that if we were to assume that time really doesn't exist, even under ordinary circumstances, we would then be left with having to conceive of the present as just one independent, momentary state among many possible states devoid of any causal relationships with one another. It was, in effect, suggested by Julian Barbour that such causally independent, momentary states may not be incompatible with our perception of the passage of time, if we assume that all that we really experience are momentary states of consciousness, which might be more appropriately described as memory states. But the problem here, again, is that even if such an explanation of consciousness as a state rather than as a process was possible (which I believe may not really be the case²²) we would then have no explanation for the fact that the present state of the universe, in which the state of our consciousness is contained, is one which is characterized by the existence of a large number of mutually consistent records of a unique lower entropy past, because such a configuration would not likely be chosen in a random trial, out of all the possibilities which would appear to exist for a momentary present state. The fact that what can be characterized as long-term records are usually preserved in what appears to be the most

²¹In the context where a satisfactory solution to the problem of the origin of thermodynamic time asymmetry that is not based on the weak anthropic principle is now available (this was the subject of section 4.9 of this report), the fact that the thermal-time hypothesis may appear suitable for an explanation of cosmological time-asymmetry based on a certain interpretation of entropy growth as a purely subjective, observer-dependent phenomenon would no longer constitute a potential advantage of a timeless interpretation of quantum cosmology.

²²Memory, as well as other basic mental faculties, are not really static events, but rather processes which require a certain duration to be experienced and if there is no duration, what one should expect to experience is not one everlasting memory, but nothing at all, which is certainly not compatible with my own experience of reality, at least.

stable structures, while short-term memories are usually preserved in more rapidly changing structures, would also remain unexplained from a timeless universe perspective.

There were many attempts at trying to explain why such present states as revealed by our personal experience of reality may not really be unexplained, even when one assumes that all that exists in the universe is an extended space without any time. But in the end, one must recognize that those proposals are inadequate and that the unlikeliness of the observed present configuration of our universe remains a complete mystery, unless one is ready to assume that what one actually observes is not really indicative of the existence of a lower entropy past, even though there is absolutely no rational motive (even of an anthropic nature) to legitimate the validity of such a conclusion. Of course, if it had actually been demonstrated without doubt that time does not exist, then we may have no choice but to assume that everything is such a strange and deceptive illusion, but this is not true and the only reasonable conclusion we are allowed to draw from our observations is that the present state of the universe, regardless of its exact nature, must be related to one single past history through the existence of unique (but not predetermined) causal relationships unfolding back in time to the state of minimum gravitational entropy that allows to explain the existence, in the present state, of mutually consistent records of a unique past.

It is usually recognized, in fact, that all that one may reasonably argue, concerning time as a quantum gravitational concept, is that it is the continuity of its flow and the existence of a unique spacetime metric signature which do not apply at the most fundamental level. Thus, if, at some point, there was such a strong desire to do away with time, it is perhaps only due to the fact that we were unable to explain the singular character of time as a dimension of spacetime, because, in the absence of guidance from the generalized theory of gravitation I have introduced in the second chapter of this report, we couldn't understand the profound significance of the homogeneity of the initial distribution of matter energy at the Big Bang, which allowed me to explain the near uniformity of the direction of propagation of effects in spacetime and therefore of the flow of time itself. In a traditional context, it was rather convenient to simply assume that time does not exist at all, given that, like space itself, time is not present in its classical form at the most fundamental level. But it must be clear that if time, or more specifically causal order, did not exist in any form at a fundamental level, then what we should definitely not experience is a dimension of time distinct from the

other dimensions of space.

Now, despite the fact that I have criticized Julian Barbour's suggestion that our experience of the passage of time may not be incompatible with a timeless description of reality, I must recognize that he, more than anybody else, is responsible for having convinced me of the validity of the concept of simultaneity hyperplanes, or more generally of space-like hypersurfaces as the basic building blocks of a dynamical theory of space that would be relevant to quantum cosmology. The only problem I have with Barbour's interpretation has to do with his insistence that those global states of the universe as a whole should all exist independently from one another and therefore cannot be causally related to one another following a unique and well-defined order (cannot be considered to form one single causal chain or to take part in one single history). But in fact, this need not be considered a requirement of a dynamical approach to quantum cosmology and as I have explained above, it would rather seem that there must exist unique causal relationships between those properly defined global states, despite the fact that there appears to be a lot of freedom in how spacetime can be sliced into such space-like hypersurfaces.

We may, therefore, retain as valid the concept that the present state of the universe as a whole, regarding, in particular, its gravitational field or spacetime curvature, is provided by the current configuration of one such space-like hypersurface, which may be represented as a point in the appropriate configuration space (say the superspace of canonical quantum cosmology), while the time variable would enter the picture as the position along the actual trajectory followed by the global state in this configuration space. This becomes a valid proposal in the context where we now have a valid explanation for how it can be that one given spacetime dimension is uniformly singled out as that along which local causality is allowed to operate (as reflected in the uniqueness of the signature that must be assigned to the metric of spacetime) and to constitute a physically significant constraint that does not apply in the case of the other three dimensions of space, even in a general-relativistic context.

To be honest, I have to mention that the conclusion that the history of a universe's space curvature can always be represented as a path in the configuration space of three-dimensional space-like hypersurfaces is dependent on the hypothesis that any solution of the gravitational field equations that contains closed time-like curves (those hypothetical configurations of the curvature of space which would make conventional time travel experiences a

reality) can be excluded. Usually, this is recognized to be possible merely if we assume without reason that the second law of thermodynamics is valid under all conditions. But given the explanation I have provided in section 4.9 for the existence of the thermodynamic arrow of time, the conclusion that closed time-like curves cannot naturally arise actually becomes unavoidable. Indeed, under such circumstances, the constraint that gives rise to thermodynamic time asymmetry must always operate in the same unique direction of time and invariably have as a consequence the diminution of entropy in the particular direction of time that points toward the initial state of minimum gravitational entropy of the Big Bang, as a requirement for the existence of causal relationships between the various elements of the universe. Therefore, a universe could not even exist, as a causally interrelated ensemble of space-like separated elementary particles, if it did not satisfy this unidirectionality constraint, which would be the case if the direction of entropy diminution could not be well-defined as a result of the curvature of space and this means that closed time-like curves are actually forbidden. From my viewpoint, it would therefore appear that it is *always* possible to represent the spatial properties of the universe and their entire history as some monotonic foliation of space-like hypersurfaces, that is to say, as a path in superspace.

It is, therefore, the existence of a unique direction in spacetime, along which effects must propagate, either forward or backward, that allows histories to be parameterized by a universal time variable (associated with a particular slicing into space-like hypersurfaces) and that enables a description of space curvature as evolving with respect to this time variable, thereby legitimating the notion of history as consisting in an ensemble of causally related global states, that is to say, a universal causal chain. What I have shown is that the *apparent* absence of a fundamental distinction between time and the other three dimensions of spacetime, which is an essential feature of relativity theory, does not constitute an insurmountable obstacle to achieving this objective, so that we are no longer justified to conclude that time is altogether absent in quantum cosmology. This is certainly a significant result for the elaboration of a solution to the problem of the interpretation of quantum theory, given that the existence of classical space and time is actually required by conventional quantum theory, for the description of histories, in the context where the various macroscopic experimental conditions which are shared by both the retarded and the advanced portions of a quantum process must be defined over one unique and classically well-defined spacetime continuum. Thus, spacetime itself must be assumed to be decoherent

under conditions where a history can be consistently defined, which means that quasiclassicality must already apply to the gravitational field in order that decoherence be observed at a higher level in the observed attributes of conventional quantum systems.

This, again, illustrates the fact that a continuous and uniformly oriented dimension of time must be allowed to emerge from a quantum theory of gravitation²³ before ordinary quantum processes can be appropriately described and conventional quantum theory itself can become a valid representation of reality, with clear and precise meaning at the most fundamental level. The problem of the emergence of time in quantum cosmology must, therefore, be recognized as constituting one particular aspect of the more general problem of the nature of the conditions necessary for the emergence of a quasiclassical world. What this means is that in order to obtain a satisfactory interpretation of quantum theory, one must first examine in which way gravitation and the curvature of space could be subjected to the same time-symmetric description as would apply to more conventional physical attributes under ordinary conditions. Achieving such an objective will allow me to identify additional constraints from which both the decoherent nature of spacetime and the persistence of quasiclassicality that characterizes all observed aspects of physical processes can be expected to arise, even in the context where quantum theory is assumed to be valid under all circumstances. What those considerations will demonstrate is that it is not just general relativity which really is a theory of the universe as a whole, as is usually recognized, but that quantum theory, from the viewpoint of its most accurate interpretation, is also essentially a cosmological theory.

5.12 Universal causal chain and quasiclassicality

We are now finally in position to examine how it is exactly that quantum theory can be extended, so as to become fully consistent from both a log-

²³Of course, even on the astronomical scale, the spatial uniformity of the flow of time is only an approximation, because the metric properties of space and time are influenced by the presence of positive-energy matter and by the inhomogeneities which are present in what remains of the negative-energy matter distribution, which means that, even from the viewpoint of the approach favored here, there is still no universally valid measure of the passage of time.

ical and an experimental viewpoint. It is here that all the breakthroughs achieved in the preceding chapters of this report, as well as in the preceding portions of the present chapter, while trying to provide a better understanding of so many aspects of physical theory associated with time directionality will converge to produce their most significant outcome: a logically consistent interpretation of quantum theory that is valid at absolutely all levels of description. It is certainly a positive development, already, that, in the preceding section, I have been able to conclude that time is still relevant to a description of our universe in a quantum-mechanical context. Under such conditions it becomes appropriate to define the intrinsic space curvature over a particular three-dimensional slice of spacetime at one particular moment as consisting of a single point in superspace. The role of time then emerges quite straightforwardly as being that of relating those global states of the universe to one another into some kind of universal causal chain, while establishing the sequential (chronological) order of events.

What's remarkable is that the existing mathematical framework by which this particular approach can be formalized, which originates in the ADM formalism²⁴ [61], allows history itself to be described as one particular trajectory in superspace [62] [63] [64]. Time, therefore, must be conceived of as the global variable to which are related the multiple local measures of change that take place as the curvature of space evolves along such a trajectory in superspace. This allows to fulfill Reichenbach's vision of time as reducing, in its most essential form, to the general concept of a causal chain, that would allow to establish and maintain the invariant local topological ordering properties of spacetime, even when its metric properties are subject to local variations. From my viewpoint, however, it would not be appropriate to consider a traditional concept of causal chain that would involve irreversibility at a fundamental level, as Reichenbach contemplated, because irreversibility is a property that must rather emerge from the particular boundary conditions which existed at the Big Bang.

²⁴Despite the fact that the ADM formalism of quantum gravity (based on ADM variables) has been replaced by the more appropriate formalism of loop quantum gravity based on a formulation of general relativity theory in terms of connection, or Ashtekar variables, I think that it is still appropriate to describe the general concept of a dynamical theory of space using the original approach to quantum cosmology based on superspace, which allows to visualize in a more intuitive way the phenomena involved and to more easily understand how the realist interpretation of quantum theory developed here can be applied in a cosmological context, as far as the discreteness of space can be neglected.

In any case, it must be clear that it is the network of *local* relationships that varies as we move along a trajectory in superspace, because, from the viewpoint of its total energy content, the universe, as the ultimate isolated system, would appear to remain in the same state without any change actually taking place (this is what motivates the unsubstantiated claim that time may not be relevant to quantum cosmology, as I explained in section 5.11). It must also be emphasized that what is provided by the concept of space-like hypersurface is not a unique and absolutely defined characterization of reality, because, even when a universal time variable is allowed to emerge, there are still many equivalent ways by which spacetime can be sliced into three-dimensional simultaneity hyperplanes (because simultaneity itself is a relative concept), which would appear to require a history of the universe's space curvature to consist, not in a unique trajectory in superspace, but rather in a given surface in the same infinite-dimensional configuration space, formed of the many equivalent trajectories which are associated with the same unique history of spatial curvature²⁵. What must be clear, then, is that even if many equivalent possibilities exist for such a trajectory, they all provide alternative descriptions of the same causal chain, to which corresponds one unique history. Once again, the freedom that surrounds the choice of a suitable slicing of spacetime must not be considered to reflect the irrelevance of time for a description of the dynamics of space on the cosmological scale, as it is merely a reflection its relational nature.

Now, from the perspective of the developments introduced in the first portion of this chapter, it would appear that a quantum-mechanical description of the metric properties of space, relating to the universe as a whole, cannot merely involve adjoining a wave function to some boundary conditions defined over superspace, under the assumption that all possible histories compatible with those conditions happen, all at once, as different branches, in the same universe. The purpose of a quantum cosmology would rather be to estimate the probability of observing a global state of intrinsic space curvature (represented as a point in superspace) when another such global state has been observed at a certain time in the past, by summing-up the (positive and negative) *probabilities* associated with all the different ways by which those two points can be joined together as a result of the global state

²⁵The modern spin-foam quantum theory of gravitation allows to more appropriately deal with this freedom and to formulate the approach discussed here in a fully covariant way with the additional benefit of providing a discrete, or quantized description of space and time.

evolving, once forward and once backward in time, along two possibly distinct trajectories in superspace for which the local curvature of space itself could differ, as long as those differences remain unobservable, that is to say, without irreversible consequences.

Here, again, we face the mystery of the existence of two interfering histories occurring in parallel, which would appear to merely complicate the causal chain picture of the universe's history by actually requiring bidirectional causality to operate in opposite directions along two otherwise similar portions of history. But, even though this aspect of a quantum-mechanical description of the universe is certainly convenient, given that it allows to explain quantum non-locality, it nevertheless remains unexplained. In order to begin to understand why this dual character of quantum reality is not as arbitrary and superfluous as it may seem, one must first examine how it is that causality would operate if there was no advanced portion to the history of the universe.

It only became clear to me what the organizing principle is that allows to clarify this situation when I began working on the problem of time travel and closed causal chains. It is at this point that I realized that, if the history of the universe was described by one universal causal chain, freely unfolding in the appropriate configuration space, along one particular direction corresponding to unidirectional time (say, that along which entropy is growing globally), there would need to be external causes that would determine how the universe began to get going along the particular trajectory over which it is found to have propagated in this configuration space (which, for now, may be assumed to be superspace, even though, ultimately, one would need to consider a more general kind of configuration space that would also encode non-gravitational degrees of freedom). This is a very important point, as an external cause is precisely what must be considered forbidden by the constraint of relational definition of the physical attributes of the universe, which basically implies that there should be no 'first cause' that would need to be attributed to some external agent that is not part of the causal structure of the universe and that may not be governed by the same physical laws²⁶.

The reader may recall the problem associated with so-called knowledge paradoxes, that would arise from the viewpoint of unidirectional time when

²⁶The same inconsistency would arise if the condition of continuity of the flow of time along a particle world-line, which was introduced in section 4.3, was allowed to be violated and therefore this previously discussed constraint can be understood to really be a condition for the *local* continuity of all causal processes.

a time traveler would take a copy of some complex and highly valuable work of art, which happens to exist in the future, back to a time in the past before which it did not yet exist, thereby allowing it to be created instantaneously, without any apparent cause, so that the invention is allowed to exist in the future, which is necessary if it is to be brought back in time. I have explained in section 5.4 that such a phenomenon is not impossible in principle, but is simply very unlikely to occur, because it would actually require entropy to increase in the past direction of time, while the time traveler would be in the process of bringing back information from the future, which would constitute a violation of the second law of thermodynamics given that it would involve a decrease of entropy in the future.

What can be learned from such a thought experiment is that, if the phenomenon described here is extremely unlikely, it would not, however, constitute a violation of the fundamental (time-symmetric) principle of local causality, because it would only involve a diminution of entropy that would be apparent from a unidirectional-time viewpoint, but would not require a real discontinuity in the flow of information along the direction in which the time traveler would be progressing in time. Indeed, as I explained in section 5.3, it must be recognized that there is no absolute difference between causes and effects at a fundamental level and this means that the future can influence the past just as much as the past is allowed to influence the future, even in the same portion of history (as long as no inconsistency develops), which is what actually happens when an elementary particle is propagating backward in time (in which case it behaves as an antiparticle). But if the present state of the universe was determined by a certain cause (located either in the past or in the future) that is not itself determined by an earlier or later cause that also belongs to the universe itself, but that would be necessary to set the universe on its course along one particular trajectory in superspace (with which is associated one particular, initial, global state and one particular information content), then a real problem would emerge, because, under such conditions, bidirectional causality would definitely be violated. Indeed, even if time was to actually begin at the moment when matter emerges from the past singularity, in the initial Big Bang state, there would still be a discontinuity in the causal chain trajectory and this is why I argued in section 4.5 that it is not desirable that matter be created out of truly nothing at the Big Bang.

But how could one avoid the conclusion that there needs to exist an external cause that would determine the initial (or the final) state of the

universe, in the arbitrarily far past (or the arbitrarily far future), that is to say, how could one explain what determined the information contained in the extended three-dimensional space-like hypersurface that constitutes the starting point along the universal causal chain that evolved into the present one? I believe that the truth is that we have no choice and that we must admit that a certain hypothesis, which may at first appear gratuitous and arbitrary, actually constitutes an absolutely essential condition that needs to be imposed in order that our quantum-mechanical description of reality be free of logical inconsistencies when it is applied to a description of the universe as a whole. It is at this very precise point that quantum theory ceases to be baffling and that its most incomprehensible aspects become essential elements of a fully comprehensible representation of reality. What emerges from the original perspective developed in this chapter is that the history of the universe is nothing but an elongated, *closed* causal chain that unfolds in superspace (or in some generalized configuration space where matter degrees of freedom would also be represented). There is no first cause. The initial impetus that sets the universe on its course is provided by the universe itself, as all later states of the universe also constitute earlier states along this closed, universal causal chain. The universe truly brings itself into existence by providing the cause of its own present condition as being nothing but a remote effect of this very same present condition.

Perhaps that you remember my earlier discussion of the closed-circuit analogy from section 5.2. What I explained is that most electrical circuits are really closed circuits and if they may not seem so under ordinary circumstances, it is simply because the circuits are usually extended in one particular direction and can only be recognized for what they really are by the fact that the cables in which they are confined are always composed of a pair of polarized wires, which betrays the fact that this unique path that seems to extend from source to sink is actually formed of the two branches of a closed circuit in which the current flows in opposite directions. Well, I believe that one must come to accept as unavoidable that this is a very good analogy of what is described by the quantum-mechanical version of the history of our universe. This history is a closed trajectory in superspace that is stretched to near infinite proportion along the direction relative to which unidirectional time unfolds, as allowed by the solution I have provided in the previous section to the problem of the origin of the differentiation between space and time.

What I suggested, in the previous section, is that the existence of a time

dimension distinct from the other three dimensions of space is an outcome of applying to the initial maximum-density state of the Big Bang a constraint of global entanglement, as a requirement for the existence of relationships of local causality on the quantum gravitational scale, for the universe as a whole, which has for consequence that the same unique direction in spacetime is selected throughout the universe for the propagation of causal signals. But such a distinction between space and time (which is made apparent by the unique signature that must be attributed to the metric of spacetime) is what allows to consistently describe the history of the universe as consisting of a trajectory in superspace. What a time-symmetric, quantum-mechanical description of the same reality allows, then, is for this trajectory to be a ‘polarized’ version of history, in the sense that it actually consists of two parallel histories which share the same observable macroscopic conditions, but whose corresponding segments are being propagated in opposite directions of time. Although this pairing of history and this polarization would remain a complete mystery from a conventional viewpoint, in the context of the above discussion it becomes a natural and essential feature of physical reality that should actually have been expected all along, if only we had recognized that, from the viewpoint of logical consistency, causal self-determination is not an optional requirement for the universe.

Indeed, if causality is of any relevance to cosmology, it is certainly due to the fact that it imposes two essential conditions on the universe in order that it be allowed to simply exist in any possible way. The first of those two conditions is that all elementary particles present in the universe must be causally related to one another as a result of having been in local contact with one another at least once in the history of the universe. As I explained in section 4.9, this must be considered necessary in order that all particles be allowed to actually consist of different elements of the same universe. The existence of such a condition, which is responsible for the low gravitational entropy of the initial Big Bang state of maximum matter density, is what allows me to assume that the history of the universe is, in effect, described by one unique trajectory in superspace, rather than by multiple, distinct (nonequivalent) and unrelated trajectories, which would really constitute the histories of many different universes, not causally related to one another. But as I just mentioned, this global entanglement constraint is also responsible for the fact that time actually exists as a dimension distinct from the other three dimensions of space on a global scale; which is responsible for giving rise to the very causal structure of spacetime (the near-uniform alignment

of light-cones). What's more, I will argue, below, that this condition is also necessary to explain the classical nature of reality, under conditions where some dynamic attribute of a quantum system becomes entangled with irreversibly evolving degrees of freedom of the environment in which the system evolves.

The second condition would then be that which I have just identified and which is that the universe must be self-determined from the viewpoint of causality. This can be satisfied when the history of the universe consists of a closed causal chain, represented by a closed trajectory in superspace, which requires the universe to eventually return to the exact same (but partly unobservable) state in which it currently is, as it evolves along this trajectory. This condition is what explains that it is necessary, in order to obtain the right correlation probabilities, to take into account the existence of two otherwise independent histories evolving in opposite directions of time, which is the distinctive feature of the realist, time-symmetric interpretation of quantum theory developed in the preceding sections of this chapter. What defines a universe, therefore, is not just the fact that all of its constituent elements (the particles) are causally related to one another despite the *spatial* distances that separates them, but also the fact that the *global* configurations of those constituent elements are all causally related to one another and to nothing else (they from a unique causal chain). When history consists of a closed trajectory in the space of all possible configurations, every single, global state can be in local 'contact' (through time) with both a preceding and a succeeding state and this is what allows all global states to be causally related to one another, regardless of the 'distance' that separates them *in time*. Thus, the multiverse is not merely the ensemble of all possible, causally independent universes (those which may be characterized by distinct values of their global states of space curvature and other physical attributes at arbitrarily-chosen times), it is really the ensemble of all possible, universal causal chains which exist as nonequivalent, closed trajectories in superspace.

What is essential to grasp is that, despite what would seem to be implied by the progress that had already been achieved towards the elaboration of a consistent time-symmetric interpretation of quantum theory, even though there appears to be two causally independent, but interfering histories to every process, from a cosmological viewpoint there is, in fact, only one history, but it feeds back on itself, so as to form a closed causal chain. But for some reason to be discussed below, the trajectory in superspace that corresponds to this causal chain goes through different, but observationally

indistinguishable states, once by progressing in some direction of superspace as time goes, and then once again by progressing in the opposite direction (the state vectors corresponding to those two portions of history may differ, because the retarded state vector is determined by past conditions, while the advanced state vector is determined by future conditions, but the observable macroscopic constraints themselves do not differ). Thus, there is no quantum system in a state of superposition, going at once, and in the same universe, through all possible histories. There is one unique history, the details of which remain in part unobservable to any observer, that unfolds as a closed causal chain in superspace (or some generalization of it), subject to the condition that this evolution happens along two mostly parallel trajectories, joined together at each extremity, as if two histories were occurring in opposite chronological orders, whose corresponding segments share all observable physical properties.

But it is, in effect, only in this particular sense that we may assume history to be unique, because, even though we are always taking part in only one history, two independent histories are unfolding all at once (in opposite directions of time) relative to *unidirectional* time, which are merely required to share the same observable macroscopic conditions, despite the fact that they differ with respect to most unobservable physical attributes. This interpretation allows to explain the fact that the interfering realities are not in causal contact with one another locally, because, even if the two portions of history share the same macroscopic conditions, they do not really happen at the same epoch and therefore the particles present in the retarded portion of a process cannot interact with those which are present in the advanced portion of what only appears to be the same process. As a result, it is no longer necessary to assume without justification that the particles that take part in different histories do not interact with one another, in order to avoid the contradiction that emerges, in the context of a more conventional approach, when it is assumed that all branches of reality coexist in the same portion of the universe's history.

It must be clear that while this approach allows the state of the universe, which is specified by a particular point along the trajectory in superspace, to be characterized by simultaneously well-defined values of conjugate attributes, observable data would still be subjected to quantum indeterminacy, because no observer that is part of the universe can determine the exact states of both a dynamic attribute and its conjugate counterpart by imposing one particular set of experimental constraints. Thus, when the ex-

trinsic curvature of space (associated with the rate of change of the intrinsic curvature of space along a trajectory in superspace) or the particle momenta are determined with arbitrarily good accuracy, the intrinsic curvature or the particle positions become totally undetermined (from an observational viewpoint), even if there always exists a definite state of intrinsic space curvature that corresponds to the relevant point on the trajectory in superspace and this is reflected in the fact that the unobserved attribute is allowed to have completely different, interfering values in the retarded and advanced portions of history. The only difference between this situation and that which would appear to exist from a more conventional quantum-mechanical viewpoint is that it can now be assumed that there does exist a unique reality, in each portion of history, for the unobserved attribute associated with a given set of observational constraints, even though this reality can differ for the retarded and the advanced portions of history and cannot, as a matter of principle, be subjected to direct experimental knowledge.

But despite the enormous clarification and simplification which are made possible by the adoption of such a viewpoint, it would remain to explain why it is, exactly, that we are allowed to expect that the intrinsic curvature of space, as well as other experimentally determined macroscopic conditions, do not differ much, most of the time, for those two portions of history, that is to say, we still need to explain why it is that, under proper circumstances (when a given dynamic attribute is under observation), the same unique, classical path is shared by the processes which unfold in both the retarded and the advanced portions of history, as required for the conventional mathematical framework of quantum theory to be compatible with what is observed on a sufficiently large scale. As I previously mentioned, this is a particular aspect of the quantum measurement problem, or the problem of the origin of the quasiclassical nature of observed reality.

Actually, as I will explain below, it is the fact that the history of our universe consists of one single, closed trajectory in superspace that allows quasiclassicality to naturally emerge as a property of the physical world, under appropriate conditions, and therefore it will be apparent that the fact that our world is, in effect, classical on a sufficiently large scale allows to confirm the validity of the hypothesis that the history of the universe constitutes a circular process that feeds back on itself. Thus, if it was not for the closed, or circular nature of quantum-mechanical history, we would really need to assume that, for some reason, two independent, quantum-mechanically interfering processes are taking place at the same time, all the time, despite the

fact that it would then be impossible to explain why it is that the retarded and advanced portions of a process actually share the same experimental conditions or the same metric properties of space (because the criterion of consistency, specified by the consistent-histories interpretation of quantum theory and enforced by decoherence, is insufficient to achieve such an outcome, as I explained in section 5.10).

One thing should be clear already, though, and it is that if the history of the universe consists of a closed causal chain, then the retarded and advanced trajectories in superspace must be smoothly joined at some point in what appears to be the future from the viewpoint of unidirectional time and also at a certain point in what appears to be the past from the same unidirectional-time viewpoint. As a result, no two points on the universal causal chain can be absolutely characterized as ‘earlier’ or ‘later’. But it is also clear that the directions of propagation in time along two corresponding segments of the closed causal chain (those which appear to be in the same macroscopic state at a given instant of time) have significance merely as relationally defined properties (only the difference between those two directions has physical meaning), because no direction of propagation can be attributed absolute significance. Thus, the direction of time associated with one or another portion of history along the universal causal chain is not a direction in configuration space, but a relationally defined property of the universal causal chain itself. Yet, the growth of (gravitational) entropy does allow for the existence of an objectively defined arrow of time, to which can be compared the direction of propagation in time along any given segment of the universal causal chain, and it is from the viewpoint of this unidirectional time parameter that the universal causal chain would eventually appear to close, at which point time would come to an end.

Now, it may appear that the hypothesis that the superspace trajectories associated with the retarded and advanced portions of history must be smoothly joined at a certain point in the future could never be proven right, given that, from the unidirectional-time viewpoint, the closure of the universal causal chain cannot be observed unless it has already occurred, in which case we would no longer be there to acknowledge this fact. But, as I mentioned above, the validity of the theoretical requirement of closure can actually be confirmed by the observation that reality is of a quasiclassical nature, for reasons I will soon explain. The existence of such an end of time, however, must be distinguished from that which would occur as a result of an interruption of the trajectory of the universal causal chain in super-

space, which, despite what one might be tempted to assume, is not really a possibility, given that it would not constitute a simple bifurcation point in unidirectional time, but would involve a causal discontinuity, even from a bidirectional-time viewpoint.

One important aspect that needs to be emphasized here is that the situation of a universe whose trajectory in superspace is submitted to a condition of closure is not the same as the situation of a universe which would evolve, as a result of Poincaré recurrence, to the exact same *observable* state in which it was at an earlier time (a state defined by the same observable macroscopic conditions, to which may still correspond two different state vectors), which, in principle, could be satisfied even if the superspace trajectories associated with the retarded and advanced portions of history do not merge at any point along the shared, coarse-grained trajectory that would take the universe to its earlier observable state. In the present case, it must be assumed that, when the universal causal chain closes in the future, this will be due to the fact that the retarded portion of history has, by chance, found itself in the exact same state as that in which the advanced portion of the process turned out to be, not just at an observable level, but even for what has to do with the unobservable state of those physical attributes which are the subject of quantum interference.

The evolution along the closed causal chain, therefore, will not merely take the universe to a state that is similar to that in which it once was, but eventually to the exact same point it once occupied in superspace, from which any further evolution would take the universe into the exact same history through which it once went, despite the random nature of this evolution (because this would actually be the same history). Yet, if an observer is present when the bifurcation point is reached, she would not be able to experience the same history she once experienced, but in reverse, because what would happen is indeed a reversal of the direction along which the causal chain is unfolding with respect to *unidirectional* time, which means that the thermodynamic arrow of time would reverse, while reality can only be experienced in the opposite direction along the closed trajectory (the same direction as that in which the observer experiences reality in the retarded portion of history). Under such conditions both an observer that is part of the retarded process and its counterpart that evolves as part of the advanced process would experience the same reality and each of them would simply cease to exist at the bifurcation point, because consciousness is a thermodynamic process that necessarily takes place along the direction of time in

which entropy is rising globally. But it must be understood that the closure of the universal causal chain does not take place in position space, but really in superspace and therefore it does not involve an annihilation of the particles present in the retarded portion of history by those present in the advanced portion of history and it is not limited by the requirement of energy conservation that would otherwise need to apply with respect to unidirectional time, because, in such a case, continuity is only required to apply from the bidirectional-time viewpoint.

Thus, the point in the future at which the retarded and advanced trajectories of the universal causal chain would merge, from the unidirectional-time viewpoint, does not need to have extraordinary properties and could be any instant of time. Again, though, it must be clear that the closure of the universal causal chain is a phenomenon that takes place in superspace and therefore it may appear to violate the principle of local causality by occurring all at once in position space. Indeed, even if the bifurcation process may appear to take place at different times, in distant regions of the universe, from the viewpoint of certain observers, once it happened in one region of the universe it would have to occur in all the other regions, as the condition that is responsible for the continuous decrease of entropy in the past direction of time does not allow for oppositely directed thermodynamic arrows of time to be present simultaneously in the same universe. Also, if time can be extended to instants past the initial Big Bang singularity, then the moment in the past at which the universal causal chain would close would not necessarily need to be that at which the Big Bang itself occurs, but would likely be a time, arbitrarily distant in the past, prior to the Big Bang, when, by chance alone, the retarded and the advanced trajectories would meet in superspace. But it must be clear that the advanced portion of the known history is not the trajectory that unfolds prior to the Big Bang. Both the current history and that which may have taken place (with entropy growing in the opposite direction of time) before the Big Bang (as apparently allowed by certain quantum gravitation theories) have their own retarded and advanced portions of the same closed causal chain and could actually be very different histories, involving different sets of observable events.

If I believe that it is not *a priori* necessary to assume that the retarded and advanced trajectories in superspace are likely to meet at the Big Bang it is because, despite the uniformity of the matter distribution and the minimum gravitational entropy that characterizes the initial maximum-density state, it remains that, even under such conditions, the retarded and advanced states

could, in principle, be different in their unobservable, quantum-mechanically interfering details. This is especially true in the context where it must be assumed that the information contained in the microscopic state of the gravitational field grows with the density of matter, for reasons I have explained in section 4.7, so that the probability of an exact correspondence of the retarded and advanced states is as small during the first instants of the Big Bang as it is at any other time (given that the universe has the same information content and the same number of microscopic degrees of freedom during the first instants of the Big Bang as it has at any other time). Thus, it would only need to be assumed that a meeting of the retarded and advanced configuration space trajectories occurs at the Big Bang if it turned out that it is impossible for bidirectional time to be extended past the initial maximum-density state, as would be the case, from a classical viewpoint, in the presence of an initial spacetime singularity.

In any case, if the hypothesis that the universal causal chain must be closed is justified, then it becomes possible to confirm the validity of the conclusion stated at the end of section 5.8, to the effect that the sign of energy of the particles which can be observed to propagate backward in time (with respect to unidirectional time) in the advanced portion of history must be opposite that of the same particles which are propagating forward in time in the advanced portion of history, while their sign of action remains unchanged. Indeed, if the superspace trajectories associated with those two portions of history are smoothly joined at a remote point in the past, as well as in the future, then the particles which propagate forward in time in one portion of history must reverse their direction of propagation in time from the unidirectional viewpoint at both the past and the future bifurcation points, due precisely to the fact that they do not reverse their direction of propagation from the viewpoint of *bidirectional* time. But this means that the energy signs, which necessarily remain unchanged relative to the direction of time in which those particles propagate (given that there is no reason to assume that their action signs would reverse), would appear to be reversed in comparison with those of the corresponding particles which propagate in the opposite direction of (unidirectional) time in the current portion of history, in agreement with the conventional description of advanced-wave phenomena.

The only difference between the hypothetical phenomenon of advanced waves, as they are conventionally described, and that of the advanced portion of history discussed above, has to do with the fact that the signs of all the non-gravitational charges carried by the particles present in the advanced portion

of history would now appear to be reversed, from the unidirectional-time viewpoint, given that it is explicitly assumed that they remain unchanged, from the viewpoint of bidirectional time, when the trajectory of the universal causal chain bifurcates in the future and in the past (which should have been expected, given that the alternative definition of the time-reversal operation I introduced in chapter 3 involves a reversal of charge from the unidirectional viewpoint). This is without consequences, however, because the relevant force fields that provide the experimental conditions observed in the advanced portion of history all have their polarities reversed as well.

The circular nature of history is also what allows to explain that quantum interferences do occur, even in the context where we are assuming that the retarded and advanced portions of a quantum process actually take place at two very distant epochs along the configuration space trajectory; which would appear to imply that they can have no effect on one another (locally). I believe that if there are quantum interferences between the many possible paths allowed for the retarded and advanced portions of history, it is because the circular nature of history imposes a condition of continuity on the phase of the wave function which is equivalent to that I have identified in section 5.8 when discussing the significance of the negative probabilities which occur in the context of a time-symmetric formulation of quantum theory.

Indeed, given that what one would need to estimate, ultimately, is the probability of observing a certain history of the universe that comprises a detailed description of all the individual sub-processes (decoherent or not) which occur in the course of that history, then one must recognize that the phase of the wave function is actually a shared property of the unique superspace trajectory that provides the most accurate account of the history of the universe as a whole (or at least of its local space curvature over an arbitrarily large space-like hypersurface). But, in the context where the entire history actually consists of a closed causal chain that feeds back on itself, there actually exists a constraint which imposes that all contributions, by intermediary sub-processes, to the evolution of the phase of the wave function associated with the complete cosmological process (along the closed configuration space trajectory), be such that they allow the phase to end up, after a complete turn, in the exact same state in which it was at the point of the trajectory that constitutes both its initial and its final boundary condition. If the universe does, in effect, evolve back to the exact same state in which it once was, then it is certainly appropriate to assume that this state cannot itself be different from what it actually is, even for what regards unobservable

physical properties like the phase of the quantum wave function, otherwise the concept would have no real significance.

It had, in fact, already been realized [65] that there is a single phase associated with the wave function of the whole universe, that is equivalent to one very rapidly rotating clock hand. But in the context where time itself must be considered to constitute a periodic phenomenon, it follows that the cosmic wave function must be similar to that which applies to quantum systems submitted to periodic boundary conditions (like an electron in orbit around a hydrogen nucleus, whose wave function must necessarily involve an integer number of wavelengths). There may, thus, be something true to the previously discussed results from canonical quantum cosmology, which appear to indicate that the wave function of the universe is of the stationary kind, even if, in the present context, this no longer means that time is irrelevant to quantum cosmology. In any case, this continuity condition is what allows me to explain why it is, in effect, appropriate to impose on the unobservable phase of the wave function that it does not end up, in the course of an ordinary time-symmetric process, in a state that would be incompatible with that in which it initially was. I believe that such a requirement can be enforced in the context where certain time-symmetric histories (with which are associated negative probabilities) are allowed to diminish the probability that a process would occur by following any possible path, as a result of the influence they exert on the very conditions (both initial and final) necessary for their own occurrence.

It appears, in effect, that it is the requirement of continuity of the phase of the wave function along the closed, universal causal chain that explains that even individual time-symmetric processes have a larger probability to occur when they leave the quantum phase invariant, because it is necessary to impose on the phase that it remains unchanged in the course of each individual time-symmetric process, even if the real constraint applies to the universal causal chain that describes the evolution of the whole universe and on which is imposed the closure condition. Indeed, when one calculates the probability for an individual time-symmetric process to happen, one implicitly assumes that the observable, macroscopic conditions imposed on the rest of history are such as to leave the phase of the wave function invariant, whenever those which are imposed on the process itself leave it unchanged, because this is the only way to assess the likeliness that a process will arise independently from the rest of the history of the universe (which must necessarily be assumed to occur if the process itself does). Therefore, if the phase change

associated with the observable, conditions imposed on an individual time-symmetric process is maximally destructive, then it means that the process cannot occur, because the whole history in which the process is embedded could not itself happen (given that it would not leave the phase invariant upon a complete turn around the closed causal chain). This is why there are consequences to a variation of the phase that would take place in the course of an individual time-symmetric process, even when the condition of invariance actually applies on a global scale (of both space and time). Thus, it is now possible to understand why it is that there are quantum interferences between multiple different histories for ordinary quantum processes, even in the context where we assume that only one (circular) history actually takes place.

Those are very significant conclusions, given that, in the present context, observation of quantum interferences is the only way by which the advanced portion of history can be deduced to exist by an observer that is present in the retarded portion of history. Such an explanation of the origin of quantum interferences would also appear to confirm that quantum non-locality is a consequence of the non-trivial topology of the configuration space trajectory (which is here assumed to be the trajectory in superspace). Indeed, according to Hans Reichenbach [21] (p. 58), when faced with unexpected non-local correlations, one can either invoke ‘preestablished harmony’ in the form of instantaneous couplings of distant events that would violate the principle of local causality, or else recognize that one is dealing with a compact topological structure in which periodicity naturally arises. What I have tried to explain is that the history of the universe is just such a structure and therefore its circular nature is what most naturally explains quantum non-locality as a phenomenon involving the entanglement of quantum phases.

Now, as I mentioned in section 5.6, it has been argued by certain detractors of the earlier, more conventional time-symmetric interpretations of quantum theory that the problem with any such interpretation is that it is not possible to distinguish between situations where interferences among different histories must be assumed to exist and situations where they can actually be ignored. I have already explained that this erroneous conclusion arises merely when we fail to recognize that decoherence must occur, even from the viewpoint of a time-symmetric formulation of quantum theory, under the same conditions where it would be expected to happen according to a many-worlds interpretation, despite the fact that the phenomenon has a different meaning in the context of a time-symmetric interpretation. But, as

I mentioned in section 5.10, there are two problems that one must face before one can conclude that decoherence does, in effect, provide the mechanism by which the quasiclassical character of macroscopic phenomena arises, even in the context of a time-symmetric formulation of quantum theory.

The first of those problems has to do with the fact that it may never be possible to assume that decoherence itself constitutes a truly irreversible process. There is no reason, in effect, to reject the possibility that, given enough time, the processes giving rise to decoherence could eventually be reversed on an arbitrarily large scale, so that the many variables of the environment with which a quantum system has become correlated could be submitted to quantum interference, even without deliberate intervention and long after a measurement would normally be assumed to have occurred. Indeed, it appears that it is merely the improbability of such an evolution that explains that we do not feel compelled to recognize that measurements may not be definitive processes and could actually be overturned in the future, which would affect the validity of theoretical predictions concerning ongoing phenomena. One may be tempted to argue that this is not a real problem, because the potential for entropy growth may be unlimited in the future and this may allow one to expect that, as the effects of a measurement spread irreversibly into an ever larger portion of the environment, the possibility that quantum interferences involving all those correlated variables would occur becomes ever more insignificant. Indeed, I have provided arguments in section 4.9 to the effect that the growth of gravitational entropy may be unlimited in our universe, due to the presence of negative-energy matter, which would appear to provide support for the conclusion that decoherence is truly irreversible, even in the context of a conventional interpretation of quantum theory.

The problem, however, is that, given an infinite amount of time, even such a continuously decreasing probability may not prevent fluctuations from eventually giving rise to quantum interference on a very large scale. Therefore, it would seem that one cannot avoid the conclusion that decoherence is not definitive, which should have significant consequences at the present epoch. Now, given that I have argued in section 4.7 that the expansion of the universe does not take place with a real growth in the amount of microscopic structure or information, due to the variation of information associated with the diminishing strength of local gravitational fields, it would appear that the probability that the universal causal chain closes at some point in the future (which would happen when the exact unobservable state of all the

microscopic degrees of freedom in the universe would happen to be the same in the retarded and advanced portions of history) is not diminishing with time. This conclusion may perhaps appear to be irrelevant to the problem discussed here, but that is not the case, because what it actually means is that unidirectional time will, *by necessity*, eventually end at some point in the future, however distant this event might be. But if history does not last forever, then the probability that decoherence may be reversed on a very large scale at some point in the future, in a universe with ever growing entropy, actually becomes null. In other words, what we are now allowed to expect is that the universal causal chain will eventually close in the future, before decoherence has the chance to be reversed on a large scale, which means that decoherence does not merely eliminate quantum interferences for all practical purpose, but must be assumed to give rise to classical outcomes of measurement as a matter of principle, and this conclusion remains valid even in the context were we do not postulate that irreversibility arises at a fundamental level. I believe that this constitutes the decisive argument that allows one to make sense, at long last, of the observation that quantum measurements, once effected, produce definitive outcomes which are never overturned.

The second problem one must confront is perhaps more significant. I have explained in section 5.10 that certain relatively well-known developments [55] appear to indicate that the criterion of ‘consistency’ (in the sense of a consistent-histories interpretation of quantum theory) would not be constraining enough to allow one to expect that the quasiclassical nature of reality would persist following a quantum measurement (conceived as an irreversible process during which decoherence is taking place), even when one is allowed to assume that decoherence itself is definitive and is not likely to ever be overturned. I can now explain why it is that the realist, time-symmetric interpretation of quantum theory I have developed is more appropriate for predicting the emergence of a quasiclassical world that remains classical once the consequence of one or another outcome of a quantum measurement irreversibly propagates into the environment. What holds the key to a complete and effective solution to this particular aspect of the quantum measurement problem and to an explanation of the quasiclassical nature of ‘macroscopic’ reality is, once again, the acknowledgment that the property of closure of the universal causal chain is not optional and must, according to the arguments provided above, be imposed as an absolutely essential consistency requirement. It is only when I recognized the unavoidable nature of this condition

that I was able to understand that, in the context where there must exist both a retarded and an advanced portion to every quantum process, additional constraints exist which only become apparent during processes which can be qualified as measurements.

So, what is it, in effect, that characterizes a process that can be described as a quantum measurement? The essential ingredient of decoherence itself appears to be irreversibility (dissipation to be more specific), but as I mentioned above, decoherence can only be part of the solution. So, what happens, as a consequence of irreversibility, that does not take place under those conditions where quantum interferences exist? To answer this question, it may help to consider what would be necessary for measurement *not* to occur and quantum interference to exist, even after a quantum system becomes entangled with its environment. It is obvious that what would be required is that the state of the quantum system along with that of the immediate environment to which it has become correlated do not become entangled with an even larger portion of the environment. In other words, there would need to be no traces, in the larger portion of the environment, that would allow one to tell through which history the system and its immediate environment actually went. The point at which irreversibility enters the picture, therefore, is through the making of a record of the events involved (conceived precisely as the kind of process during which the effects of one or another of several alternative outcomes of the evolution of a microscopic system is amplified to macroscopic proportions). Only when the state of each physical attribute whose determination would allow one to tell what the history of the system and its immediate environment was is submitted to quantum interference before this information has the time to spread into the environment, can interferences actually be observed.

It would, therefore, appear that if irreversibility is, in effect, necessary for the elimination of interferences, it is because the making of a record can only occur when future evolution takes place irreversibly. What happens, when a record is produced, is that one unique cause in the past leaves multiple recognizable and mutually consistent traces of its occurrence in the future. A long-lasting record is one whose mutually consistent traces themselves each produce multiple recognizable and mutually consistent traces in the future that can all be traced back to the same unique original cause in the past. What happens when a quantum measurement, conceived as a particular, but general instance of such a record making process, comes into effect, therefore, is that a unique, particular outcome of the evolution of a quantum system

produces a recognizable effect on a multitude of other events in the future, which would all have been affected in recognizably different ways had that original outcome been different or nonexistent. What must then be responsible for the elimination of quantum interferences that follows decoherence is the fact that a growing number of observable variables become correlated with one unique, specific outcome of the evolution of a microscopic system, while *all* of those variables would have evolved differently if another outcome had been obtained for the same measurement, in the past.

Now, the important point in all of this is that the spreading of effects does not take place with respect to an arbitrarily-chosen dynamic attribute, but always relative to position space. Indeed, as I have emphasized in section 5.7, at the most fundamental level, reality appears to consist of elementary particles, which are objects that are localized in position space and which allow the propagation of effects through local contact, again in position space. There is, thus, something very particular with position space for what has to do with unidirectional causality and the irreversible propagation of effects and this is apparent in the fact that the spreading of wave fronts always occurs in position space and not in configuration space. The singular status of position space is made even clearer by the fact that the particular boundary condition which I have identified as being responsible for the asymmetry of the evolution in time of systems with a large number of independent, microscopic degrees of freedom is a condition that is imposed on the *spatial* distribution of matter in the first instants of the Big Bang. Indeed, it is the homogeneity of the spatial distribution of positive- and negative-energy matter particles in the maximum-density state of the Big Bang that allows the universe to evolve irreversibly toward a state of larger gravitational entropy, characterized by a greater inhomogeneity of the two matter distributions, as space expands, in the future direction of time. But as I explained in section 4.9, this condition is what allows one to assume that the cosmic horizon, which limits the scale on which effects were allowed to propagate since the Big Bang, actually grows with time, from the minimum value it had in the initial singularity.

What is allowed to happen, on a smaller scale, as a result of this particular condition, is for an irreversible spreading of effects into an ever larger volume of space to take place in the future direction of time, as elementary particles freely propagate in either the retarded or the advanced portion of history (this is particularly apparent in those cases where dissipation is involved). In fact, as I explained in the previous section, the same constraint of

global entanglement which gives rise to thermodynamic irreversibility is also responsible for allowing time to differentiate from the other three dimensions of space and therefore for giving rise to the causal structure of spacetime that is described by relativity theory and which is responsible for the fact that effects necessarily spread in space, either forward or backward in time. But what characterizes *unidirectional* causality is not only the fact that it operates relative to a unique dimension of space-time, but also the fact that it does, indeed, give rise to an irreversible spatial spreading of effects in the future direction of time, which is actually what the principle of local causality is usually considered to be all about. Thus, as time goes, a growing number of independent, microscopic degrees of freedom can be affected in recognizable ways by unique causes located in the past, while the reverse phenomenon is never observed to happen and this is really a property that is unique to the evolution of position states.

We are now very near a solution to a very old problem. What I have just explained is that the making of a record is the essential condition for a quantum measurement to take place and that what it entices is the production of a multiplicity of correlated effects involving very many, otherwise independent, variables which could all have evolved differently in the future, had the outcome of this measurement itself been different. A multitude of correlated effects as the outcome of one single quantum measurement. It is not very difficult to realize that, as time passes, the observable difference between the consequences of one single past measurement and what would have been the consequences of obtaining a different result for the same measurement becomes ever more significant. But in the context where one recognizes that the universal causal chain must, as a matter of principle, form a closed trajectory in superspace, then this remark becomes highly significant.

This is because, in a world that would have been quasiclassical on a macroscopic scale until now, if a measurement performed on the retarded state of a quantum system was to give rise to an outcome that is different from that which was obtained as a result of a similar measurement performed on the advanced state of the same system by a measuring device whose irreversible evolution actually also takes place in the future direction of time, then, as time goes (in the future), an exponentially growing number of independent variables from the environment of the system that takes part in the retarded portion of history would be allowed to differ from those of the same system that takes part in the advanced portion of history. This means that the two trajectories in superspace, which until now had always been very similar to

one another, would begin to diverge in a way that would actually make it *increasingly* less likely that they could ever merge with one another at some point in the future, because of this property of the record making process which is to produce an accumulation of recognizable changes in the states of an innumerable number of independently evolving degrees of freedom, as a consequence of one little change in the past.

It is the requirement of closure, that applies to the universal causal chain, that constrains the future evolution of the retarded and advanced portions of history to not diverge in any *observable* way, from the unidirectional-time viewpoint, because, if this condition was not obeyed, the number of independent variables, from both portions of history, that would need to change together in the same recognizable way at some point in the future, so as to allow a merger of the two trajectories, would become too large for the closure requirement to ever be fulfilled. As a result, the universal causal chain must be stretched into two similar trajectories evolving side by side in superspace, along the unidirectional direction of time, for the whole duration of history, as if two indistinguishable versions of history were taking place in parallel, all the time, without ever interacting with one another. But the constraint of non-divergence need not be any more restrictive than that, because what remains unobserved does not give rise to the formation of a record and has no irreversible consequences and therefore is not required to correspond, for the two portions of history, by the requirement of closure of the universal causal chain. Quantum interferences are not forbidden altogether, they merely become increasingly more unlikely as the entanglement of a quantum system with its environment becomes more significant and this is exactly what is required from an observational viewpoint.

It must be clear, however, that despite the unique role played by position in giving rise to the formation of records, the dynamic physical attribute of a quantum system that is known with perfect accuracy is not necessarily always its position. The privileged status of position space only means that even when the measured attribute is *not* position, it is nevertheless a *spatial* distribution of macroscopic constraints that allows such a measurement to be performed, because it is concerning those constraints that information is available in the form of records. This means that there is no freedom in deciding which dynamic attribute is classically well-defined in any particular situation where we have knowledge of a specific set of macroscopic conditions (while in fact such conditions are always present for one and only one dynamic attribute, as I mentioned in section 5.10). On the other hand, the

dynamic attribute of a quantum system about which only a *minimum* amount of information is available as a result of the existence of records concerning the position states of various parts of a measuring device (the environment degrees of freedom), is the attribute that may go through any possible trajectory (not necessarily in position space) during both the retarded and the advanced portions of history, thereby giving rise to quantum interferences.

What's important to understand is that, given that it is for position space observables that the making of a record of past events can take place, then it follows that the constraint of non-divergence of the retarded and advanced superspace trajectories, or more specifically of the observable retarded and advanced states of a time-symmetric quantum process, is a constraint that applies only to the dynamic attribute of a system whose state is restricted to a subset of values as a result of being submitted to experimental conditions of such a nature. But such a constraint does not only give rise to non-interfering outcomes of measurement following decoherence, but really to a quasiclassical evolution that persists in time for the same family of consistent histories (the physically relevant set of histories).

It had already been remarked, in effect, that, from a phenomenological viewpoint, decoherence, even as it is traditionally conceived, appears to select position as the relevant collective observable (that which becomes correlated with the microscopic system under study), at least for mechanical systems, in the presence of dissipation. It was conjectured that this is merely a consequence of the fact that the laws of physics (particularly in a quantum field theoretic context) are invariant under a change of reference system. In the present context, however, this could only be understood to mean that the selection of position as the relevant collective observable for decoherence is indeed a consequence of the fact that unidirectional causality (the irreversible spreading of effects) operates in position space, because what emerges, as a result of relativistic invariance, is the causal structure of spacetime, which, under appropriate conditions (when evolution is irreversible), gives rise to unidirectional causality and therefore to the existence of persistent records of past events. The fact that the phenomenon of dissipation merely consists in one particular instance of irreversible spreading of effects in position space would therefore appear to confirm that it is the closure requirement (that must be applied to the universal causal chain) that allows quasiclassicality to emerge and to persist for those attributes of a quantum system whose states are restricted by macroscopic conditions of a spatial nature.

There should be no doubt that the existence of such an objectively de-

finer, preferred basis is absolutely necessary from an observational viewpoint, because if none arose, it would be impossible to determine what causes the persistence of the quasiclassical nature of reality (even under the assumption that the universal causal chain must close at some point). Indeed, if reality was classical with respect to one family of consistent, coarse-grained histories at a given time and then relative to another such family at a later time, as would be allowed in a more conventional context, then this reality would no longer appear classical from the first viewpoint after the transformation has occurred. But when quasiclassicality is the outcome of imposing a requirement of closure to the universal causal chain and irreversibility is a feature of the spreading of effects in position space, it follows that a preferred basis (a preferred choice of dynamic attribute to represent quantum states) is naturally selected for the elimination of quantum interferences and it is from the viewpoint of the records which are available concerning the constraints (of a spatial nature) that select this dynamic physical attribute that the world necessarily appears to remain classical following a measurement. I believe that those conditions, therefore, allow to satisfy Dowker and Kent's requirement for an additional, purely quantum-mechanical principle that would allow one to select a particular set of (consistent) histories as being of particular physical significance, without having to rely on solipsistic arguments.

So here we are, having actually explained why it is that, in practice, one never observes quantum superpositions involving macroscopic states of measuring apparatuses. If we never experience histories in which a cat is alive and dead all at once (following an experiment of the Schrödinger's cat type), it is because, if it was not the case that the cat was either alive or dead in the retarded and the advanced portions of history alike, this would change the future in ways which would render impossible an eventual meeting of the retarded and advanced trajectories in (some extended version of) super-space, while this is necessary for the universe to be self-determined from the viewpoint of causality. The identified constraint simply makes it extremely unlikely (as unlikely, in fact, as the growth of entropy that took place while the retarded and advanced states became distinct is important) that such an evolution could ever be experimentally deduced to have occurred. The essential characteristic sought by Von Neumann and which would differentiate a measuring apparatus from the system it measures is simply the possibility that exists for the measuring device to generate a record of the particular evolution it goes through, which has decisive consequences in the context

where reality is a causal chain that must close at some point in the future. From that viewpoint, of course, quantum interference of macroscopic states is not completely impossible, but even if such an unlikely phenomenon was to happen, then one would not *see* a cat that is both alive and dead at the same time (despite the fact that one would then have to be in a state of superposition as well), because one is always confined to directly perceive only the portion of history (either the retarded or the advanced part) in which one happens to be located and in any such a history there is always a unique set of causally related facts.

But this does not mean that a state of superposition involving a macroscopic portion of reality would have no apparent consequences, because if the advanced state was to become distinct from the retarded state on a large scale, then the estimation of transition probabilities for future processes would be affected in dramatic ways from the viewpoint of those observers which are part of the process while it is under way, which means that their future would actually become unpredictable unless they assume that such a divergence from classicality has indeed occurred and this is how they would actually gain knowledge of the existence of such a distinction between the current retarded and advanced states. But if the condition of closure of the universal causal chain has the consequences I'm expecting it would have, then the observers which were part of such a process would not be allowed to remember through which history they went on either the retarded or the advanced portion of the process, after quantum interference is over, as otherwise this knowledge could spread into the environment²⁷.

The point that is perhaps the most difficult to understand concerning what I believe would qualify as an appropriate account of experiments of the Schrödinger's cat type, in which there would be interferences between macroscopic states, is that in the final state of such an experiment the cat would have to be neither in a live-with-no-poison-in-its-blood state, nor in

²⁷This observation cannot constitute the basis of an alternative explanation of thermodynamic time asymmetry, because, if one does not assume that there exists a constraint for the retarded and advanced states not to diverge that is made necessary by the independent condition of low gravitational entropy at the Big Bang, which from my viewpoint is responsible for time irreversibility, then one has no reason to expect that the retarded and advanced portions of history should converge back to the same macroscopic state after having diverged on a large scale and this means that our memory of the particular history that actually took place would not need to vanish and therefore its persistence would not need to be correlated with a history where entropy grows in the future.

a dead-with-poison-in-its-blood state, even though it is true that the animal may no longer exist in a recognizable form, because this is not the same as a cat that is dead due to having absorbed the poison released as a result of the measurement on the quantum particle having produced a negative result, even if it does mean that the cat may no longer be alive in the final state. What is required, therefore, is that it be impossible to tell, from the information that is present in the final state, whether the cat was killed by the poison or whether it might have been alive without any poison in its blood before the final measurement was performed that would have revealed the existence of quantum interferences, so that, even if the cat no longer exists in the final state, it would not be correct to say that it was killed as a result of the particular outcome of the particle disintegration process. In any case, given that no complex macroscopic system, such as a cat, was ever subjected to any reproducible experiment in which quantum interferences would have been observed, then it would appear that the requirement of closure, which I suggest must be imposed on the universal causal chain, is well-founded, because it does allow one to expect that macroscopic objects, which can never be completely isolated from their environment, should practically never be found in states of quantum superposition.

An additional advantage of the approach proposed here is that it allows one to understand how it is that global consistency would be enforced, in the context where a classical time travel experience would occur and the course of history could potentially be altered so as to give rise to an alternate future. Indeed, when the effects of a future measurement can be propagated backward in time (as a result of the existence of an advanced portion of history) and there is a condition for the retarded and advanced portions of history to share the same observable macroscopic conditions in the future (so that the universal causal chain can close at some point), it follows that the present can only be influenced by the future to be such as to give rise (through forward-in-time causation) to classical outcomes of measurement, rather than to a retarded state that would differ from the advanced state. In other words, the present cannot be influenced by the future in such a way that it would be likely to evolve toward a different future. Thus, even if the second law of thermodynamics could be temporarily violated in a local region of space, perhaps as a result of a formidably improbable fluctuation, and information about the future would become available, no violation of the principle of global consistency could arise that would involve observable phenomena.

It is, therefore, the circularity of the causal process and the existence of a thermodynamic arrow of time arising from the requirement that all the particles in the universe be causally related to one another which allow global consistency to be preserved in the context of a time symmetric interpretation of quantum theory. The conclusion that global consistency would always be preserved in a quantum-mechanical context, therefore, need not depend on the hypothesis that all histories are followed all at once and that a ‘splitting of branches’ occurs whenever an alternate reality is produced, as is often assumed, because it can be derived much more naturally by recognizing that for the universe to be causally self-determined, its history must consist in a closed causal chain. Yet it does seem appropriate to assume that it is quantum theory that would ultimately be responsible for the impossibility of even a classical time travel paradox, as I suggested in section 5.4, because the limitation discussed here is made unavoidable as a result of the time-symmetric nature of quantum reality, which enforces consistency on a global scale (as necessary for the existence of non-local correlations) without violating the principle of local causality.

Now, it must be clear that, even though no record of the future can exist in the context where entropy only rises in this direction of time, reality would necessarily remain quasiclassical relative to the same family of consistent, coarse-grained histories in the past direction of time as well, because if entropy rises continuously from as far back in time as the first instants of the Big Bang, then the condition imposed on future evolution by the requirement that the universal causal chain closes at some point in the future imposes that history be classical right back to the initial singularity. Indeed, if the property of quasiclassicality is required to be valid for the entire duration of history, as a result of the constraint that applies on future evolution, then despite the fact that no record of the future exists which would constrain past evolution to remain quasiclassical (in the context where the universal causal chain must also close at some point in the past), the condition will nevertheless also apply to past history, as a result of the condition that applies on the future, as long as entropy is actually growing in the future throughout this entire history. This is why we are allowed to expect that a unique past does exist that is compatible with our observation of the existence of mutually consistent records of past events.

But if bidirectional time extends past the initial singularity following a hypothetical quantum bounce, then it would be the condition that the universal causal chain closes *in the past*, prior to the initial singularity, that would

require history to remain quasiclassical in the past (and therefore, again, also in the future), right from the instant at which matter emerges from the ‘initial’ singularity, because entropy would then be growing in the past (for reasons I have discussed in section 4.9) and records would only exist about the future (which would then be similar to our past, from a thermodynamic viewpoint). Just as is the case for the future, it is not possible to say when it is exactly that a meeting of the retarded and advanced configuration space trajectories would occur in the past²⁸, because the only condition that must be imposed on the closure of the universal causal chain in the past is that it does not occur before the time at which the initial singularity is formed (on our side in time of the Big Bang), because global entanglement must have had the time to occur, as otherwise the universe would not have been allowed to exist as an entity formed of causally related elements.

Up to this point, I have only discussed the emergence of quasiclassicality as it arises in a conventional quantum-mechanical context, where the metric properties of spacetime constitute a common, unique background over which both the retarded and advanced portions of a process unfold, either with or without quantum interference, depending on whether or not the particular history of the particles propagating over this background space gives rise to the making of a record. But what right do we have to assume that the metric properties of spacetime themselves should always be shared by the retarded and advanced portions of history, if all other physical quantities can, under appropriate circumstances, differ and interfere for the two trajectories of the universal causal chain? If the other macroscopic conditions which are shared by both portions of history are so determined merely as a result of the fact that they give rise to an irreversible spreading of effects, then why would the metric properties of spacetime which are shared by both portions of history be simply given once and for all in their classical form, instead of being subjected to the same rules that govern the other physical attributes of our universe? The truth, of course, is that the metric properties of space are not always classically well-defined and that they may differ and interfere for the

²⁸One should note that it is not possible to assume that the universal causal chain closes at the Big Bang and yet that there is a history taking place in reverse, prior to the Big Bang, otherwise the meeting of the retarded and advanced trajectories in superspace would no longer have any meaning, even for the future, because, when the closure condition would be met, history could nevertheless continue to take place as if nothing had actually happened.

two portions of history.

It is already understood, in fact, that macroscopic changes to the gravitational field are a very potent way by which decoherence can be triggered, as confirmed by the fact that the motion of planets is one of the phenomena for which the absence of quantum interferences is the most conclusive and the most persistent, while it was shown that this is not unrelated to the magnitude of the gravitational fields involved. Now, I have already mentioned that, in a quantum gravitational context, what we would be dealing with are situations where the intrinsic curvature of space would be allowed to differ in the retarded and advanced portions of history. I may now add that this would occur whenever information, in the form of records, would only be available about the extrinsic curvature of space associated with the rate of change of intrinsic curvature along the actual trajectory that is followed in superspace. Indeed, the intrinsic and extrinsic curvatures of a space-like hypersurface are the quantum gravitational equivalent of position and momentum and therefore they constitute conjugate physical attributes whose states cannot be determined together with arbitrarily high precision using one unique set of experimental constraints. But this does not mean that all histories involving distinct intrinsic curvatures are followed all at once when the extrinsic curvature is known with high precision, but merely that, under such conditions, the intrinsic curvature may be different for the corresponding retarded and advanced portions of history, because information about the actual history, is available (in the form of records) only for the extrinsic curvature.

The situation we normally experience (outside the quantum gravitational regime) is one where the curvature of space in general is classically well-defined (knowledge is available about both the intrinsic curvature and its rate of change) and there are no quantum interferences arising from the curvature of space being potentially different for the retarded and advanced portions of a process (even when space is not flat locally), as is necessary for conventional quantum theory to provide a viable description of reality. But that need not always be the case and indeed, in situations where we would try to determine the extrinsic curvature of space with a very high degree of precision, by measuring the rate of change of the gravitational field over a very small time interval, then the intrinsic curvature of space would be subjected to quantum interference, as its state would no longer be constrained to be the same in the retarded and advanced portions of history, for reasons I already mentioned. Under such conditions, it would no longer be possible to estimate

transition probabilities while using one unique set of metric properties, that is to say, by assuming the existence of one single classical spacetime over which particles would propagate in both portions of a process and it would be necessary to take into account the possibility that the metric properties of space, themselves, could evolve differently in the two portions of history, along with other unobserved dynamic attributes. To determine which metric properties are likely to emerge upon observation, one would then need to take into account the existence of quantum interferences between the many possible histories of space curvature. When interference would happen to be constructive, a given curvature would have more chances to be observed and when interference would be destructive, the very boundary conditions necessary for the observation of such a curvature would be unlikely to have existed in the first place.

From such considerations, it transpires that, if time itself can be subjected to quantum interferences or superpositions, it is only in the sense that, on a sufficiently small scale, time may flow faster, or slower, locally, for the two portions of a quantum gravitational process, due to the fact that the curvature of space may not be the same in both portions of history and may therefore give rise to differing durations for otherwise similar propagation processes. But given that the constraint of global entanglement that is responsible for selecting the particular signature of the metric of spacetime that gives rise to a universally valid distinction between time and the other three dimensions of space only applies to the initial state at the Big Bang, it may be possible for the light-cone structure to be altered to such an extent that the causal order of events would be reversed along a time-like interval, on a sufficiently short scale, in the context where the gravitational field itself can be subjected to quantum indefiniteness, because the existence of closed time-like curves is only forbidden on a time scale for which thermodynamic time asymmetry is required to apply.

It is not true, though, that there is no definite space and time in the quantum gravitational regime. A unique curvature of space does exist throughout history, only, it can differ for the two corresponding portions of history along the universal causal chain, to the extent that there may, in fact, no longer be a simple correspondence between those two portions of history on a very small scale. Reality always remains a unique, closed causal chain, even though, on a smaller scale, the regularity of the particular trajectory followed by the state vector in superspace may be altered, given that the metric properties of space and the gravitational field may themselves no longer remain unaffected

by the inherent randomness of quantum-mechanical evolution, which is then allowed to give rise to a divergence of the retarded and advanced trajectories in superspace, as long as no record is available regarding what those metric properties actually are.

What is significant for a quantum-mechanical description of gravitation and space curvature, from the viewpoint of the developments introduced in the first part of this section, is that there must be a level at which the intrinsic curvature cannot remain superposed (as it necessarily is on the quantum gravitational scale) and must give rise to a quasiclassical evolution and this turning point would be determined by the availability of information concerning the metric properties of space. It is, in effect, precisely when the consequences on the propagation of elementary particles of a particular curvature of space irreversibly spreads into the environment and gives rise to the formation of mutually consistent records of a particular history, that the metric properties involved must begin to evolve quasiclassically, because the requirement of closure of the universal causal chain can only be satisfied when such an evolution is observed, just as is the case in a more conventional context and this means that irreversibility would be an essential condition for a classical spacetime structure to emerge. It is only when the state of the gravitational field becomes observable that it is no longer subjected to interference effects and that it is no longer allowed to affect the propagation of matter particles differently for the retarded and advanced portions of a process.

It would, therefore, appear that the existence of a decoherent spacetime is itself dependent on the existence of unidirectional time, which emphasizes just how important it is that there exists an independent constraint, of the kind I have previously identified, for the emergence of irreversibility, because, in a quantum gravitational context, when the irreversible character of time itself does not emerge from the underlying theory, decoherence cannot alone give rise to the classical spacetime structure. What will be very important for the argument that will be developed in the concluding section of this chapter is the observation that, if random fluctuations of the metric properties of space exist which have no observable effects of the kind that would require the gravitational field to actually have the exact same configuration in both the retarded and the advanced portions of history, then those fluctuations might be allowed to exert an unexpected influence on the propagation of elementary particles, even on a scale well above that at which gravitation becomes as strong as the other interactions.

5.13 A possible role for gravitation

I must immediately warn the reader that the developments that will be the subject of this concluding section of the present report will probably be considered more speculative than other portions of my analysis and I would not myself consider such a judgment entirely inaccurate. Yet I believe that it is important to explain what I have learned concerning how it can be that, even when a quantum system is submitted to the same macroscopic boundary conditions from one trial to another, many different possibilities are allowed for the one particular (unobserved) time-symmetric history of the system that actually occurs in the course of any given process. What's significant here is the fact that those developments are motivated by the same desire to uphold the validity of the principle of local causality that motivated the approach I followed in dealing with other problems in cosmology and quantum mechanics. Despite the fact that this discussion comes last, it is actually based on results I had obtained in the earliest portion of my research program, while I was still working on the problem of elaborating a generalized, classical theory of gravitation that would describe the interaction of positive- and negative-energy matter.

It is one of those strange turns of fate that, while I was searching for a paper in the immense science and engineering library at McGill University, at the very beginning of my research career, I came upon an article in a very old volume of an obscure research journal that sought to explain the randomness of quantum measurement results as being caused by perturbations attributable to the interaction of a quantum system with a background of gravitons present in its environment. As I now understand, this was a particular instance of *classical* hidden-variables theory which was inadequate mainly as a result of the fact that it was incompatible with the requirements imposed by quantum entanglement and non-locality. Yet, for some reason, I had the strong intuition that the idea that gravitation was involved in explaining certain aspects of the quantum-mechanical description of reality was generally valid and should be further explored. This imperative remained in the back of my mind as a guiding principle as I explored other problems in fundamental theoretical physics and even though I soon realized how such a proposal could be made viable, it is only much later that I came to understand that there is actually something unavoidable with the hypothesis that gravitation must become an integral element of a truly consistent formulation of quantum theory.

In the previous section I suggested that quantum theory, as it is currently interpreted, is incomplete, given that it does not explicitly require history to be described by a closed, universal causal chain, while, as I have explained, such a concept is essential if we are to obtain a realist theory that allows for the emergence of a maximum quasiclassical domain. But at this point, it was still possible to argue that the current formalism of quantum theory (in its most appropriate form) is compatible with this more complete version of the theory. However, given that the interpretation I have proposed is dependent on the assumption that there exists a unique reality, and in a certain sense, a unique history behind all quantum-mechanical processes, even in the presence of quantum interferences (a hypothesis which is necessary in order to maintain agreement with the uniqueness of the outcome of every quantum measurement), then it transpires that if our understanding of the theory is to be considered complete, then one cannot avoid having to examine how this unique unobserved history may come to be determined from a causal viewpoint.

What should be clear, first of all, is that, while the closure requirement that must be imposed on the universal causal chain is constraining enough to predict that classical outcomes follow measurement, the decoherence process does not select one unique outcome of measurement, but rather leaves all potentialities on an equal footing. Thus, it should be clear that it is not decoherence or the closure requirement imposed on the universal causal chain that require that only one outcome be observed following a measurement on a previously interfering dynamic attribute of a quantum system. Yet the uniqueness of the outcomes of quantum measurements is what imposes on the reality that must exist in between observations that it is also unique. The problem is that if this unobserved reality is unique, while it is also allowed to vary from one process to another under unchanged experimental conditions, then it would seem that something essential remains unexplained by the theory. As a result, here again, one must face the possibility that quantum theory is incomplete, but now in a way that would appear to require that it be reformulated. Indeed, even in the context of the realist, time-symmetric interpretation of quantum theory I have proposed, it would appear that the question of completeness can only be positively answered once one allows for a further extension of the formalism from the viewpoint of which the uniqueness of the unobserved portion of history would not constitute an additional problem, but would rather provide a hint as to what goes on when a particle propagates in the space of its unobserved physical attributes.

Those remarks become particularly significant once one recognizes that it is not possible to explain the unique and random nature of measurement results themselves by simply postulating that all the interfering branches of history, which are usually assumed to occur all at once in the same universe, remain equivalent to one another following measurement, because this would require one to assume that, despite all the evidence, history is not, in fact, unique. But I have also explained that it would be inappropriate to argue that an attribute that is indefinite, in the quantum-mechanical sense of the word, could be objectively indefinite, in the sense that it would not satisfy the requirements of scientific realism in any possible way. It would, therefore, appear necessary to conceive of a unique reality of some kind, such as that which emerges from the above discussed analysis, where, even in the absence of direct observation, particles always follow unique, but possibly different paths in the retarded and advanced portions of history. In such a context, however, the question necessarily arises as to what determines which path is actually followed by a particle in between measurements?

You may recall that I have argued in section 5.10 that the unpredictability of quantum measurement results cannot be a consequence of the measurement process itself. It rather seems necessary to assume that there is already randomness before a particle meets a detector, while it is still propagating in the two unobserved portions of history, and what remains unexplained is the variable nature of this evolution, which applies even for physical systems prepared in the exact same way. What is it, indeed, that determines the particular evolution of a certain physical attribute that takes place in between measurements and which must merely be compatible with the outcomes of those measurements? What I have realized is that, in order to answer this question, it is necessary to recognize that the current theory is merely an idealization and that it must be reformulated to give rise to a more elaborate, but statistically equivalent model, in which the unique unobserved evolution that takes place in the absence of observational constraints would be a natural consequence of the existence of fundamentally unobservable, random causes, whose existence is inevitable and does not have to be postulated on purpose in order to solve the above described problem.

A related question one may ask is whether the concept of objective chance, which is usually assumed to be implied by the fundamental unpredictability of quantum measurement results, itself constitutes an appropriate notion, in the context of a realist interpretation of quantum theory? In other words, if objective indefiniteness is to be rejected, must one also reject the associated

concept of objective chance? The conclusion to which I have arrived is that this depends on what we mean by objective chance. If we are asking whether the unpredictability of measurement results can be circumvented given a more precise assessment of the microscopic state of a quantum system, then the answer would definitely be no. But if what we understand by objective chance is the idea that the unique unobserved path of a quantum system might be ‘determined’ by nothing at all, instead of being the outcome of fundamentally unobservable causes, that is to say, if we are asking whether it is possible for an unobserved, variable feature of reality to have no identifiable cause, then the answer could only be provided in light of what we already know about reality at the level where it can be observed and by taking into account any possibility that there may be for such a variable feature to actually be causally determined (in the time-symmetric sense of the word). Only if we decide that an absence of causes is not physically unacceptable and if we can be confident that no influence exists that would provide such unobservable causes, can we argue that such a *strong concept of objective chance* is still applicable at the most fundamental level of description of physical phenomena.

It is often remarked that the concept of objective chance conflicts with common sense, but that this merely reflects another failure of our intellect to grasp the essentially distinct and counter-intuitive nature of quantum reality. Again, however, I would like to argue that this is not all there is and that, from the mere viewpoint of logical consistency, there is actually something problematic with assuming that a reality can differ and yet that such a difference need not be the result of any known cause, even of a fundamentally unobservable nature. What is easy to overlook is that, when we assume that a difference could exist that would have no ‘cause’, then we actually allow for a violation of the requirement that the physical attributes of the objects which are present in our universe are to be describable by referring only to aspects of reality which are an integral part of this universe. Indeed, if one assumes that it is acceptable, for certain variable aspects of reality which would exist beyond the observable portion of physical phenomena, to have no identifiable causes (even of a random nature) originating from within the universe in which those phenomena arise, then it may no longer be possible to avoid the conclusion that those particularities actually are the product of external intervention, which would simply mean that our universe is an incomplete instance of reality. I believe that a physical model that would offer a complete account of what happens inside any given universe must,

therefore, avoid postulating an absence of causes for physically distinct aspects of that reality. This is probably the purest form of the principle of local causality.

There is, thus, something rational in our aversion for a reality that would differ without any identifiable (even if potentially unobservable) causes, that is to say, there are good motives to doubt that a strong concept of objective chance is relevant to our description of physical reality. No distinctive feature of our universe should have as a cause ‘nothing’. If events are, in general, related, in statistically significant ways, to other events of a similar nature through what we call causality, then we are justified to expect that there should be no event that would be related to something of an entirely different nature which we call nothing, but which could actually be anything at all. That does not mean, however, that we have to reject the notion that reality is fundamentally unpredictable, as I already mentioned, because even a causally-determined world would, in the context of the existence of closed causal chains and backward-in-time causation, involve an irreducible randomness, given that the cause of an event can be influenced backward in time by this very same event, despite the fact that no information is allowed to flow backward concerning that future event (so that it necessarily remains unpredictable), as I explained in section 5.4. What this means is that, even if unobservable causes were to be found to exert an influence on unobserved aspects of a quantum process, reality would remain fundamentally random and not just unpredictable, even if it is causally determined in every way. This is the exquisite beauty of time-symmetric causality: it allows for causal determination without giving rise to complete determinism²⁹.

What allows the wave function to evolve deterministically, but only until a measurement occurs, even in the context where one must assume that the underlying evolution is of a random nature, is the fact that we are dealing with a unique reality for which what happens in the future contributes to determine what happens in the past. In such a context, the outcome of a measurement on a quantum system at time t_2 can change what happens to the system as far back as the time t_1 when the system was prepared, which allows the evolution that takes place immediately after t_1 to agree with the

²⁹Such a conclusion would seem to confirm that a time-reversal operation that would apply to the present state of the whole universe defined over a given space-like hypersurface would not necessarily give rise to the exact same history in reverse, but could potentially give rise to an entirely different and genuinely unpredictable evolution, as I suggested in section 4.6.

outcome of a measurement that takes place at time t_2 , despite the fact that this evolution is taking place randomly on a local level. Thus, the fact that the wave function evolved deterministically until time t_2 (at which decoherence took place and the state vector was reduced), is not incompatible with the hypothesis that the system evolved randomly before that measurement, because this random evolution was influenced all along by what happened at a later time, when the measurement was performed, given that, from the viewpoint of time-symmetric causality, the system is required to obey constraints which may be determined by what happens in the future, as a result of that measurement. But once the potentialities are actualized, at time t_2 , the observed outcome is only required to be compatible with what actually happened in the past and this is what explains that randomness becomes apparent. Quantum evolution is always random, but a real change is actually occurring when a measurement takes place, which makes it seem like this is where randomness originates, because right until the measurement is actually performed, multiple different outcomes are still possible and the system appears to evolve indifferently toward all those final states, all at once, and this is what makes this evolution appear deterministic, as it always happens in the same way from an observational viewpoint and must always be compatible with whatever could happen in the future.

In any case, as long as the unidentified causes which may explain the variation of the unobserved paths of quantum particles that takes place in between measurements remain unobservable, reality must remain unpredictable from the viewpoint of all observers. I would therefore object suggesting that the validity of a causal theory based on the realist conception of reality developed in the preceding sections of the present chapter would imply that the wave function provides an incomplete description of the state of a quantum system, because the wave function does provide the most complete account of how a system evolves as a result of the observable constraints exerted on it, only this still leaves us with a classically indefinite state for the physical attributes which are left unconstrained by the macroscopic experimental conditions which apply to both the retarded and the advanced portions of a process. I believe that this provides an important clue as to the nature of those unobservable random causes.

What must be clear, also, is that the existence of such unobservable causes, obeying the principle of local causality, is not ruled out by the phenomenon of quantum entanglement in the context of the realist, time-symmetric interpretation of quantum theory I have proposed, because, even

if the trajectories of two particles forming an entangled pair are separately influenced by those unobservable causes, when there is as much influence of the future on the past as there is of the past on the future it is possible for the two entangled particles to evolve so as to enforce the non-local requirements imposed on the wave function as a result of their entanglement. This is why one must differentiate such an approach from the naive, realist interpretations of quantum theory which were proposed in the past and which can be appropriately called *classical* hidden-variables theories. Here it is the very concept of an objective reality that differs in essential ways, given that we are now dealing with a universal causal chain that feeds back on itself so as to give rise to two interfering, but otherwise independent versions of history for each and every process, to which must be *independently* applied the requirement of local causality. Thus, the unobservable causes are allowed to exert different effects on the retarded and advanced portions of history along the trajectories followed by any of two entangled particles, but given that the two portions of both processes interfere with one another quantum mechanically, as a result of being part of the same closed causal chain, then it becomes possible for non-local correlations to exist between the outcomes of measurements performed on the two otherwise independently evolving systems.

From my viewpoint, the reality that is causally determined is not unique in the classical sense and this is what allows even a causal theory to agree with the requirements imposed by the quantum entanglement of distant particles, without requiring complex and arbitrary non-local mechanisms of a conspiratorial nature, in contrast with all classical hidden-variables theories. The only difference between a causal theory involving unobservable causes of the kind I suggest may need to be considered and the orthodox interpretation of quantum theory would therefore be that, from my viewpoint, not only is it possible to assume that there can indeed exist a unique reality, even in between measurements of a certain physical attribute for which quantum interferences are observed, but it is also possible for this reality to be causally determined, as all observed phenomena. One of the advantages of this particular approach would, therefore, be that it naturally agrees with a much larger body of observational evidence, which clearly indicates that when there is an effect, there usually is a cause, even if its consequences may sometimes remain unpredictable.

In the second chapter of this report I have developed a generalized framework

for relativity theory that helped confirm the validity of the hypothesis that spacetime curvature really is a consequence of the existence of an interaction. Indeed, once one recognizes that local inertial reference systems and the curvature of space are dependent on the energy sign of the particles experiencing them, then one must accept that there is no such thing as a metric structure of space existing independently from the nature of the interaction that determines its properties. Thus, through an analysis of the *quantum-mechanical* concepts of bidirectional time and negative energy, I was allowed to develop an improved, classical theory of the gravitational field, which helped confirm the validity of the hypothesis that the metric properties of space and time really are the product of an interaction. What I would like to discuss now is the possibility that a better understanding of the microscopic properties of *classical* gravitational fields, which emerges from the fact that they are the outcome of a quantized interaction, could provide the basis for a reformulation of quantum theory that would allow it to be consistent with the uniqueness of history that emerged from the realist time-symmetric interpretation I have developed. It must be clear, however, that the approach I will propose does not constitute a replacement for current quantum gravitation theories (such as loop quantum gravity), but merely provides a complementary contribution to the field, similar in scope to my derivation of the number of discrete degrees of freedom that characterize the state of matter particles which are under the influence of an elementary black hole (see sections 3.10 and 4.3) or to my explanation of the emergence of a universal time variable in the initial Big Bang state (which was discussed in section 5.11).

The approach I have followed is actually the exact opposite of one that would be based on the many-worlds interpretation of quantum theory, because, instead of positing a deterministic evolution involving multiple simultaneous histories, I'm assuming a random evolution involving one causally-determined history (forming a closed causal chain). Thus, from my viewpoint, one no longer needs to assume that reality is deterministic from a theoretical viewpoint, but random from an observational viewpoint, which, all by itself, certainly constitutes significant progress. In fact, it is well-known to specialists that the many-worlds interpretation of quantum theory suffers from an additional inconsistency, which is associated precisely with the hypothesis that in general no unique outcome follows measurement. The problem is that, when all potentialities are actualized together (in the same universe), it seems that the concept of outcome probability becomes meaningless, while quantum theory is all about probabilities and nothing else,

which would appear to make this usually favored approach useless.

In any case, what should be clear already is that, if quantum systems do not always go through the same unobserved path when submitted to the same macroscopic conditions, this can only mean that, even when an optimal experimental characterization of the evolution of a physical system is available, it necessarily leaves aside fundamentally unobservable, but causally significant aspects of the process. It is only the fact that, traditionally, it appeared impossible to assume the existence of such causes without allowing violations of the principle of local causality to occur that explains that we came to believe that such an otherwise more consistent viewpoint was no longer viable, even though a realist, time-symmetric interpretation of quantum theory of the kind I have proposed actually makes such an approach perfectly sensible. Indeed, once one recognizes that, as a matter of principle, no information could ever be obtained concerning the causes which may determine the random paths of unobserved dynamic attributes, then one must conclude that no violation of the uncertainty principle could occur as a result of the existence of such causes. It is also only under the incorrect assumption that additional information could be obtained about this unobserved layer of reality (that is not already accounted for by the quantum state of a system), that one would have to conclude that information may no longer be conserved and that violations of the second law of thermodynamics may arise.

Now, even though it has long been my opinion that both the classical theory of gravitation and quantum field theory must be altered *prior* to being integrated into a quantum theory of the gravitational interaction, it is only after I realized that our understanding of *classical* gravitation leaves aside important aspects which cannot be ignored in a quantum mechanical context that I began to appreciate the fact that the quantization issue does not concern merely the general theory of relativity, but that its resolution probably requires that quantum theory itself be reformulated so as to take into account those properties of the gravitational field which arise as a consequence of the very quantum-mechanical nature of this interaction. Thus, while I do recognize that the classical theory of gravitation must be subjected to a quantization procedure on the scale at which this interaction becomes as strong as the other known interactions, I also believe that the quantized nature of gravitation would have consequences on a much larger scale where this interaction can still be appropriately described by using the approximation of a continuous force field associated with the curvature of spacetime. To be

more specific, I believe that conventional quantum theory must come to integrate a certain element of spacetime curvature, even under those conditions where we currently assume the existence of a flat and invariant spacetime. For that purpose gravitation must no longer be assumed to merely be involved in defining a constant and uniform spacetime background, but must be understood to exert a random influence that contributes to determine the unobserved paths followed by elementary particles, even in the absence of local matter inhomogeneities, as long as it remains impossible to tell, even based on the results of subsequent measurements, what was the actual path taken by a particle as a consequence of those perturbations.

Those requirements can be fulfilled once one acknowledges that the trajectory of the universal causal chain in superspace, that describes the evolution of the intrinsic or extrinsic curvature of space for the whole universe, can differ for the retarded and advanced portions of history, as a result of the fact that unobservable local fluctuations of the classical gravitational field, which are attributable to the very quantum mechanical nature of this interaction, are affecting the trajectories of matter and radiation particles in the space of those dynamic physical attributes which are not the subject of direct observation. What must emerge, therefore, is a theory where spacetime does not merely provide an additional set of macroscopic constraints, as a result of its observable nature and the irreversibility of its effects, but where the local inertial reference systems may be allowed to fluctuate in unobservable ways that may differ for the two time-reversed portions of a process, which are otherwise submitted to the same macroscopic conditions. The important point, here, is that, even though such fluctuations would indeed remain unobservable, they would nevertheless be physically significant, given that they would actually allow to explain what determines the unique and possibly distinct paths which are followed during the retarded and advanced portions of every quantum process by the dynamic attributes of a quantum system which are not submitted to observation.

From that viewpoint, even the classical spacetime continuum over which the unobservable paths of quantum particles are assumed to unfold in the absence of local matter inhomogeneities would no longer constitute a uniform and static background, but would actually fluctuate as much as the particle trajectories themselves. This proposal is merely an extension of the general-relativistic idea according to which it is no longer possible to speak of a situation where there is an absence of gravitational field. Indeed, Einstein himself reflected on the irrelevance of such a notion by noting that even in

those situations where the metric is Euclidean and no mass is present nearby, there is still a gravitational field, only it is a field that does not vary with position (while an absence of gravitational field would require that there exist no metric properties at all). Here, the idea is that, even when it would appear, from a superficial, macroscopic viewpoint, that the gravitational field does not vary with position, in fact it still exerts a decisive, random influence on the trajectories of elementary particles depicted in the sum-over-histories formulation of quantum theory.

What makes such a hypothesis necessary is the fact that even conventional quantum field theory implicitly takes into account the existence of gravitational interactions all along the unobservable trajectories of quantum particles, given that it assumes the relevance of a well-defined spacetime background over which the matter particles propagate. But, as Lee Smolin once remarked, it is difficult to imagine how a dynamical theory of spacetime (such as a background-independent quantum theory of gravitation) could actually be derived from a theory where the geometry of space is assumed to be fixed (such as conventional quantum field theory). What I'm suggesting is that, once we recognize that gravitation exerts a decisive influence, even on the scale of ordinary quantum theory, then it is also necessary to recognize that, under such conditions, the gravitational field is not constrained to have the properties of uniformity and constancy that we usually attribute to it in the absence of local matter inhomogeneities (when large measures of action are involved and classical physics is a suitable approximation). The gravitational field definitely is omnipresent and does exert an effect at every 'point' along the unobserved trajectories of elementary particles (including gravitons), but I believe that what the random nature of the paths followed in the unobserved retarded and advanced portions of any process indicates is that this classical gravitational field cannot be required to be completely uniform and to evolve deterministically on a microscopic scale, but must rather be allowed to fluctuate in ways that could differ for the retarded and advanced portions of the process, when the existence of those random fluctuations would have no observational consequences.

If there is no valid motive to reject the possibility that the gravitational field may so fluctuate in the absence of observations, then what one would have to recognize is that it is the local inertial reference systems which are allowed to vary unpredictably with position and time. In fact, I believe that this should have been expected, even independently from any consideration of a quantum-mechanical nature, given that, from a Machian viewpoint, local

inertial reference systems are an effect of the gravitational interaction of the particles experiencing them with the ensemble of matter in the universe and such effects must necessarily vary unpredictably with both position and time, as the matter distribution itself is not perfectly unchanging and uniform over the entire universe and throughout history, even if, on the average, such fluctuations should necessarily cancel out, due to the large number of individual interactions involved. What happens is that, when the trajectory of a particle is the outcome of multiple, near simultaneous, quantized interactions, such as is the case with ordinary Brownian motion, then there necessarily arise fluctuations in the number of those interactions that produce a momentum variation in one direction, that are not necessarily matched by the fluctuations that simultaneously occur in the number of those interactions that produce a momentum change in the opposite direction and this must give rise to small variations in the equilibrium of forces acting on the particle (which would here be the gravitational forces that determine the local inertial reference systems).

Those considerations are particularly significant in the context where, as I have explained in section 2.6, the absence of gravitational interactions with the matter that is missing in the direction of a void in an otherwise uniform matter distribution can actually have a considerable influence on the motion of matter particles, even if that is not always recognized. Thus, if the local inertial reference systems which determine the trajectory of a particle with a given sign of energy must ultimately be conceived as being the outcome of an equilibrium in the sum of gravitational forces attributable to the interaction of this particle with all the matter in the universe with the same sign of energy, as I explained in section 2.4, then one is certainly justified to assume that the unobservable trajectories of elementary particles should be randomly influenced by the presence of fluctuations (also unobservable) in this equilibrium of gravitational forces, given that gravitational forces are themselves conveyed by elementary particles and must, therefore, fluctuate. The crucial point is that this would be true even in the context where the approximation of a classical spacetime continuum would still be valid (and the metric would remain Euclidean locally), given that we are not concerned here with individual quantum interactions, but with fluctuations in a very large number of such interactions taking place nearly simultaneously.

What makes it possible for such unobservable fluctuations in the local equilibrium of inertial gravitational forces to have decisive consequences on the evolution of quantum systems, even outside the quantum gravitational

regime, is the fact that, even though the gravitational interaction is very weak, inertia, as a gravitational phenomenon, constitutes a very significant influence for elementary particles, given that it is the outcome of the gravitational interactions which are taking place between a given particle and all the other matter particles present in the universe, whose number largely compensates the very small probability that the particle absorbs or emits a graviton in the course of an ordinary quantum process. In such a context, it would appear that it is merely the fact that the random fluctuations of the classical gravitational field which determine the unobserved trajectories of elementary particles cancel out on the scale of action at which ordinary quantum theory itself becomes irrelevant that allows the metric properties of spacetime to be described as deterministically-evolving under ordinary circumstances.

If those considerations are valid, it would then mean that what one needs to formulate is a realist time-symmetric version of stochastic gravitational field theory (based on the generalized gravitational field equations introduced in section 2.15 and in accordance with the requirement of closure of the universal causal chain) that would apply to the retarded and advanced portions of every quantum mechanical process, independently. For this purpose, it is necessary to recognize that the constant and uniform gravitational field which is assumed to exist in the absence of local matter inhomogeneities merely constitutes an approximation that must emerge from a more accurate description where unobservable random fluctuations are present all along a particle trajectory. The classical description can, therefore, be expected to break down on the action scale associated with ordinary quantum phenomena, where random fluctuations of the metric properties of space are unavoidable. What explains that such fluctuations can usually be ignored is the fact that it is precisely on such a scale that they can be expected to remain unobservable, while they must cancel out, for the most part, when larger measures of action are involved, which could have revealed their existence. A more adequate formulation of quantum field theory, that would integrate the semi-classical description of gravitational fields envisaged here, would allow for random fluctuations in a medium that remains classically well-defined locally and would only break down on the quantum gravitational scale, where it can be expected that the approximation of a classical spacetime continuum is no longer valid.

This means that there are actually three levels of applicability to a theory of the gravitational field, because the intermediary, semi-classical level, where gravitation is usually assumed to be irrelevant, actually also involves

this interaction in a decisive way. On such a scale, gravitation may already be considered to merge with quantum theory, but merely in the sense that fluctuations of a quantum-mechanical origin must now apply to the classical gravitational field, while quantum evolution becomes causally determined as a consequence of the very gravitational nature of the forces that determine the local inertial reference systems to which elementary particles are submitted, even in the absence of observable, local perturbations of the curvature of spacetime. What makes this hypothesis significant is the universal nature of the gravitational interaction and the fact that it is allowed to affect not only the propagation of all matter particles, but also that of the particles which mediate their interactions, including the gravitational interaction itself, without having to refer to a pre-existing background structure, given that this is the interaction that determines the very metric properties of the spacetime over which all particles propagate.

Now, if local fluctuations of the metric properties of spacetime actually occur, which remain unobservable, then they would have effects which would be indistinguishable from temporary violations of the conservation of momentum and energy, given that energy would be exchanged with the gravitational field, that would be unaccounted for classically. I believe that this is what explains that virtual processes, like ordinary particle-interaction processes, involve such violations of energy and momentum conservation, which are allowed to occur merely as long as they remain within the limits of quantum uncertainty, that is to say, as long as they remain unobservable (even though they are absolutely necessary to explain the kind of phenomenon involved). Indeed, even from a semi-classical viewpoint, the reality of a particle's existence may depend on the presence of a local gravitational field or acceleration (think about the Unruh effect for instance) and in such a case all that matters is that once the presence of a particle is actually measured by a detector, even when this is made possible as a result of an exchange of energy with the gravitational field, then this event must become an established fact that is not dependent on the position or the state of acceleration of an observer. What's different, from the viewpoint of the approach advocated here, is that the virtual particles mediating an interaction can now be considered to be as real as ordinary matter particles, because what differentiates them is merely the fact that they do not exist permanently, with invariant energies, but merely as a result of energy exchanges with the randomly fluctuating classical gravitational field.

Anyhow, if the non-measurable violations of energy and momentum which

are allowed by quantum indeterminacy are taking place as a result of undetectable exchanges of energy with the fluctuating gravitational field, this would explain why it is that only the conservation of energy, momentum, and perhaps also angular momentum is allowed to be violated in such a way, while the electric and other non-gravitational charges of elementary particles (the static attributes) are always rigorously conserved, despite quantum uncertainty. In this context, the fact that the quantum indefiniteness associated with the position of a particle diminishes with the magnitude of its momentum would appear all the more natural, given that a particle with a larger energy can be expected to interact with more gravitons all at once and therefore to be less affected by individual interactions, as if it was experiencing a reduced level of fluctuation in the equilibrium between the forces attributable to all such interactions (which may help explain why the variation of the quantum phase associated with the propagation of elementary particles is dependent not only on the energy of the particles involved, but also on their mass, even when gravitation is the only macroscopic constraint involved, as is the case in the context of the classical neutron interferometer experiment in a gravitational field).

From my viewpoint, the fact that quantum indefiniteness in momentum rises as we consider increasingly smaller regions of space would also appear all the more natural, given that it can be expected that fluctuations of the classical gravitational field would rise as we consider shorter space intervals, over which the quantized nature of the gravitational field becomes more apparent, until we reach the Planck scale where (as I explained in section 3.10) every matter particle is submitted to the gravitational field of an elementary black hole and momentum (actually the direction of space intervals associated with acceleration) is totally undetermined (given that it can be either positive or negative, but with maximum magnitude in both the retarded and the advanced portions of a process). To avoid confusion, however, it is necessary to understand that, despite the fact that the degree of randomness to which are submitted elementary particles as a result of the existence of unobservable fluctuations in the gravitational field may depend on the magnitude of their energy (which determines the frequency of the wave function), interference effects cannot be considered to be an aspect of the gravitational field itself, because fluctuations in the metric properties of space merely explain what determines the unique trajectories which are followed in the retarded and advanced portions of history, while it is still the constructive or destructive nature of the interferences associated with a complete time-symmetric pro-

cess which determines whether those trajectories are likely to be followed, when one takes into account the requirement of invariance imposed on the phase of the wave function, in the context where the universal causal chain must be assumed to be closed.

An interesting outcome of this particular approach is that it allows one to more easily understand why it is that photons and other massless particles are allowed to have non-measurable (but theoretically mandatory) velocities larger or smaller than the normal speed of light in a vacuum and to travel along curved trajectories on a small scale (as Feynman diagrams for radiative corrections so appropriately illustrate), because when one takes into account the existence of unobservable, local fluctuations of the metric properties of space, it is still possible to assume that massless particles in a given energy eigenstate always travel along straight lines at their normal c velocity, as long as one recognizes that this propagation takes place along the geodesics of a locally curved spacetime. This is made possible in the context of the generalized gravitation theory I proposed in chapter 2, where matter configurations may exist which give rise, not to gravitational attraction and an apparent diminution of the speed of light (as a result of local space contraction), but to gravitational repulsion and an apparent increase of the limiting velocity experienced by massless, positive-energy particles (as a result of local space dilation). From such a perspective, the multiple possible trajectories of unobserved, dynamic quantum attributes would simply be the causally-determined geodesics of a randomly-fluctuating dynamical spacetime, rather than the random paths of particles evolving over a static and uniform spacetime background.

It must be emphasized that what I'm proposing is not that there arise stochastic perturbations of the Schrödinger equation itself, when a measurement takes place, as is sometimes proposed in order to try to explain the emergence of classicality and the random nature of quantum measurement results. Once again, it should be clear that the irreducible randomness of quantum processes cannot be assumed to be a consequence of what goes on during measurement and the improved framework, that would allow to reproduce the statistical predictions of the current theory, would differ merely in that it would allow to explain what determines the particular, unobservable trajectories followed by elementary particles in between measurements, while the absence of interferences that follows quantum measurement would still be a mere consequence of decoherence, now enforced by the requirement of closure of the universal causal chain. From that viewpoint, it would also

be decoherence and the closure requirement that would trigger the process of state vector reduction which can be expected to arise when a quantum superposition of randomly fluctuating spacetime curvatures develops that leads to specific consequences of an observable nature. As long as the randomly fluctuating curvature of space that exists in a certain retarded portion of history and that which exists in the related, advanced portion of history remain without observable consequences, they are allowed to differ as any other dynamic attribute of a quantum system. Such differences in the curvature of space may trigger decoherence, like other divergences between the retarded and the advanced portions of a process, but it is precisely the fact that they are not required to do so when they do not give rise to any observationally recognizable effect, that explains how it is possible for the unique history that takes place in between measurements of a dynamic attribute of a quantum system to vary from one virtual process to another, even under identical experimental conditions.

To sum up, I believe that, instead of simply adapting the current classical theory of gravitation to accommodate the quantum-mechanical nature of reality, we should first redefine the foundations of quantum theory to take into account a certain overlooked, but unavoidable aspect of a consistent semi-classical theory of gravitation, that would allow to explain how the unique histories, which constitute a basic feature of the formalism of quantum field theory, are randomly determined from the perspective of time-symmetric causality. It is important to mention, however, that the idea that the unique retarded and advanced portions of history which take part in every quantum process are influenced by unobservable local fluctuations in the classical gravitational field is not absolutely necessary for the validity of the solutions I have provided to other aspects of the problem of the interpretation of quantum theory. Thus, as I came to realize, this proposal is not even necessary to solve the problem of the objectification of quantum measurement results, because even if the unique reality behind interfering quantum-mechanical histories was not causally determined in every way, it would nevertheless remain unique from a time-symmetric viewpoint, which is sufficient to make it compatible with the uniqueness of the outcomes of quantum measurements. But while we may never be able to directly confirm that the unobserved quantum paths, despite their absolutely unpredictable nature, are nevertheless causally determined in every way, the fact that it is already possible to envisage the exact form of a theory that would satisfy those consistency requirements should encourage us to recognize that the only reasonable con-

clusion is that reality is not fundamentally without causes.

Chapter 6

Conclusion

I have come a long way since first asking what would happen to a negative-mass object dropped in the gravitational field of the Earth. Yet I was able to confirm that my early intuition was right and that consistency dictates that the negative mass would need to ‘fall’ upward, despite the fact that this goes against current expectations. This is a conclusion for which I have provided ample justification and even if that was all I had been able to establish, I would already be very satisfied with the outcome of my undertaking. But several other developments were introduced in this report which are all related to the issue of time directionality as a concept independent from the thermodynamic arrow of time. In fact, the hypothesis of the existence of a fundamental time-direction degree of freedom has become the vital lead that allowed me to better understand many aspects of gravitational and quantum physics. Yet, despite the fact that, originally, the main objective of this report was to provide a consistent account of the way by which the concept of negative energy that emerges from those considerations can be integrated into a classical theory of gravitation, I have also made use of those theoretical developments to provide solutions to various specific problems in cosmology and to develop a more adequate interpretation of quantum theory.

First of all, using the proposed description of negative-energy matter particles as consisting of voids in the positive-energy portion of zero-point vacuum fluctuations, I was allowed to show that the existence of negative energy-matter would not give rise to catastrophic vacuum decay and to the creation of energy out of nothing, even in the case of those negative energy states which are already predicted to occur under exceptional circumstances by conventional quantum field theory. This is one clear benefit of the ap-

proach favored here, in the context where the existence of those negative energy states must be recognized as unavoidable, even from the viewpoint of a traditional interpretation. The prediction of an absence of vacuum decay can be considered as one of the most significant result of the alternative approach to classical gravitation which was developed in this report, because this potential problem is usually considered to be the most serious affecting the hypothesis that negative-energy states can be occupied.

An important outcome of my reformulation of the discrete symmetry operations, on the other hand, was the derivation of an exact binary measure of entropy for the matter contained within the event horizon of an elementary black hole. This result is particularly noteworthy in that it actually matches the constraints set by the semi-classical theory of black-hole thermodynamics. The possibility that is allowed, in the context of the proposed interpretation of negative energy states, to generalize the analogy between classical, thermal equilibrium states and black holes, through an application of the concept of negative temperature, allows to confirm the relevance of the concept of negative-energy matter for gravitation theory. Those unexpected benefits come in addition to the solutions which were offered, in the second chapter of this report, to the more traditional problems usually associated with the concept of negative-energy matter and which allow to demonstrate the viability of a generalized gravitation theory based on the proposed, alternative interpretation of negative energy states.

But while the most significant result derived in this report will probably remain the elaboration of a quantitative framework which generalizes relativity theory in a way that increases its simplicity, rather than adding in complexity over the already elegant gravitational field equations, the most concrete results are those which were obtained by applying the lessons learned while solving the problem of negative energy states to address several long standing issues in theoretical cosmology. I believe that what this shows is that a cosmological model based on a consistent theory of negative-energy matter provides a fertile ground for understanding all sorts of astronomical phenomena in which the gravitational interaction plays a crucial role. This appears nowhere more clearly than in the case of the cosmological-constant problem. Indeed, using the proposed formulation of the gravitational field equations, I was able to show that the cosmological constant, conceived as an average, residual value of vacuum energy density, can be expected to be as small as is the difference that may exist between the scale factor determined using the metric properties of space currently experienced by positive-energy

observers and the scale factor which would be determined using the metric properties of space experienced by negative-energy observers.

What makes this possible is the fact that additional contributions to vacuum energy density, arising from those zero-point fluctuations which are directly experienced only by negative-energy observers, must be taken into account, which allow the natural value of the cosmological constant to actually be zero, rather than the very large number associated with the energy scale of quantum gravitational phenomena which is produced by more conventional estimates. It remains, however, that, in the context where energy must be assumed to be null for the universe as a whole, a non-zero initial value for the energy of matter would give rise to a non-zero value for the gravitational energy of the universe that would require space to be curved on a global scale and that would make the cosmological constant arbitrarily large in the very first instants of the Big Bang. This is why it is now possible to explain the fact that space is still perfectly flat and the observation that the current value of the cosmological constant is relatively small as being unavoidable requirements of the weak anthropic principle. In such a context, it appears that the presence in the primordial universe of negative-energy matter particles, described as voids in the positive-energy portion of the vacuum, is observationally confirmed, given that it is required to balance the initial matter energy budget, while allowing gravitational energy itself to be null, independently, for positive- and negative-energy observers, so that the rate of expansion can be set to its critical value in the initial Big Bang state. This constitutes a further proof that the alternative concept of negative-energy matter which I proposed, based on independent motives, is fully justified, even from a purely empirical viewpoint.

I must admit, however, that, for a while, I was not fully convinced that a solution as technically (although not necessarily conceptually) simple as that which I had derived (based on the hypothesis of the existence of negative-energy matter) could alone solve such a complex and difficult problem as that of flatness. What I had realized, of course, was that, if I was right, then it probably meant that inflation theory could no longer be invoked to solve other aspects of the inflation problem either and this was difficult to believe, given that inflation theory was the dominant paradigm for cosmology at the time when I obtained my first results. But I came to recognize that this is the only appropriate conclusion and that there actually exists a more natural solution to the flatness problem that merely requires one to acknowledge the reality of negative-energy matter. Thus, even though I may have preferred arriving

at a different conclusion, there is no longer any doubt in my mind that it is really the condition of null energy (imposed by the constraint of relational definition of physical attributes) and the balancing effect of negative-energy matter which allow to explain the flatness of space on the largest scale in the context where one recognizes that an observer must be present in the universe to measure a value for this parameter.

One amazing consequence of the particular approach to the problem of dark energy which was proposed in the present report is that, despite the fact that it relies on the existence of a previously ignored symmetry principle, it nevertheless allows one to understand why it is that the current value of the cosmological constant is not perfectly null. But it also allows the density of the uniform portion of vacuum energy to vary with time, because the rate of expansion of space measured by positive-energy observers can differ from that which is measured by negative-energy observers under conditions where a lesser proportion of negative-energy matter survives the early period of annihilation of matter with antimatter. But given that, in the end, it appears that the correct form of the vacuum-energy term that enters the generalized gravitational field equations does not require a non-zero average value of vacuum energy to exert an influence on the specific rates of expansion that would give rise to an even larger measure of vacuum energy, then we can avoid the conclusion that despite its small initial value the cosmological constant should become arbitrarily large in the future.

Also of importance is the conclusion that there must have existed additional gravitational attraction on visible, positive-energy matter overdensities in the early universe from the presence of underdensities in the distribution of negative vacuum-dark-matter energy. What makes this conclusion particularly significant is the fact that those forces can be expected to have accelerated the formation of large-scale structures in the primordial distribution of positive-energy matter in a way that actually helps explain the presence of well-developed galaxies at an epoch when they should not yet exist according to the currently favored cosmological models. The conclusion that no such a contribution to structure formation should exist for negative-energy matter on stellar and galactic scales, in the context where an absence of baryonic negative-energy matter is both theoretically possible and observationally unavoidable (while it would appear to be required for positive-energy matter underdensities to form on such a scale), is also of particular importance, given that it allows to explain why no observation that would provide an unmistakable evidence for the presence of gravitationally repulsive structures has

ever been performed.

It is quite remarkable, as well, that the additional source of gravitational attraction which arises from the presence of negative-energy matter underdensities in the early universe is allowed to so adequately complement the contribution to gravitational instability which is provided by ordinary dark matter, once the missing-mass effect, which is usually believed to arise solely from the presence of weakly-interacting massive particles, is understood to actually be a consequence of the presence of local variations in the density of vacuum energy, whose presence is attributable to the fact that opposite-energy observers experience different metric properties of space in the presence of matter inhomogeneities. But given that a clear distinction nevertheless exists between ordinary matter and vacuum dark matter, due to the fact that the presence of ordinary matter is equivalent to both missing vacuum energy and missing vacuum charge, while vacuum dark matter arises merely from a local divergence of the measures of maximum positive and negative vacuum energy densities effected by opposite-energy observers that do not affect the electrical or non-gravitational neutrality of the vacuum, then it is possible for those two concepts to coexist without any ambiguity. The conclusion that a much smaller portion of the missing-mass effect can be attributed to the existence of baryonic dark-matter particles carrying reversed bidirectional charges, is also significant, especially since it allows to provide an additional contribution to the fundamental, binary degrees of freedom that characterize the state of matter particles on the quantum gravitational scale, thereby allowing to explain the fact that what once appeared to constitute a fundamental unit of surface actually contains four Planck units of surface, each of which can now be associated with one discrete, quantum-gravitational degree of freedom.

Another significant outcome of the existence of matter with reversed bidirectional charges is that it becomes possible to explain exactly how the imbalance which exists in our universe between the number of ordinary matter particles and that of ordinary antimatter particles is allowed to arise. This can be achieved once we recognize that the condition of continuity of the flow of time along an elementary particle world-line, which is responsible for the absence of non-gravitational interactions between particles with opposite bidirectional charges, does not prevent those same particles from interacting gravitationally with one another. As a result, particles propagating opposite bidirectional charges in the same direction of time can be produced by pairs out of gravitational radiation, under the extreme conditions which existed

during the Big Bang and when a violation of time reversal symmetry exists, this allows more particles than antiparticles to be produced that can survive the early annihilation of matter with antimatter. Thus, it may be possible to avoid the conclusion that there would exist an absolutely characterized lopsidedness of the universe with respect to the direction of time as a result the violation of time reversal symmetry that is responsible for the presence of baryonic positive-energy matter, because, if a quantum bounce occurs, a similar asymmetry can be expected to arise relative to the opposite direction of time, in the portion of history that precedes the initial state of maximum matter density.

It is while I was trying to solve the mystery of the thermodynamic arrow of time, however, that I was led to derive the most surprising results regarding classical cosmology and to gain the essential insights which allowed me to solve virtually all remaining aspects of the inflation problem. First of all, I provided decisive arguments to the effect that temporal irreversibility is not a matter of viewpoint, because the growth of entropy can be characterized in an objective way, due to the existence of the natural definition of coarse-graining that is provided by the macroscopic parameters associated with black-hole event horizons, even when it is recognized that information is always rigorously conserved. But I also explained that there exists a usually ignored measure of information, concerning the microscopic state of the gravitational field associated with a uniform matter distribution, and that its value diminishes when the density of matter decreases below its average cosmic value locally. It is this variation that allows the total measure of information in the universe to remain constant, even in the context where the amount of missing information required to describe the microscopic state of the gravitational field must be assumed to rise as a result of the growing strength of local gravitational fields that takes place when the density of matter is increasing locally. As a result, it becomes possible to conclude that information is always conserved, despite the growth of inhomogeneities, while a similar conclusion applies in the context where expansion itself contributes a growing amount of information by continuously creating new quantum gravitational units of space in the vacuum.

Based on the notion that the thermodynamic arrow of time originates from the smoothness of the initial distribution of matter energy at the Big Bang, I was then led to propose that it is the requirement that there must exist causal relationships between all the elementary particles which are present in the universe that explains the asymmetric character of the growth of

gravitational entropy. Most people no longer hesitate to recognize that the physical properties of our universe are constrained to a very small subset of potentialities by the requirement that those properties must allow for the spontaneous development of a conscious observer. What I have proposed is that solving that oldest of all physics problems, the mystery of the origin of the arrow of time, requires taking into consideration the similarly obvious requirement that, for the universe itself to exist as a consistent whole, a certain requirement must be met, which can only be satisfied when the universe goes through a state of maximum density and minimum gravitational entropy at least once during its history, because it is only under such conditions that all of its components can actually become causally related to one another. Thus was solved that long-lasting puzzle. I believe that this unexpected outcome illustrates better than anything else the fact that serious consequences may follow when we choose to uphold, without good reasons, the validity of certain commonly held hypotheses, such as the absolute positivity of energy and the purely attractive nature of gravitational interactions, because it is as a consequence of not having been held by such a prejudice that I was allowed to solve the problem of the origin of the arrow of time.

What I'm most satisfied with having achieved, however, is having been able to actually understand quantum theory. Indeed, when I began doing research in fundamental theoretical physics, some 30 years ago, I did not suspect that some of the early ideas and insights I was trying to develop would eventually become essential for producing a consistent interpretation of quantum theory. But the hypothesis that the gravitational interaction is symmetric under exchange of positive- and negative-energy matter turned out to be indispensable to the formulation of an interpretation of quantum mechanics in which no implicit or explicit assumptions contradict one another or some observable aspects of reality, because this idea is what allows one to understand how it is possible for thermodynamic time asymmetry to emerge despite the time-symmetric nature of causality. Indeed, outside the context of a generalized gravitation theory compatible with this essential condition, it would be meaningless to assume that there must be a constraint on the emergence of a maximum quasiclassical domain imposed by a requirement of closure of the universal causal chain. Actually, it wouldn't even be possible to assume that there exists a universal time variable along which the causal chain can unfold. As a matter of fact, if we were to ignore those theoretical developments, it wouldn't be possible to assume that there is a reality at all in the absence of measurement, unless we are willing to reject some

equally unavoidable theoretical requirements derived from observation, like the principle of local causality.

Now, the most significant aspect of a quantum-mechanical description of reality is certainly the use of interfering probability amplitudes in place of conventional probabilities (or equivalently the appearance of negative probabilities for time-symmetric histories). But in the context of an interpretation of quantum theory that satisfies the requirement of scientific realism, the existence of interference effects can be understood to be a consequence of the circular nature of causality that is associated with the requirement of closure of the universal causal chain. I believe that this, again, serves to demonstrate how crucially important it is to acquire a proper understanding of certain purely cosmological aspects of reality in order to develop a consistent interpretation of quantum theory. This dependence is further emphasized by the fact that gravitation may ultimately be involved in explaining what causally determines the one unique history which is followed in the course of any particular time-symmetric process, given that a more accurate understanding of the phenomenon of inertia, which also arises from considerations of a cosmological nature, may require that quantum field theory be reformulated so as to accommodate the randomly fluctuating nature of the classical gravitational field and in such a way relieve the theory from its dependence on the concept of a constant and uniform spacetime background.

It is also the notion that causes cannot be restricted to propagate only in the future direction of time, as the classical principle of causality would appear to require, that made unavoidable a picture of quantum reality involving two corresponding, time-reversed, but non-interacting histories for each process. The understanding that this is made necessary when all causes are required to belong within our universe then made possible the elaboration of the first complete solution to the quantum measurement problem. Indeed, I have explained that it is the requirement of a relational description of reality that imposes a condition of continuity on the universal causal chain which can be most naturally satisfied when causality is appropriately conceived of as a circular phenomenon in which time plays a role similar to that which would be played by space in a closed universe, while such a closure requirement is what allows to explain the persistence of quasiclassicality following decoherence. In such a context, it becomes clear that there is no real difficulty associated with the assumption that there does exist a unique reality at all times, as long as one recognizes that this reality does not consist of one single classical history propagating in one single direction of time at

all times. Once this is understood, it is no longer necessary to retreat into complicated and confused philosophy in order to try to explain the simplest and most elementary phenomena which are taking place right in front of us, all the time.

While I was progressing toward a better understanding of quantum mechanics, I realized that my position concerning the many-worlds interpretation of quantum theory is somewhat similar to my position regarding the weak anthropic principle. Indeed, while I do believe that both anthropic selection and the existence of a multiplicity of causally independent universes are necessary concepts, I also provided decisive arguments to the effect that the many-worlds interpretation, which is often considered to be a multiverse theory, is not viable as a realist interpretation of quantum theory, from both a logical and an observational viewpoint. But I have also explained why the quantum measurement problem is not to be considered a mere idealistic issue in a world where the emergence of quasiclassicality would be a subjective notion, associated with the biased nature of the perception of reality that would be characteristic of our conscious experience. Such a subjective approach, however, could only be made legitimate on the basis of the validity of the weak anthropic principle, whose relevance is therefore diminished by the developments I have introduced in the last portion of this report. It is somewhat ironical, therefore, that the weak anthropic principle was once considered to be bad science on the basis of the fact that it would require the existence of multiple universes, whose existence could not be confirmed by any *other* means, because, as I previously mentioned, this stubbornness is actually a form of solipsism which, in the above described context, would be supported by the weak anthropic principle, which would therefore require the existence of a multiplicity of universes.

Concerning the realist conception of quantum reality developed in this report, it is perhaps appropriate to note that, while it can be expected that the most virulent objections to such an interpretation would probably have to do with the ‘hypothesis’ of a unique reality behind interfering quantum-mechanical histories, I think that this resistance is not merely an undesirable by-product of the long tradition of instrumentalism that emerged from the Copenhagen interpretation of quantum theory, but also constitutes an unfortunate consequence of the more profound inadequacy of a philosophical position that originates from Descartes’ desire to free himself from the ‘superfluous’ hypothesis that his mind may not be all that there is in the world. I must emphasize, once again, that it is my strong belief that the most sig-

nificant challenge currently facing fundamental theoretical physics and the development of a consistent philosophy of the natural world is that of overcoming the psychological barrier associated with the reluctance to accept as real what one cannot perceive directly and to realize the sterility and the inadequacy of the opposite viewpoint, when what one wants to assess is the nature of reality itself. Here it may not be consistency alone which is at stake, but the very meaningfulness of the whole exercise, that which embodies the quest for the ultimate representation of reality.

6.1 Historical perspective

The significance of the developments introduced in this report can be better appreciated by describing the progress achieved from a historical perspective. If we start with general relativity, I think that what the theory allowed us to understand is that all motion, including acceleration, is relative and that the state of motion of an object must be defined in relation to the state of motion of the rest of the matter in the universe. Thus, relativity theory embodied in its structure the requirement of a relational definition of physical properties. But it also failed to integrate the requirement of the relativity of the sign of energy. The common belief which existed, since the creation of the general theory of relativity, is that energy must be considered positive definite, because, otherwise, apparently insurmountable problems would arise. Now, what quantum field theory allowed us to understand is that negative energies are unavoidable for properly estimating the probability of all possible transitions involving particles and antiparticles. But the current interpretation of this theory also failed to accommodate the fact that no constraint exists that would justify assuming that those negative energies are only useful for computational purposes and do not show up as properties of real matter particles, distinct from ordinary particles and antiparticles, when gravitation comes into play. What I have tried to achieve in the first portion of this report is to generalize relativity theory to produce a fully relativistic theory, compatible with the requirement that the sign of energy should also be a relative property. What motivated those developments was a better understanding of the relationship between the sign of energy of a particle and its direction of propagation in time, which again arose from applying the requirement of relational description of physical properties. In such a context, it appeared, in effect, that a concept of negative energy distinct

from that which is usually assumed to be relevant to quantum field theory was not only allowed, but was required by a truly consistent classical theory of gravitation.

Once it had been shown that the difficulties usually associated with negative-energy matter can be solved without rejecting the physical relevance of the whole concept of negative energy, there appeared to no longer be any rational motive for rejecting the possibility that negative energies can propagate forward in time and give rise to gravitational phenomena distinct from those involving exclusively positive-energy matter. It thus became clearly inappropriate to attribute a preferred status to positive-energy matter and this, in turn, meant that we are no longer justified to assume (as some authors did) that, even as it becomes integrated with general relativity, the concept of gravitationally-repulsive matter cannot involve negative energy, but must merely give rise to the notion of an observer-dependent metric devoid of any theoretical justification. Indeed, it has been clearly emphasized in this report that it is only when the concept of negative energy is well integrated to classical gravitation theory, by considering the equivalence between the presence of negative-energy matter and an absence of positive energy from the vacuum that a consistent theory (for which all measures of energy are relative) emerges which agrees with all experimental and observational constraints. The original approach which was developed in the preceding chapters is thus unique in that it actually allows to account for the very existence of the phenomenon of inertia, despite the fact that both positive- and negative-mass matter must be present on the largest scale. It alone also enables the success of the standard model of cosmology at predicting the rate of expansion of positive-energy matter to be reproduced in a bi-metric theory.

It must be clear that the concept of negative energy already existed before the developments I proposed in order to make it a consistent notion were introduced. But negative energy was always defined in an absolute or non-relational manner which, as I have shown, leads to serious difficulties, in particular because it would give rise to violations of the principle of inertia. Indeed, the idea that energy could be negative in an absolutely defined way and should therefore gravitationally repel all matter, regardless of its energy sign, as if this repulsion was a distinctive property of negative-energy matter itself, was here shown to give rise to undesirable effects, even aside from the plain logical inconsistency it would involve. The alternative interpretation of negative energy states which I have proposed has allowed to avoid those problems, while also making unnecessary the hypothesis that only positive-

energy matter can exist in stable form, because it explains why matter in a negative energy state is unobservable from the viewpoint of observers made of positive-energy matter and why even negative vacuum-dark-matter energy is mostly absent, at the present epoch, in regions of the universe occupied by positive-energy stars and galaxies. It has also become possible to explain why it is that weak gravitational lensing experiments have not revealed the presence of gravitationally repelling negative-energy matter overdensities, because it is now possible to assume that the average density of baryonic negative-energy matter is currently much smaller than that of baryonic positive-energy matter, while the absence of such matter implies that negative vacuum-dark-matter overdensities must be virtually absent on the scale of individual stars and galaxies and can only exist in smoothly distributed form in the largest voids in the distribution of positive-energy galaxies. Thus, it was actually explained why negative-energy matter has remained mostly out of reach of astronomical observations, so that this property no longer constitutes a profound mystery.

Concerning quantum reality now, the problem that there was traditionally is that we regarded its distinctive non-local character as a mere curiosity and we were convinced that it did not constitute a challenge to our conventional understanding of causality, simply because we could not see how the difficulty could be resolved if it is, in effect, real. The fact that the mathematical framework of quantum theory nevertheless allowed to produce accurate predictions, while the kind of non-locality involved did not allow information to be transmitted instantaneously, appeared to legitimize this position and this is what explains that people stopped searching for a solution to the problem of the apparent incompatibility between quantum entanglement and the constraint imposed by relativity theory on the propagation of effects. But all along, we continued searching, with more and more sophisticated experiments, for possible loopholes that could explain quantum non-locality as being an outcome of conventional unidirectional causality, just like people kept searching for experimental evidence of our motion relative to absolute space over a century ago. This happened because we were not willing to accept the conclusion that our concept of causality is inadequate in some ways, given that it does not allow to explain facts without requiring the propagation of effects at faster-than-light velocities.

It is the fact that the notion of time directionality remained so poorly understood, even after the progress which was achieved in this area by the creation of quantum field theory, that explains that there was so much confu-

sion over what constitutes an appropriate definition of the discrete symmetry operations (from the viewpoint of both clarity and consistency) when I began studying the subject. It is indeed the stubbornness to consider time from a traditional viewpoint, where only one direction is allowed for this degree of freedom, that explains that the time-reversal operation was never appropriately described and that the time-symmetric nature of causality, which allows one to make sense of quantum non-locality, was never properly assimilated. This was allowed to occur despite the clues arising from the discovery of antimatter and the successful description of antiparticles as particles propagating backward in time. The commonsense feeling inherited from our experience of thermodynamic time is so strong that it is still commonly believed that antiparticles are merely identical particles which happen to have opposite charges, rather than being the same particle propagating backward in time, as seems to be required from a mathematical viewpoint. This is what explains that time reversal was never considered to involve a reversal of charge, as I have shown to actually be required. But once this was recognized, the possibility opened up to explain other facts. It is, in effect, by using this insight that I was able to propose an explanation for the fact that a finite number of discrete degrees of freedom, which is proportional merely to the area of a black hole, allows to completely specify the microscopic state of the elementary particles which were captured by the gravitational field of such an object. In such a context, it can no longer be argued that the notion of backward-in-time propagation is merely an expedient for facilitating the calculations of probability amplitudes. Our notion of time direction has been irretrievably altered and there is no going back.

Acknowledgements

I would like to thank the government of Québec and its taxpayers which, through the generous social programs they offer, have allowed me to benefit from a steady source of income during the years in which I was working on the present project without any support from academia or the industry. They have not only allowed me to benefit from the conditions necessary to achieve the depth of knowledge required to conduct this research, but they have also saved my life at times when studying physics was the only activity that still gave enough meaning to my existence that it actually felt endurable.

Bibliography

- [1] P. A. M. Dirac. A theory of electrons and protons. *Proc. Roy. Soc. (London)*, A126:360, 1930.
- [2] R. P. Feynman. The theory of positrons. *Physical Review*, 76:749, 1949.
- [3] M. J. Pfenning. *Quantum Inequality Restrictions on Negative Energy Densities in Curved Spacetimes*. PhD thesis, Tufts University, Medford, Massachusetts, 1998. URL arxiv.org/abs/gr-qc/9805037.
- [4] H. Epstein, V. Glaser, and A. Jaffe. Nonpositivity of the energy density in quantized field theories. *Nuovo Cimento*, 36:1016, 1965.
- [5] H. B. G. Casimir. On the attraction between two perfectly conducting plates. *Proc. K. Ned. Akad. Wet. Ser. B*, 51:793, 1948.
- [6] L. H. Ford and T. A. Roman. Averaged energy conditions and quantum inequalities. *Phys. Rev. D*, 51:4277, 1995. URL arxiv.org/abs/gr-qc/9410043.
- [7] L. H. Ford and T. A. Roman. Restrictions on negative energy density in flat spacetime. *Phys. Rev. D*, 55:2082, 1997. URL arXiv.org/abs/gr-qc/9607003.
- [8] M. J. Pfenning and L. H. Ford. Quantum inequalities on the energy density in static Robertson-Walker spacetimes. *Phys. Rev. D*, 55:4813, 1997. URL arxiv.org/abs/gr-qc/9608005.
- [9] M. J. Pfenning and L. H. Ford. Scalar field quantum inequalities in static spacetimes. *Phys. Rev. D*, 57:3489, 1998. URL arxiv.org/abs/gr-qc/9710055.

- [10] D. R. Brill and S. Deser. Variational methods and positive energy in general relativity. *Ann. Phys. (New York)*, 50:548, 1968.
- [11] D. R. Brill and S. Deser. Positive definiteness of gravitational field energy. *Phys. Rev. Letters*, 20:75, 1968.
- [12] S. Deser. Timelike character of gravitational field energy-momentum. *Nuovo Cimento*, 55B:593, 1968.
- [13] D. Brill, S. Deser, and L. Faddeev. Sign of gravitational energy. *Phys. Lett.*, 26A:538, 1968.
- [14] R. M. Schoen and S.-T. Yau. On the proof of the positive mass conjecture in general relativity. *Comm. Math. Phys. (Germany)*, 65:45, 1979.
- [15] R. M. Schoen and S.-T. Yau. Proof of the positive-action conjecture in quantum relativity. *Phys. Rev. Lett.*, 42:547, 1979.
- [16] R. M. Schoen and S.-T. Yau. Positivity of the total mass of a general space-time. *Phys. Rev. Lett.*, 43:1457, 1979.
- [17] R. M. Schoen and S.-T. Yau. The energy and the linear momentum of space-times in general relativity. *Comm. Math. Phys.*, 79:47, 1981.
- [18] R. M. Schoen and S.-T. Yau. Proof of the positive mass theorem. ii. *Comm. Math. Phys.*, 79:231, 1981.
- [19] M. M. Nieto and T. Goldman. The arguments against ‘antigravity’ and the gravitational acceleration of antimatter. *Physics Reports*, 205: 5, 1991.
- [20] H. Bondi. Negative mass in general relativity. *Reviews of Modern Physics*, 29:423, 1957.
- [21] H. Reichenbach. *The Philosophy of Space & Time*. Dover Publications, 1957.
- [22] P. J. E. Peebles. *Principles of Physical Cosmology*, pages 108, 296. Princeton University Press, 1993.
- [23] T. Piran. On gravitational repulsion. *Gen. Rel. and Grav.*, 29:1363, 1997. URL arxiv.org/abs/gr-qc/9706049.

- [24] G. D. Birkhoff. *Relativity and Modern Physics*, page 253. Harvard University Press, 1923.
- [25] S. Hossenfelder. A bi-metric theory with exchange symmetry. *Phys. Rev. D*, 78:044015, 2008. URL arxiv.org/abs/0807.2838.
- [26] J. C. Lindner. Theoretical basis for a solution to the cosmic coincidence problem, 2006. URL arxiv.org/abs/gr-qc/0603005.
- [27] A. Perez. Statistical and entanglement entropy for black holes in quantum geometry. *Phys. Rev. D*, 90:084015, 2014. URL arxiv.org/abs/1405.7287.
- [28] S. Weinberg. The cosmological constant problem. *Rev. Mod. Phys.*, 61:1, 1989.
- [29] I. Labbé, P. van Dokkum, E. Nelson, R. Bezanson, K. A. Suess, J. Leja, G. Brammer, K. Whitaker, E. Mathews, M. Stefanon, and B. Wang. A population of red candidate massive galaxies ~ 600 myr after the Big Bang. *Nature*, 2023. URL arxiv.org/abs/2207.12446.
- [30] M. Zumalacarrégui and U. Seljak. Limits on stellar-mass compact objects as dark matter from gravitational lensing of type Ia supernovae. *Phys. Rev. Lett.*, 121:141101, 2018. URL arxiv.org/abs/1712.02240.
- [31] S. S. McGaugh, F. Lelli, and J. M. Schombert. Radial acceleration relation in rotationally supported galaxies. *Phys. Rev. Lett.*, 117:201101, 2016. URL arxiv.org/abs/1609.05917.
- [32] M. Meneghetti, G. Davoli, P. Bergamini, P. Rosati, P. Natarajan, C. Giocoli, G. B. Caminha, R. B. Metcalf, E. Rasia, S. Borgani, F. Calura, C. Grillo, A. Mercurio, and E. Vanzella. An excess of small-scale gravitational lenses observed in galaxy clusters. *Science*, 369:1347, 2020. URL arxiv.org/abs/2009.04471.
- [33] J. Moreno, S. Danieli, J. S. Bullock, R. Feldman, P. F. Hopkins, O. Çatmabacak, A. Gurvich, A. Lazar, C. Klein, C. B. Hummels, Z. Hafen, F. J. Mercado, S. Yu, F. Jiang, C. Wheeler, A. Wetzel, D. Anglés-Alcázar, M. Boylan-Kolchin, E. Quataert, C.-A. Faucher-Giguère, and D. Kereš. Galaxies lacking dark matter produced by close

- encounters in a cosmological simulation. *Nature Astronomy*, 6:496, 2022. URL arxiv.org/abs/2202.05836.
- [34] E. Tryon. Is the universe a vacuum fluctuation? *Nature*, 246:396, 1973.
- [35] R. Penrose. Singularities and time asymmetry. In *General relativity: An Einstein centenary survey*. Cambridge University Press, 1979.
- [36] S. M. Carroll and J. Chen. Spontaneous inflation and the origin of the arrow of time, 2004. URL arxiv.org/abs/hep-th/0410270.
- [37] T. Tröster, A. G. Sánchez, M. Asgari, C. Blake, M. Crocce, C. Heymans, H. Hildebrandt, B. Joachimi, S. Joudaki, A. Kannawadi, C.-A. Lin, and A. Wright. Cosmology from large-scale structure: Constraining Λ CDM with BOSS. *Astronomy & Astrophysics*, 633(L10), 2020. URL arXiv.org/abs/1909.11006.
- [38] M. Asgari, C.-A. Lin, B. Joachimi, B. Giblin, C. Heymans, H. Hildebrandt, A. Kannawadi, B. Stölzner, T. Tröster, J. L. van den Busch, A. Wright, M. Bilicki, C. Blake, J. de Jong, A. Dvornik, T. Erben, F. Getman, H. Hoekstra, F. Köhlinger, K. Kuijken, L. Miller, M. Radovich, P. Schneider, H. Y. Shan, and E. Valentijn. Kids cosmology: Cosmic shear constraints and comparison between two point statistics. *Astronomy & Astrophysics*, 645(A104), 2020. URL arXiv.org/abs/2007.15633.
- [39] L. Sklar. *Physics and Chance*, page 395. Cambridge University Press, 1993.
- [40] D. Deutsch. Quantum mechanics near closed timelike lines. *Phys. Rev. D*, 44:3197, 1991.
- [41] R. P. Feynman and J. A. Wheeler. Interaction with the absorber as the mechanism of radiation. *Rev. of Mod. Phys.*, 17:157, 1945.
- [42] P. L. Csonka. Advanced effects in particle physics. *Phys. Rev.*, 180:1266, 1969.
- [43] R. B. Partridge. Absorber theory of radiation and the future of the universe. *Nature*, 244:263, 1973.

- [44] J. G. Cramer. The arrow of electromagnetic time and the generalized absorber theory. *Foundations of Physics*, 13:887, 1983.
- [45] O. Costa de Beauregard. No paradox in the theory of time anisotropy. In *The study of time I*, page 132. Springer Verlag, 1972.
- [46] J. G. Cramer. The transactional interpretation of quantum mechanics. *Rev. of Mod. Phys.*, 58:647, 1986.
- [47] Y. Aharonov, P. G. Bergmann, and J. L. Lebowitz. Time symmetry in the quantum process of measurement. *Phys. Rev. B*, 134:1410, 1964.
- [48] Y. Aharonov and L. Vaidman. The two-state vector formalism: an updated review. In *Lecture notes in physics*, volume 734, page 399. Springer, 2007. URL arxiv.org/abs/quant-ph/0105101.
- [49] M. Proietti, A. Pickston, F. Graffiti, P. Barrow, D. Kundys, C. Branciard, M. Ringbauer, and A. Fedrizzi. Experimental rejection of local observer independence. *Science Advances*, 5(9), 2019. URL arxiv.org/abs/1902.05080.
- [50] R. P. Feynman. Negative probability. In *Quantum implications: Essays in honor of David Bohm*, page 235. Routledge, 1987.
- [51] A. Zeilinger, R. Lapkiewicz, P. Li, C. Schaeff, N. K. Langford, S. Ramelow, and M. Wieśniak. Experimental non-classicality of an indivisible quantum system. *Nature*, 474:490, 2011. URL arxiv.org/abs/1106.4481.
- [52] Y. Aharonov, F. Colombo, S. Popescu, I. Sabadini, D. C. Struppa, and J. Tollaksen. The quantum pigeonhole principle and the nature of quantum correlations. 2014. URL arxiv.org/abs/1407.3194.
- [53] E. Joos, H. D. Zeh, C. Kiefer, D. Giulini, J. Kupsch, and I.-O. Stamatescu. *Decoherence and the Appearance of a Classical World in Quantum Theory*, page 21. Springer-Verlag, second edition, 2003.
- [54] H. D. Zeh. *The Physical Basis of the Direction of Time*, page 81. Springer-Verlag, second edition, 1992.

- [55] F. Dowker and A. Kent. On the consistent histories approach to quantum mechanics. *Journal of Statistical Physics*, 82:1575, 1996. URL arxiv.org/abs/gr-qc/9412067.
- [56] R. Omnès. *Understanding Quantum Mechanics*. Princeton University Press, 1999.
- [57] R. B. Griffiths. Consistent histories and the interpretation of quantum mechanics. *Journal of Statistical Physics*, 36:219, 1984.
- [58] R. Omnès. Logical reformulation of quantum mechanics. I. Foundations. *Journal of Statistical Physics*, 53:893, 1988.
- [59] M. Gell-Mann and J. B. Hartle. Quantum mechanics in the light of quantum cosmology. In W. H. Zurek, editor, *Complexity, Entropy and the Physics of Information*, volume VIII of *SFI Studies in the Science of Complexity*. Addison Wesley, 1990.
- [60] S. Goldstein and D. Page. Linearly positive histories: Probabilities for a robust family of sequences of quantum events. *Phys. Rev. Lett.*, 74:3715, 1995. URL arxiv.org/abs/gr-qc/9403055.
- [61] R. Arnowitt, S. Deser, and C. W. Misner. The dynamics of general relativity. In L. Witten, editor, *Gravitation: An Introduction to Current Research*. Wiley, 1962.
- [62] J. A. Wheeler. Geometrodynamics and the issue of the final state. In C. DeWitt and B. DeWitt, editors, *Relativity, Groups, and Topology: 1963 Les Houches Lectures*. Gordon & Breach, 1964.
- [63] J. A. Wheeler. Superspace. In R. P. Gilbert and R. Newton, editors, *Analytic Methods in Mathematical Physics*. Gordon & Breach, 1970.
- [64] A. E. Fischer. The theory of superspace. In M. Carmeli, S. I. Fickler, and L. Witten, editors, *Relativity*. Plenum, 1970.
- [65] W. K. Wootters, A. Aleksandrova, and V. Borish. Real-vector-space quantum theory with a universal quantum bit. *Phys. Rev. A*, 87:052106, 2013. URL arxiv.org/abs/1210.4535.

Index

- absolute space, 521, 644
- action at a distance, 493, 516, 526, 538, 546–548, 599
- advanced waves, *see* advanced solutions, wave equations
- advanced-waves problem, 479–484, 535n
- Aharonov, Yakir, 487, 537
- alternative counterfactual, 465
- anthropic principle, *see* weak anthropic principle
- antimatter
 - baryon-antibaryon annihilation, 290n, 294–298, 324, 365–367, 409, 449
 - baryon-antibaryon asymmetry, 325
 - early matter-antimatter annihilation, 285, 292–298, 303, 324, 332, 342, 343–347, 357, 367–372, 408, 417, 438, 444, 636, 637
 - matter-antimatter asymmetry, 290, 320–325, 438, 637
 - pair creation and annihilation, 124–128, 285, 290, 306–308, 322–325, 411, 444, 594
- arrow of time, *see* thermodynamic arrow of time, time irreversibility
- artificial super-intelligence, 491
- ATP interferometer experiment, 537
- axioms
 - negative energy, 151–157, 177–181
 - negative mass, 151–157, 177
 - negative-energy matter, 151–158, 177–181
- backward-in-time signaling, 516–521, 543–547
- Barbour, Julian, 579–581
- bare cosmological constant, 282
- baryonic matter, 290, 291, 297, 302, 303, 329, 334, 335, 342–345, 388, 398, 438, 451, 637
 - inhomogeneities, 327
- Bekenstein bound, 230, 255–266, 386
- Bell’s inequality, 515–517, 538–541
- Bell, John, 566
- Bergmann, Peter, 487
- bidirectional charge, 306–316, 321–325
 - non-reversed-bidirectional-charge particle, 311–316, 325
 - reversed-bidirectional-charge particle, 306–325, 337, 338, 637
- Big Bang
 - boundary conditions, 275, 409, 452, 603, 608n
 - energy conservation, 353–355, 361
 - hot Big Bang, 371, 412, 440
 - initial conditions, 290, 332, 350–355, 372, 373, 409, 410, 416–440, 448–454, 462–466, 473–

- 475, 489, 490, 519, 556–558, 572–574, 580–589, 595, 621, 635
- initial expansion rates, 295–298
- initial singularity, *see* past singularity, Big Bang
- light element abundances, 434
- matter-dominated era, 290–292, 298, 342–344, 367, 410
- maximum energy densities, 127, 285, 351, 352, 357–362, 372, 373, 413, 414, 421–440, 448, 449, 479, 548, 595, 603, 637
- maximum matter density, *see* maximum energy densities, Big Bang
- nucleosynthesis, 319
- past singularity, 127, 135, 351, 359, 368–371, 408, 413–440, 595, 603, 610
- pieces of evidence, 434
- quantum-bounce continuity condition, 371, 436, 437
- quark-hadron transition, 347
- radiation reflector, 480
- radiation-dominated era, 290–295, 343
- singularity theorems, 413
- time before the Big Bang, 321, 371, 428, 436–438, 490, 595, 610
- trapped surface, 413
- uniqueness, 435
- variation of expansion rate, 351, 352, 372
- Big Bounce, *see* quantum bounce, quantum gravitation
- black body radiation problem, 24, 277
- black hole, 146–148, 229, 387, 388
- angular momentum, 230, 239–244
- charge, 230, 239–241, 254, 255
- conservation of information, 230, 231, 259
- decay, *see* evaporation, black hole
- discrete radiation degrees of freedom, 249n
- electric charge, 252, 339
- electromagnetic field, 253–255
- entropy, 9–12, 148, 182, 229–261, 266–270, 277, 336–340, 375, 385–389, 393n, 394, 395, 415, 430, 634
- entropy density, 409, 416
- evaporation, 231, 248–250, 259, 268, 415
- event horizon, 22, 148, 229–233, 239, 244–249, 258–268, 337, 369, 375, 385–395, 408–415, 423, 428–442, 448–451, 573, 574, 638
- event-horizon degrees of freedom, 230–240, 249–261, 336–340, 386, 400
- future singularity, *see* spacetime singularity, black hole
- generic past singularity, 436
- gravitational collapse, 147, 239–257, 337, 386, 393, 414
- gravitational potential, 259
- Hawking radiation, 231, 248–250, 255–261, 270, 271, 386, 496
- information, 12, 229–244, 249–266, 336–340, 375–377
- information availability limitation, 260, 261
- information-loss paradox, 250
- macroscopic event horizon, *see* event

- horizon, black hole
- macroscopic parameters, 230, 231, 241, 242, 259, 260, 387–391, 638
- mass, 230–241, 251, 260, 271, 386, 387, 393
- mass reduction, 148–151, 248
- matter absorption, 148–151
- matter degrees of freedom, 182, 229–244, 250–268, 336–340, 389, 621, 634, 644
- merger, 268
- momentum, 239–241
- negative-energy matter, 239–247, 255–272
- negative-energy radiation, 270
- non-rotating, 244
- particle bidirectional charge sign, 339, 340
- particle charges, 252–257, 336–338
- particle energies, 250, 251, 257–266, 336–338
- particle entanglement, 250
- particle handedness, 252–257, 336–340
- particle momenta, 244–258, 263–266, 336–340
- particle space intervals, 336–338, 629
- particle time intervals, 338, 339
- potential barrier, 428
- primordial black hole, 408–412, 443
- redshift, 246
- Schwarzschild radius, 235
- semi-classical theory, *see* thermodynamics, black hole
- space contraction, 249, 259
- spacetime singularity, 135, 146, 147, 229, 240–250, 255–259, 266, 336–340, 386, 392–396, 414, 428–431, 450
- stable state, 240, 251, 256–268, 336–340
- surface area, 148, 232–235, 245, 252–254, 260, 386, 394–396, 436, 638
- surface gravitational field, 234–238, 252–257, 264, 270, 336–340, 394–400
- temperature, 251–255, 270, 271
- thermal radiation, *see* Hawking radiation, black hole
- thermodynamics, 9, 229–240, 249–257, 267–272, 336–340, 386, 392, 634
- time dilation, 246–250, 259, 260
- Bohr, Neils, 499, 505, 548
- Boltzmann’s constant, 235
- Boltzmann, Ludwig, 381, 550
- Boolean logic, 494, 501, 507, 525, 530, 531, 567
- broken symmetry, 297
- causal chain
 - closed, 459, 468–475, 528–531, 576, 586–590, 609, 619–622, 640
 - invariant direction in time, 466
 - particle-antiparticle loop, 469
 - time-reversed, 463–466
 - traditional concept, 584
 - unidirectional, 463–466
 - universal, *see* universal causal chain
- causality, 100, 423
 - absolute time order, 469, 520, 521
 - backward causation, 418, 419, 424, 460–477, 482, 488–490, 516–

- 521, 528–547, 555, 567, 587,
609, 619, 620
- bidirectional, *see* time-symmetric,
causality
- causal circularity, *see* closed, causal
chain
- causal discontinuity, 587, 593
- causal order, 613
- causal ordering postulate, 466
- causal process continuity, 586n, 593,
594, 640
- causal relationships, 460, 477, 573,
619
- causal signals, 588
- causal structure of spacetime, 12,
462, 482, 515, 589, 603–606
- causal theory, 620, 621
- causally independent branches, 514
- causes and effects, 424, 454, 460–
471, 463n, 481, 519–521, 557,
587, 602, 603, 621
- classical, *see* unidirectional, causal-
ity
- classical spacetime, 425
- definite time order, 467
- direct contact, 425–431, 440, 460
- external cause, 586, 587, 618
- faster-than-light communication,
497, 516–519, 539n, 547, 644
- faster-than-light propagation, 501,
516, 526, 546
- final causes, 467–470, 519–522
- final conditions, 465, 483
- first cause, 587
- global consistency requirement, 465–
477, 490, 514, 528–531, 554,
587, 609
- local, 417, 424, 435, 462, 488, 538,
546, 573–575, 609
- local contact, 603
- local perturbations, 442
- non-local, 507
- ordinary causation, 418, 419, 424
- past and future light cones, 462,
482, 517, 573, 574, 589
- principle of local causality, 10, 55,
274, 307n, 455–459, 483, 497–
502, 509–518, 535, 546, 572,
595–599, 615–623, 639
- psychological expectation, 467, 487
- relative time order, 469, 520, 521
- relativistic speed limit, 466, 482,
514–518, 538, 644
- relevance to cosmology, 589
- signals, 423–431
- spreading of effects, 603–606, 611–
614
- teleological character, 464
- time-symmetric, 10, 189, 425, 454,
460–469, 463n, 477, 508–518,
528, 536, 547–549, 574, 586,
587, 617–620, 631, 639, 644
- time-symmetric causality violation,
587
- two essential conditions, 589
- unidirectional, 418–425, 432–436,
454, 460–475, 481, 482, 487,
519–521, 542, 574, 603–606,
640–644
- unique direction in spacetime, 574–
582
- wave function, 620
- chaotic systems, 383
- strong non-linearity, 384
- classical curvature of space, 612
- classical gravitation theory, 10–12, 20,

- 98, 623
- classical gravitational field, 623
 - microscopic properties, 621
- classical history, 485–490, 499, 640
- classical instability, 379, 383
- classical measuring device, 562
- classical mechanics, 374–377
- classical neutron interferometer experiment, 628
- classical particle, 485, 492, 505
- classical path, 554
- classical physics, 492, 493, 625
- classical probability, 512, 523–533, 640
- classical reality, 506–512, 526, 536, 537, 563, 589
- classical scale, 626
- classical space and time, 580–582
- classical spacetime, 250, 573–582, 592, 611–614, 624
- classical state, 515, 525, 544, 561
- classical trajectory, 485, 486, 495, 537–547, 545n, 592
- classical wave, 485–492, 506, 523–526, 535
- classically meaningful probability, 529, 530, 550, 561–565
- classically unique reality, 554, 555
- closed-circuit analogy, 458, 459
 - pairs of polarized wires, 459, 588
- compact topological structure, 599
- complex number, 534
- complex-number weighting coefficients, 509
- constraint of relational definition
 - absence of external cause, 307n, 640
 - absolute direction, 184–186, 207, 570n
- absolute lopsidedness, 184–188, 320, 321, 637
- absolute space, 75–79
- center of mass of the universe, 95, 96
- completeness, 185
- direction of propagation in time, 458, 463–467, 489, 522
- directional asymmetry, 184–187, 217
- discrete symmetry operation, 184–188
- energy of the universe, 274, 354, 366–368, 635
- fundamental lopsidedness, *see* absolute lopsidedness, constraint of relational definition
- gravitational force, 43, 56–62, 114–118, 158, 177, 413
- imbalance, 186
- interaction field energy, 117
- lower energies, 130–133
- metaphysical elements of reality, 570n
- momentum of the universe, 354
- physical attributes, 20, 21, 325, 354, 472, 570, 618, 635
- physical attributes of the universe, 354, 586
- polar asymmetry, 184, 185
- principle of relativity, 75–78, 87, 88, 186
- relativity of acceleration, 72–88, 642
- reversal of energy, 184
- reversal of momentum, 184, 199, 226
- self-determination, 185

- sign of charge, 184, 185, 199, 208–210, 316, 458
- sign of energy, 8, 26, 102–106, 114, 151, 157, 158, 177, 199, 226–228, 282, 299, 642, 643
- sign of mass, 55–69, 77, 88, 151, 157, 158
- space and time directions, 187, 199, 204–207, 226, 324, 354, 637
- space and time reversals, 184–187, 199, 208
- time, 572, 584
- time-direction-dependent property, 51
- universe, 20, 21, 55, 56, 81, 131, 184–187, 228
- conventional rules of logic, *see* Boolean logic
- converging wave front, 479–481
- coordinative definition, 41, 58, 308
- Coriolis force, 84
- correlated phases, 479
- correlation probability, 460–465, 470, 477, 529, 576, 585, 590
- cosmic horizon, 362, 417, 422–427, 433–435, 440–443, 603
- cosmic microwave background, 282, 347, 408, 409, 417, 434, 442, 447–451
 - angular scale of fluctuations, 444
 - epoch of decoupling, 443, 444
 - epoch of last scattering, 347, 412, 440–445
 - epoch of recombination, 344
 - gravitational wave signal, 446
 - matter density estimates, 444
 - negative-energy matter, 444
 - polarization, 446
- spectrum harmonics, 443
- temperature fluctuations, 347, 417, 418, 442–446
- temperature homogeneity, 416, 440–442, 448
- temperature modifications, 416
- cosmological principle, 79, 95, 440
- cosmology, 9–14, 19–23, 454
 - S_8 tension, 445
 - biasing hypothesis, 344
 - cold-dark-matter model, 348, 349
 - cosmic coincidence problem, 276
 - cosmological-constant problem, 18, 274–281, 299, 374, 446, 634
 - cosmology problem, 18
 - current difficulties, 409
 - dark-matter problem, 300n, 309
 - fine-tuning problems, 374
 - flatness problem, 18, 274–276, 350–374, 446–450, 635
 - horizon problem, 18, 274–276, 350, 369, 440–451
 - inflation problem, 9, 274–276, 350, 374, 441, 446, 447, 452, 635–638
 - inflationary cosmology, 443
 - initial value equation, 356, 357
 - local variations of the flow of time, 583n
 - matter creation problem, 371, 446
 - negative-energy-matter cosmology, 452
 - outstanding problems, 273–276, 446, 447, 452, 633, 634
 - problem of dark energies, 273–276
 - problem of time asymmetry, 374, 375, 380–385, 447, 550, 578n
 - smoothness problem, 274, 416, 442–

- 446
- standard Big Bang model, 18, 276, 350–352
- standard model, *see* standard Big Bang model, cosmology
- time as a measure of relative changes, 572, 584
- topological defects problem, 274, 441, 446
- uniform flow-of-time direction, 580
- uniform time flow, 12, 574, 575, 583
- universal time variable, 570–577, 582–584, 621, 639
- weak gravitational lensing, 303, 325, 343, 348, 444
- Costa de Beauregard, Olivier, 480
- Cramer, John, 484–487
- critical density, 282, 304, 330–332, 350, 351, 368–372, 451
- critical expansion rate, 350–353, 362–373, 420, 447–449, 635
- d’Espagnat, Bernard, 566
- dark energy, 9, 18, 273–277, 282n, 297, 327, 636
 - quintessence, 277
- dark matter, 9, 18, 25, 47, 273, 281, 282, 300–309, 316–318, 327, 341, 349, 637
 - absence of, 336n
 - baryonic, 300, 309, 319, 335–338, 637
 - clumping, 336
 - cold, 303, 342–344, 349
 - colliding clusters, 319, 336
 - computer simulations, 336n
 - concentration around cluster elements, 335
 - constant positive and negative energies, 332–334, 372
 - correlation between gravitational accelerations, 335, 336
 - dwarf galaxies, 336n
 - galaxy mergers, 336
 - gravitational lensing, 335
 - gravitational repulsion, 327
 - homogeneous distribution, 302, 332–334
 - hybrid form of matter, 332
 - inappropriateness of conventional interpretation, 335
 - local variations of vacuum energy, 398–402
 - MACHOs, 319
 - nature of dark matter, 273, 319, 335, 336
 - negative energy, 302, 303, 318, 325, 330–334, 342–348, 372, 407, 417, 444, 445, 636
 - non-baryonic, 325
 - overdensities, 303
 - reversed-bidirectional-charge particles, 191n, 215, 309, 310, 315–325, 338, 637
 - self-interacting, 319
 - spherical halos, 319
 - vacuum dark matter, 329–336, 341–348, 372, 388, 396–403, 417, 444–446, 496, 636–638, 643
 - vacuum dark matter overdensity, 345
 - vacuum energy, 273, 300, 329–336
 - variation of density, 332
 - weakly-interacting massive particles, 300n, 329, 637

- weakly-interacting particles, 305
- delayed choice experiment, 488, 511
- density parameter, 351, 447
- Descartes, René, 641
- determinism, 563, 619
- deterministic evolution, 501, 559, 560, 576, 619–622
- deterministic theory, 470, 622
- Deutsch, David, 471, 472
- Dirac, Paul, 28, 48, *see* Dirac's solution, negative energy
- discrete symmetry
 - violation, 182–188, 206–209, 216, 226–228, 320–325, 343, 438, 637
 - weak interaction, 184, 213
- discrete symmetry operations
 - action-sign degree of freedom, 221–227, 238–243
 - alternative formulation, 12, 182, 183, 197–209, 215, 229, 241, 268, 634
 - angular momentum, 183, 184, 194–219, 225, 241–243, 252
 - antimatter, 182, 191, 198, 206, 207, 320, 337, 338
 - backward motion, 183
 - basic action-reversal operation M_I , 220, 225–227
 - charge conjugation C , 182–186, 197–227, 238, 254, 321–325
 - classical equations, 218, 219
 - combined operations, 186, 207, 214–217, 225–227
 - conjugate attributes, 197, 212, 218
 - CPT theorem, 186, 215
 - currents, 207
 - dependencies, 188, 209
 - electric field, 207, 208
 - enantiomorphic equivalent, 210
 - equivalent operations, 215, 216
 - fermion quantum phase, 216
 - fermion wave function, *see* fermion quantum phase, discrete symmetry operations
 - field polarities, 596
 - fundamental degrees of freedom, 233–243, 252, 257, 264, 338, 386, 387, 637, 644
 - gravitation, 182, 228
 - handedness, 184, 198–204, 213, 241–243, 252–257, 336–339
 - identity operation I , 221–226
 - invariance of the sign of action, 197–214, 220–224, 242
 - joint variation, 197
 - kinematic representation, 183–189, 195, 206–210
 - magnetic field, 207, 208
 - microscopic state, 12, 277
 - momentum, 183, 184, 191–213, 218–228, 241–243, 252–254, 265, 266
 - parity P , *see* space reversal P , discrete symmetry operations
 - PTC transformation, 186, 206, 214–216, 226, 227, 321
 - quantum field theory, 182–186, 199
 - quantum operators, 199
 - reversal of action M , 9, 182, 197, 202, 220–228, 237–243, 279
 - reversal of motion, 189, 195, 203–210, 224–228
 - semi-classical viewpoint, 182
 - sign of charge, 182–192, 198–201, 206–215, 224, 241–243, 252–

- 257, 305, 336–338, 596, 644
- sign of energy, 182, 191, 197–211, 220–228, 239–243, 258, 266, 305, 338
- space and time coordinates, 191–215, 220
- space intervals, 183–213, 218–222, 228, 238, 243–246, 251, 252, 258, 265, 266
- space reversal P , 9, 182–227, 237–242, 251, 322
- space rotation, 216
- space-related properties, 197–199, 204, 212–225, 237
- spacetime reversal, 182, 211–215
- spin, *see* angular momentum, discrete symmetry operations
- time intervals, 184–192, 197–212, 218–222, 228, 238–243, 252–254
- time reversal T , 9, 12, 182–228, 237–242, 254, 290, 320–325, 336, 343, 379, 438, 454, 577, 596, 619n, 637, 644
- time-related properties, 197–199, 204, 212–225, 237
- traditional conception, 182–188, 194–206, 213–215, 379, 644
- double slit experiment, 484–488, 506–508, 522–533
- Dowker, Fay, 556, 564, 565, 606
- dynamic equilibrium of forces, 63–68
- Einstein's elevator experiment, 366n
- Einstein, Albert, 20, 75, 83, 170, 502, 508, 514, 548
- electromagnetic waves, 478–482
- electrostatic field pair creation, 255
- elementary particle, 11, 384, 425–434, 448, 638
 - composite particles, 315
 - conjugate attributes, 31n, 513
 - flavor changing interaction, 316
 - gravitational interaction, 316–318
 - indistinguishable, 502
 - interactions, 311–316, 497
 - localized nature, 495, 505, 506, 545, 603
 - physical attributes, 31
 - quantum reality, 492–498, 504–508, 512, 521
 - reality through interaction, 568
 - relevance of time, 578
 - static attribute, 243n, 502, 628
 - theories, 315
 - unified charge, 243n, 400n
 - wavelength, 495
- energy-out-of-nothing problem, 54, 112, 113, 123–128, 133, 143, 144, 633
 - work and useful energy, 145, 481
- epistemology, 14
- equations of state, 105
- equivalent gravitational field, 62–80, 114, 139
- Everett III, Hugh, 551
- expansion rate
 - isotropic, 367–369
 - observer dependence, 277, 293–295
- false vacuum, 297
- fermion, 215, 307–310
- Feynman diagrams, 504, 505
- Feynman, Richard, 19, 49, 116, 278, 460, 479, 480, 492, 497n, 503–505, 505n, 532, *see* Feynman's

- interpretation, time-direction
 - degree of freedom
- free-will, 461–465, 470, 476, 477, 528
- Friedmann potential, 356
- fundamental principles, 20, 409, 434
- Gell-Mann, Murray, 546, 561–564
- general relativity, 13–19, 584n
- general-relativistic theory
 - absence of gravitational field, 624
 - additional variables, 165
 - alternative notation, 165n
 - alternative proposals, 163–165
 - assumptions, 160
 - average stress-energy tensors, 168, 179
 - bi-metric theories, 8, 159–170, 177–180, 643
 - Cartan tensor, 256
 - classical approximation, 231, 627
 - closed time-like curve, 147, 581, 613
 - conjugate metrics, 160
 - consequences, 160
 - conservation of energy, 135, 178, 179
 - constant and uniform gravitational field, 249, 625–627
 - cosmological term, 175, 282, 327
 - curvature of space, 623
 - curvature tensors, 161–163, 177
 - deterministically-evolving metric properties, 626
 - differentiation of time from space, 572–575, 580–582, 588, 589, 603, 613
 - distinctive features, 159
 - dynamical theory, 570, 575, 581, 584n, 625
 - earlier interpretations, 160, 161
 - earlier publications, 159, 160
 - Einstein tensor, 165, 172, 256
 - energy sign convention, 158
 - Euclidean metric, 624
 - extrinsic curvature of space, 591, 612, 624
 - fluctuating metric properties of space and time, 627–630, 640
 - foundations, 631
 - general covariance, 162
 - generalized gravitation theory, 9–11, 158, 159, 278, 300, 336, 352, 366, 446–454, 575, 580, 615, 621, 630–634, 639–643
 - generalized gravitational field equations, 8, 151, 158–181, 228, 282, 288–299, 327–331 327n, 366, 413, 621, 627, 634–636
 - geodesics, 161
 - gravitational field energy, 135–138, 179, 247
 - gravitational field equations, 356, 570, 581
 - intrinsic curvature of space, 583–585, 591, 612, 624
 - irregular stress-energy tensors, 168–180
 - justifications, 160
 - local growth of space curvature, 334
 - local topology correspondence, 285, 286
 - local variations of light-cone structure, 573, 574, 613
 - locally Euclidean metric, 626
 - maps, 165

- mathematical requirement, 165
- mathematical structure, 13, 158–177
- metric conversion factors, 165–176, 283–288, 293–296, 326, 327, 355
- metric properties of space and time, 333, 584, 585, 611–614, 621, 627
- natural vacuum-stress-energy tensors, 173, 174, 181, 282–287, 293, 327, 351, 413
- natural viewpoints, 158
- Newtonian limit or approximation, 180
- observer-dependent energy sign, 158–167, 179
- observer-dependent gravitational field, 70–72, 90, 158–162, 177–181, 285, 329
- observer-dependent metric properties, 20, 74, 158, 159, 165–175, 181, 283–287, 295–297, 326–335, 360, 634–637
- omnipresent gravitational field, 625–627
- physical requirements, 168–172
- proper time interval, 574
- pull-overs, 165
- redefined energy ground state, 164–174, 179
- relativistic invariance, 20
- Riemannian spacetime, 494
- semi-classical theory, 631
- simultaneity hyperplanes, 581–584
- slicing of spacetime, 570–572, 581–584
- smooth gravitational field, 262
- solution of gravitational field equations, 573
- space contraction, 326, 327
- space dilation, 326, 327, 630
- space-like hypersurface, 570–576, 581–587, 597, 612, 619n
- spacetime curvature as the outcome of an interaction, 621
- spacetime foliation, 570, 571
- spacetime torsion, 255, 256, 336–340
- spin angular momentum tensor, 256
- stress-energy tensors, 62, 160–179, 252–256
- time contraction, 326, 327
- time dilation, 326
- unique metric signature, 12, 570–573, 581, 588, 613
- vacuum stress-energy tensors, 172–176, 181
- vacuum-energy terms, 167–176, 282–299, 327n, 366, 636
- variational principle, 170
- global entanglement constraint, 426–443, 450–452, 460–466, 473–475, 480–482, 489, 490, 516–519, 547, 548, 558, 572–576, 588, 589, 603, 610–614
- global inertial reference system, 84–88, 366n, 572n
- Grand Unified Theories, 36, 278–280, 441
- gravitational lensing
 - arcs of light, 348
 - blobs of light, 348
 - repulsive, 348
- gravitational physics, 23, 454, 455

- gravitational potential energy, 135–138
- gravitational repulsion, 8, 53, 113–116, 146–151, 162, 178–180, 228, 240–244, 261, 281, 299, 405, 412–416, 428–430, 450, 458, 630–633, 643
 - antimatter, 115, 143–146
 - antimatter experiment, 143–146
 - from matter overdensities, 301–303, 341–350
 - from missing positive vacuum energy, 100, 153, 178, 318, 394–399
 - from voids in a matter distribution, 91–98, 153, 178, 300n, 346–350, 390–399, 405–409
 - uncompensated, 110, 301
 - uncompensated gravitational attraction, 94–100, 107, 119–121, 135, 151–154, 178
 - weakness, 138
- Griffiths, Robert, 561–565
- Hartle, James, 561–564
- Heisenberg, Werner, 499, 506
- Hilbert space
 - linearity, 510
- Hossenfelder, Sabine, 159
- Hubble constant, 351, 352
- idealism, 499
- inconsistent probabilities, 530, 531
- inertial gravitational force, 62–73, 79–84, 114, 141, 180
 - from opposite-energy matter distributions, 79–87, 109, 180
- inflation theory, 18, 274, 350, 351, 359, 369, 418–421, 444–452, 635
 - accelerated expansion, 350, 369, 418, 435, 446–450
 - bubble universe, 452
 - eternal inflation, 452
 - exponentially accelerated contraction, 418–420
 - fine-tuning, 449
 - free parameters, 452
 - inflation process, 359, 368–374, 418–422, 432, 439–452
 - inflaton field, 449
 - megaverse, 451, 452
 - observational evidence, 447–451
 - primordial quantum fluctuations, 443
 - reheating, 371, 448–450
 - required initial conditions, 420, 421, 441, 447, 448
 - weakening gravitational repulsion, 444
- information, 376, 387
 - availability, 260, 429
 - average matter density, 405, 406
 - average particle energy, 405, 406
 - bidirectional charge sign, 338–340, 400n
 - conservation, 249, 250, 261, 377, 385–406, 393n, 474, 554, 623, 638
 - cosmic expansion, 385, 399–403, 402n, 600, 638
 - creation out of nothing, 474
 - energy inhomogeneities, 407
 - expansion of space, 406, 407
 - general surface, 230–234, 240, 253–268, 338, 386
 - global decrease, 399–407
 - global growth, 400–403, 638

- global variation, 401
- gravitational field, 377, 385–394, 399–407, 600, 638
- homogeneous matter distribution, 399–407, 638
- homogeneous matter-energy distribution, 407
- invariant measure, 377, 595, 600
- local decrease, 389–408, 638
- local growth, 386–399, 393n, 404–408, 638
- loss, 231, 249, 259, 337, 338, 376, 377, 386–392, 408
- matter inhomogeneities, 405, 406
- matter overdensity, 261, 389–399, 406–408
- matter temperature, 405
- matter underdensity, 389, 403–406
- microscopic gravitational field configuration, *see* microscopic state of gravitational field, information
- microscopic gravitational fields, 406
- microscopic matter inhomogeneities, 406
- microscopic state of electromagnetic field, 338
- microscopic state of gravitational field, 254, 255, 262, 267, 333, 385–411, 595, 638
- microscopic state of interaction field, 336–340, 393n
- microscopic state of matter, 238, 262
- microscopic state of torsion field, 256, 336–340, 400n
- microscopic structure, 250, 375–377, 388, 403–405
- minimal coarse-graining, 377
- negative information, 404
- perfectly uniform fluid, 405
- vacuum dark matter, 407
- vacuum energy, 406
- void in the distribution of vacuum energy, 397–399
- void in the matter distribution, 390–399, 407, 408
- zero level, 405
- initial density fluctuations, *see* primordial inhomogeneities, initial matter-energy distribution
- initial matter-energy distribution
 - compressions and rarefactions, 443
 - conventional smoothing processes, 274, 416–422, 440
 - Harrison-Zel’dovich’s spectrum of fluctuations, 443, 444, 451
 - homogeneity, 322, 332, 347, 361–371, 408–443, 448–450, 463, 572, 573, 580, 595, 603, 638
 - microscopic inhomogeneities, 332, 333
 - oscillating fluctuations, 443
 - primordial inhomogeneities, 276, 282, 333, 344, 345, 362–365, 409–418, 423, 428–430, 439–451, 574
 - scale-invariant spectrum of fluctuations, 344n, 443, 444, 451
 - smaller-scale fluctuations, 444
- instrumentalism, 22, 641
- interaction boson, 128–130, 153, 314–316, 431, 497
- continuity of the flow of time, 310–315
- decay, 305

- neutral interaction, 311, 317, 325
 - spin, 317n
 - spin-two graviton, 118n
- interaction vertex
 - continuity of the flow of time, 310, 311
 - mixed action signs, 129, 130
- interference pattern, 536
- interferometer experiment, 488, 539, 544
- intuitive leap, 508, 521
- invariance of the sign of action, 596
- irreducible randomness, *see* quantum chance

- Kent, Adrian, 556, 564, 565, 606
- kinetic energy, 137, 138, 150, 247
 - average kinetic energy, 333
- kinetic theory of gases, 381

- large-scale structure, 9, 18, 281, 301–303, 346–350, 416, 439, 636
 - bubble-like pattern, 348
 - computer simulations, 346
 - dissociation of the matter distributions, *see* polarization of the matter distribution, large-scale structure
 - fluctuations in matter distribution, 445, 625
 - homogeneous matter distribution, 127, 302, 333, 372, 412, 439
 - Local Sheet, 346
 - Local Void, 346
 - polarization of the matter distribution, 412–416, 438, 450
- larger-than-one probability, 531–533
- least-action principle, 493

- Lebowitz, Joel, 487
- light chemical elements, 282n
- Liouville’s equation, *see* Liouville’s theorem
- Liouville’s theorem, 377, 554
- local gravitational fields, 326, 387–398, 406–414, 428, 573–577
- local inertial reference system, 63, 77–90, 95–98, 496
 - energy sign dependence, 621
 - equilibrium of gravitational forces, 625–627
 - interaction with all matter in the universe, 626
 - Machian viewpoint, 625
 - unobservable fluctuations, 624, 625
- Lorentz transformation, 521

- Mach’s principle, 95
- Mach, Ernst, 72
- massive neutrino, 25
- material nature of gravitation, 82–86
 - electromagnetic field analogy, 83
- mathematics, 15
- matter creation, 113
 - Big Bang, 322
 - cosmic expansion, 322
 - favorable conditions, 123
 - gravitational radiation energy, 322, 637
 - observational evidence, 113, 123
 - out of nothing, 126, 127, 355, 368–371, 448, 587
 - quantum gravitational scale, 322
- Maxwell equations, 478, 479
- Maxwell, James Clerk, 478
- memory state, 579
- meson, 245

- missing-mass effect, 273–276, 300–309, 318–320, 326–335, 346, 347, 637
- modified gravitational dynamics, 336
- MOND, *see* modified gravitational dynamics
- multiverse, 21, 185, 451, 452, 472
 - causally independent universes, 472, 553, 554, 589, 590, 641
 - ensemble of possible universal causal chains, 590
- mutually consistent records, 564, 579, 580, 602, 610–614
 - long-term records, 579, 602
 - short-term memories, 579
- naked singularities, 441
- negative energy
 - advanced waves, 478
 - antiparticles, 28–40, 48–51, 102–105, 112–115, 128, 129, 143, 144, 157, 642
 - attractive force field, 46, 47, 114–117, 140, 153–156, 247
 - black hole, 396, 415, 428, 450
 - bound systems, 46, 47, 114–116, 140–142, 156, 181
 - concept, 454
 - Dirac’s solution, 28–30, 48, 101–105
 - energy conditions, 44, 45
 - filled energy continuum, 102–105
 - in general relativity, 10, 24, 74
 - in quantum field theory, 10–12, 25–28, 50, 51, 111, 121–123, 128, 133, 143, 642, *see* negative densities, vacuum energy
 - interaction constraint, 26, 121–123
 - kinetic energy, 149, 150
 - motivations, 10, 118, 272
 - myths, 151
 - negative action, 37–42, 52, 58, 111, 124, 134, 157, 224, 278, 458, 478n, 538, 633, 642, 643
 - negative frequencies, 25, 26
 - negative pressure, 273, 276, 288, 330
 - positive-energy theorems, 47, 48
 - propagation constraint, 25–30, 308
 - steady state cosmology, 31
 - the problem of, 8, 18–31, 44–51, 172, 182
 - traditional interpretation, 10, 43, 113–115, 134, 142, 147–151, 160, 643
 - transition constraint, 10, 29, 49, 123, 129
 - versus negative charge, 27
- negative mass, 52–90, 633
 - absolute gravitational force, 55, 643
 - absolute inertial mass, 139–142, 156
 - acceleration, 64–73, 84, 114, 139, 140, 156
 - generalized Newton’s second law, 64–77, 114, 139
 - gravitational mass, 52–62, 69, 74, 114, 139–142, 152
 - motivations, 8
 - negative inertial mass, 43, 52–76, 114, 139–142, 152–156, 177
 - Newton’s second law, 64
 - Newton’s third law, 55
 - Newtonian gravitation, 180
 - positive inertial mass, 60–62, 69

- principle of inertia, 53, 54, 60, 61, 76–79, 643
- traditional concept, 8, 43, 52–64, 156, 633, 643
- negative probability, 531–534, 540, 565, 566, 585, 597, 598, 640
- negative temperatures, 149, 150, 269, 270
 - black hole, 634
 - decrease of entropy, 269, 270, 396
 - energy levels, 269–272
 - heat exchange, 405
 - infinite temperature, 269–271
 - Kelvin scale, 270n
 - negative energy, 270–272
 - negative heat, 270, 394–397, 404, 405
 - negative-energy black hole, 270–272, 394–396, 404, 405
 - positive heat, 396, 405
 - spin system in magnetic field, 269–271
 - void in the positive-energy matter distribution, 396
- negative-action matter, *see* negative-energy matter
- negative-energy matter, 8, 9, 277, 454, 463, 473, 519, 548, 558, 574–576, 600–603, 615, 621
 - absence of baryonic matter, 28, 122, 303, 324, 325, 342–348, 417, 444, 636, 643
 - absence of interactions with, 117–130, 135–138, 144, 153, 154, 163, 177, 178, 283, 296, 305, 407, 431, 643
 - accumulation, 27, 28
 - annihilation to nothing, 126, 127, 144
 - antimatter, 222, 321–324
 - antimatter experiment with, 144, 145
 - baryonic matter, 303, 324, 329, 342, 343, 438
 - Big Bang, 352, 428
 - colliding opposite-energy bodies, 134–138, 178
 - concentrations, 122
 - conservation of energy, 112–114, 129–138, 144–150, 360
 - conservation of momentum, 114, 123, 134–138, 366n
 - cosmological models, 175, 180
 - current average density, 302, 318
 - dark matter, 8, 28, 116–121, 153, 350
 - discrete symmetries, 222–227, 239, 268
 - dominant paradigm, 116, 643
 - force field energy, 116, 153
 - gravitational potential energy, 137, 138, 145, 178
 - heat, 149, 272
 - homogeneous distribution, 107–109, 154, 155, 165, 179, 302–304, 331, 332, 341, 352, 360, 367, 407, 413, 430, 444
 - hot dark matter, 347
 - implicit assumptions, 115, 299
 - inhomogeneities, 122, 154, 304, 329, 347, 360, 411, 430, 444–446, 583
 - momentum direction, 88, 123, 134
 - negative heat, 149–151
 - neutrinos, 347
 - non-baryonic, 122

- nonexistence theorems, 142n
- observational evidence, 25–29, 58, 111, 116, 121, 122, 636, 643
- outstanding problems, 111–115, 142–151, 160, 178, 634, 643
- overdensity, 91, 154, 175–179, 281, 301, 325, 330, 342–346, 404, 643
- pairs of opposite-action pairs, 126
- positive-definite energy exchanges, 54, 138, 149
- potential energy, 145, 146
- primordial abundance, 122, 302, 303
- rarity, 122, 123
- ratio of average densities, 180, 285, 290–295, 302, 303, 324, 355, 367–372, 438, 643
- requirement of exchange symmetry, 11, 42, 58–60, 131, 151–154, 160, 167, 175–177, 228, 278–284, 299, 318, 636–639
- reversed bidirectional charge, 343
- temperature, 149–151
- thermal energy, 149, 150, 405
- traditional concept, 8, 282n, 352, 413
- transformation into, 124
- underdensity, 91–93, 154, 165–169, 175–179, 301–305, 318, 326–330, 341, 346–348, 404, 636, 637
- uniform distribution, *see* homogeneous distribution, negative-energy matter
- universal expansion, 180
- vacuum energy contributions, 278, 279
- voids in positive vacuum energy, 99–109, 119–126, 135, 136, 153–155, 165–178, 283, 296, 297, 318, 331, 351, 360, 367, 393, 403–405, 413, 496, 633–637, 643
- voids in vacuum charge distribution, 120, 153, 331, 496, 637
- neutrino, 308
- neutron’s electric dipole moment, 208
 - direction of dipole, 208
 - precession movement, 208
- Newton, Isaac, 492
- Newtonian physics, 574
- nucleus, 255
- observer selection effect, 365–373, 382, 412, 421, 447, 448, 635–638
- observer-dependent average densities, *see* specific densities
- observer-dependent expansion rates, *see* specific expansion rates
- Omnès, Roland, 558–561, 567
- open questions, 17, 633
- opposite-action particles
 - pair annihilation, 112, 113, 123–128, 144, 305
 - pair creation, 112, 113, 123–127, 305
- opposite-bidirectional-charge particles
 - absence of non-gravitational interactions, 322, 637
 - creation and annihilation, 323, 343, 637
 - gravitational interaction, 322–325, 637
- pairs of linearly polarized photons, 545n

- angles of polarization, 545n
- particle accelerator, 244
- particle beam, 244
- particle horizon, *see* cosmic horizon
- particle physics, 14–17, 374
- Penrose, Roger, 375, 408, 450, 566
- perception of the passage of time, 579–581
- perpetual motion problem, 115, 143–146
- Petit, Jean-Pierre, 160n
- philosophy, 13–15, 641
- photon’s angle of impact, 539–548
- physical nature of geometry, 84
- Planck, Max, 24, *see* quantum gravitation
- polarity, 458
- polarization, 215
- Popescu, Sandu, 537
- principle of equivalence, 43, 74–90, 494
 - Einstein’s elevator experiment, 65, 76–83, 88
 - entangled system, 141, 142
 - equivalent source, 65–68
 - relativized, 74, 88–90, 156, 181
 - violation of the, 61, 74–81, 87–91, 114, 139–142
- probability distribution, 509
- probability of occurrence, 484–489, 499, 528–533, 540
 - negative contribution, 528, 532, 533
 - whole universe history, 597
- propagator, 199
- quantization hypothesis, 24
- quantization of electromagnetic radiation, 480
- quantum chance, 379–383
- quantum chromodynamics, 245
- quantum cosmology, 10, 556, 562, 570, 575, 576, 581–588
 - ADM formalism, 584
 - boundary conditions over superspace, 585
 - canonical, 571, 581, 598
 - circular history, 592, 597–599
 - closed trajectory in superspace, 590–592
 - configuration space, 581, 586, 587
 - diverging superspace trajectories, 604, 605, 613
 - dynamic elements, 576
 - equivalent superspace trajectories, 584
 - global state, 581–590
 - histories of space curvature, 612, 613
 - identity of initial and final boundary conditions, 597
 - irrelevance of space, 572n
 - irrelevance of time, 570–572, 577–584, 598
 - meeting of superspace trajectories, 593–596, 604–610
 - monotonic foliation of space-like hypersurfaces, 581
 - network of local relationships, 584
 - non-gravitational degrees of freedom, 586
 - pair of indistinguishable trajectories in superspace, 590
 - periodicity, 599
 - periodicity of time, 598
 - point in superspace, 585, 591–594
 - predetermined relationships, 576

- relevance of time, 575, 576, 583
- rotating-clock-hand analogy, 598
- shared coarse-grained superspace trajectory, 594
- superspace, 581–587, 603
- timeless, 578n
- topology of superspace trajectory, 599
- trajectory in superspace, 575, 581–599, 604, 613, 624
- uniqueness of history, 580–584, 590, 591, 598
- Wheeler-DeWitt equation, 571
- quantum field theory, 19, 20, 278, 279, 309, 383, 472, 479, 480, 485, 495–498, 503, 504, 558, 575–577, 623
 - conventional formulation, 640
 - Dirac’s equation, 482
 - fermion loops, 558
 - Feynman diagrams for radiative corrections, 630
 - fixed geometry of space, 625
 - position-dependent interaction probability, 497n
 - radiative correction terms, 558
 - reformulation, 627
 - relativistic invariance, 490, 606
 - renormalization procedure, 558
- quantum formalism
 - complex conjugation as time reversal, 487
 - consistent histories, 487, 530, 561–565, 606–610
 - extension of formalism, 616
 - Heisenberg’s matrix mechanics, 503n
 - incomplete description, 620
 - linear equations, 550
 - linearly positive histories, 565
 - most adequate, 512
 - path integrals, 503
 - relation to gravitation, 615
 - Schrödinger equation, 489–491
 - squaring of the wave function, 487, 529
 - stationary Schrödinger equation, 571
 - sum-over-histories formulation, 490, 499–505, 512, 521, 624
 - time-reversal invariant, 565
 - time-symmetric formulation, 10, 307n, 520, 538, 559, 570, 597–599
 - traditional formulation, 503, 625
 - two-state-vector formalism, 487–490, 519–521, 547
 - unitarity, 505n
 - unitary evolution, 487, 519, 558–561
- quantum gravitation, 12, 171n, 229–232, 250, 268, 299, 338, 371–374, 388, 389, 425–427, 454, 595, 612–614, 623–629
 - area gap, 236n, 339
 - Ashtekar variables, 584n
 - background-independent theory, 570, 575, 625–627
 - Big Snap, 402n
 - causal structure of spin foams, 572, 573
 - covariant theory, 575, 584n
 - discontinuous time flow, 580
 - discrete elements of structure, 576
 - discrete space, 232–234, 277, 584n
 - elementary black hole, 9, 236–238, 245, 252, 257, 264, 265, 336–

- 340, 369, 386, 400, 423, 430, 431, 621, 629, 634, 644
- elementary unit of area, *see* elementary unit of surface, quantum gravitation
- elementary unit of space, *see* elementary unit of surface, quantum gravitation
- elementary unit of surface, 233–236, 252, 257, 265–267, 338–340, 351, 385, 386, 400–402, 402n, 425–431, 637, 638, 644
- embryonic concept of time, 570
- embryonic element of causal order, 573, 580
- embryonic element of time directionality, 572
- embryonic notion of space, 576
- emergence of continuous time, 583
- emergence of space, 576
- emergence of time, 10, 570–575, 583, 621
- energy fluctuations, 267, 276–280
- fluctuating gravitational field, 262n, 264–267, 613, 614
- fluctuating metric properties of space, 613, 614
- four fundamental degrees of freedom, 339, 637
- four-dimensional boundary conditions, 576
- fundamental element of causality, 575
- fundamental scale, *see* Planck scale, quantum gravitation
- fundamental unit of surface, 339
- graviton decay, 325
- gravitons, 236, 317
- irrelevance of time, 571, 577–579
- loop quantum gravity, 232n, 236n, 248–250, 339, 584n, 621
- maximum energy densities, 332, 351, 352, 370–373, 413, 635–638
- microscopic black hole, 264–268
- minimum distance, 247, 431
- minimum time interval, 245, 351
- momentum orientation, 265, 266
- negative energy, 267
- non-unique spacetime metric signature, 580, 613
- opposite-energy solutions, 575
- phenomenological unit of area, 339
- Planck area, 235, 236, 257, 265, 339, 340, 427
- Planck energy, 171–175, 236, 245–247, 264, 277, 322, 339, 427, 428
- Planck length, 235, 245, 261, 339, 425, 426
- Planck mass, 236, 264, 271, 431
- Planck momentum, 247
- Planck scale, 12, 89, 171–173, 232–234, 236n, 241–245, 249n, 250–253, 259, 264–267, 280, 299, 332, 339, 369, 570–575, 588, 613, 614, 623–629
- Planck time, 233, 372, 425–428, 473
- probability of graviton emission, 626
- quantization of space, *see* discrete space, quantum gravitation
- quantization procedure, 623
- quantum bounce, 127, 248–250, 255–257, 321, 336, 354, 368–

- 371, 424, 436–438, 448, 490, 610, 637
- relevance of time, 576
- scale of distance, 322, 338, 339, 425, 426, 637
- semi-classical approximation, 236n
- semi-classical description, 234, 339, 340, 495, 577, 627
- spin network, 576
- spin network edges, 232n
- spin-foam quantum gravity, 575, 584n
- superposed space curvatures, 614
- timeless theory, 12, 577–581
- unit of distance, *see* Planck length, quantum gravitation
- unit of time, *see* Planck time, quantum gravitation
- unsolved issues, 574
- quantum gravitational scale, *see* Planck scale, quantum gravitation
- quantum handshake process, 484–486, 535
 - confirmation wave, 485
 - offer wave, 485
- quantum measurement
 - absence of interference, 525, 630
 - absence of quasiclassicality, 568
 - actualization of potentialities, 384, 509, 527, 543–560, 576, 619–622
 - amplification process, 558
 - classical outcome, 600, 609, 616
 - coarse-grained history, 533, 561–567
 - coarse-graining, 529, 530, 564
 - consciousness, 569
 - conspiracy theory, 554
 - correlated attribute, 544, 545, 545n
 - correlated effects, 604
 - decoherence, 456, 506, 514, 530, 543, 549–569, 582, 592, 597–606, 612–619, 630, 640
 - decoherence effect, 384
 - decoherent branches, *see* splitting branches, quantum measurement
 - decoherent space and time, 578–583, 614
 - delocalized phase relations, 557
 - diagonalization of reduced density operator, 551
 - discontinuous, 497, 509
 - distinct dynamical laws, 552–556
 - emergence of quasiclassicality, *see* persistence of quasiclassicality, quantum measurement
 - ensemble entropy, 554
 - entangled particles, 528, 539–548
 - environment degrees of freedom, 456, 543, 549, 556–567, 589, 600–608, 614
 - EPR-type experiment, 535–538, 544–548
 - essential condition, 604–607
 - evidence of past quasiclassicality, 564–568
 - expectation of future quasiclassicality, 564–568
 - factual definiteness of reality, 567
 - family of coarse-grained histories, 530, 563, 606–610
 - gravitational field, 612
 - graviton perturbations, 615
 - identically prepared systems, 482, 513, 617

- interaction-free, 485, 562
- irrelevance of complexity, 558
- irreversibility, 384, 456, 557–559, 558n, 566, 567, 600–602
- likeliness of alternative conditions, 532
- likeliness of initial conditions, 532, 533, 540, 598
- likeliness of macroscopic conditions, 532, 533, 540
- logical consistency of history, 558, 565
- macroscopic experimental constraints, *see* observable macroscopic conditions, quantum measurement
- maximum quasiclassical domain, 10, 19, 455, 514, 558–564, 570, 593, 601–606, 615, 639
- meaningless probabilities, 534, 622
- measuring device, 556, 568, 569, 604
- metric properties, 612
- mixed quantum state, 561
- mutually exclusive macroscopic constraints, 513
- non-local, 497–500, 514, 546, 555
- non-subjective, 557
- non-superposed outcomes, 557
- observable macroscopic conditions, 483–490, 487n, 495, 502–513, 522–543, 557–562, 582, 588–598, 605–615, 620–624
- observable outcome, 458
- observed history, 534
- observed physical attribute, 502–514, 542, 551, 557–561, 582, 606
- observer, 569
- observer awareness, 486
- observer-dependent knowledge, 543
- origin of randomness, 511, 619
- overturning, 600
- persistence of quasiclassicality, 10, 19, 455–459, 498, 549–570, 575, 583, 592, 599–601, 606–614, 640, 641
- physically relevant set, 562–566, 606
- position measurement, 547
- post selection, 487–489, 511, 519, 520, 537, 547, 548, 555
- preexisting correlation, 546
- preferred basis, 606
- probability of future outcome, 533
- probability sum rules, 565
- quasiclassical gravitational field, 582
- quasiclassical reality, 456–459
- quasiclassical world, *see* maximum quasiclassical domain, quantum measurement
- quasiclassicality as an evolutionary advantage, 564
- random outcome, 516, 615–619, 630, 631
- relative time order, 520, 545, 548n
- relevant collective observable, 606
- Schrödinger’s cat experiment, 607, 608
- splitting branches, 471–476, 514, 527, 546, 551–555, 609
- spreading of effects, 600
- state vector reduction, 520, 521, 546, 555–559, 566–569, 619, 630
- statistical distribution of outcomes,

- 506, 523–526, 536, 567n
- stochastic perturbations of the Schrödinger equation, 630
- subjectivity, 543
- subjectivity of quasiclassicality, 641
- summed-over aspects, 561–567
- superposition of macroscopic observables, 550–552, 562–567, 567n, 607, 608
- time symmetry, 558
- unique datum, *see* unique outcome, quantum measurement
- unique outcome, 615–617, 622
- unique preexisting state, 554
- uniqueness of experimental facts, *see* uniqueness of measurement results, quantum measurement
- uniqueness of measurement results, 492, 501, 507–518, 527, 551–557, 568, 631
- universe as a whole, 556
- unnatural coincidences, 546
- unpredictability, 559, 560, 617
- variable outcomes, 511–513
- quantum reality, 483, 631
 - absence of causes, 617–619
 - alternative conception, 495
 - backward-evolving process, 490
 - backward-evolving state, 488, 489
 - causal determination, 616–622, 627–631
 - classical aspects, 530, 531n
 - classical hidden variables, 535, 536, 554
 - conjugate physical attributes, 502–513, 554, 561, 591, 612
 - contradictory nature, 494, 507–517, 525, 526, 548n, 551–553
 - correlations, 537
 - created by observation, 511, 639
 - criterion of consistency, 498, 530, 531, 548n, 563–566, 592, 601
 - decoherent branches of history, 513
 - dependence on experimental conditions, 502
 - determined by measurement, 498
 - dynamic attribute, 499, 507, 515, 550, 564, 589, 603–606
 - electron spin state, 525, 526
 - energy eigenstate, 630
 - energy state, 497
 - entangled pair, 515, 520, 521, 528, 538–546, 620
 - entangled state, 539–543
 - entangled systems, 516, 528, 536–539, 546
 - entanglement, 10, 455, 495, 514–519, 535, 536, 545–548, 599, 615, 620, 621, 644
 - expanding spherical wave function, 546, 547
 - fine-grained history, 561–565
 - forward-evolving state, 489
 - fundamental hypothesis, 483
 - fundamental randomness, 470
 - future entanglement, 548
 - generalized consistency conditions, 565
 - history, 492, 560
 - holistic, 545
 - inadequate representations, 484
 - initial value of phase, 531–533, 540
 - interferences, 455, 482–495, 501–509, 515–518, 523–534, 539–544, 545n, 550–569, 594–615, 621, 628, 629, 640

- interfering histories, 471, 488–490, 497–518, 523–530, 539–547, 553, 565, 571, 597–599, 611, 631, 641
- interfering states, 499, 549, 559
- irreconcilable requirements, 503
- irreducible randomness, 619, 630
- large-scale interferences, 600
- locality assumption, 517
- logical consistency, 531–533, 563–566, 587, 588, 618
- meaningful aspects, 530
- minimally coarse-grained histories, 487, 530–533
- momentum eigenstate, 497, 506
- momentum state, 495, 505–509
- multiple coexisting trajectories, *see* multiple-branches hypothesis, quantum reality
- multiple-branches hypothesis, 508–528, 546, 551–559, 567, 585, 590, 591, 609–616, 622
- mutually exclusive representations, 505
- naive conception, 507, 515, 536
- non-classical nature, 22
- non-classical uniqueness, 621
- non-interfering histories, 566
- non-local correlations, 455, 515–521, 526, 535–549, 539n, 599, 609, 620
- non-local hidden variables, 499–501, 511–518, 535
- non-locality, 455–459, 495, 500, 514–521, 526, 535–539, 546–549, 555, 586, 599, 615, 644
- non-objective, 517
- objective chance, 501, 617–619
- objective indefiniteness, 501, 616, 617
- objective reality, 22, 526, 620
- observer-dependent facts, 517
- pair of causally independent histories, 489, 490, 521–527, 535, 546, 590–592, 597, 605–607, 620, 640
- pair of histories unfolding in opposite time directions, 487–490, 521–533, 543–546, 586–590, 640
- pair of interfering histories, 508, 561, 586, 592, 620
- particle trajectory, 467, 492–499, 523–525, 624–626
- periodic evolution, 497
- phenomenological model, 508
- position state, 495–498, 509, 518
- predetermination, 535
- preestablished harmony, 599
- probability amplitude, 455, 484–487, 497, 497n, 512, 518, 523–534, 554–561, 640
- quantum field, 497
- quantum indefiniteness, 425, 506
- quantum indeterminacy, 510, 591
- quantum phase, 512, 531–533, 539–541, 557, 628
- quantum state, 470, 487, 525, 546, 623
- quantum strangeness, 505–508, 526, 531, 555
- quantum superposition of random space curvatures, 630
- quantum uncertainty, 425, 495, 505, 510, 628
- randomness, 487n, 511, 559, 560, 576, 594, 617–625

- realist conception, 22, 495–525, 530–538, 576, 620, 641
- realist description, *see* realist conception, quantum reality
- relational description, 568
- relative notion, 517
- retarded and advanced portions
 - of history, 483, 489, 490, 522–548, 555–557, 566, 582–631
- retarded and advanced states, 542, 554, 590–595, 604–609, 630
- shared quantum phase, 543, 597
- state indefiniteness, 506, 620
- state superposition, 509, 518, 525, 536, 542, 550–554, 561, 567, 590
- state vector, 487n, 489, 490, 509, 526, 550, 590–594
- subjective wave function, 543
- time symmetry, 527, 552–559, 567, 609
- time-symmetric conception, 530, 583
- time-symmetric history, 533, 615, 640
- time-symmetric process, 528–541, 598, 606, 629
- understandable, 493
- undetermined state, 525, 561
- unique history, 498–515, 521, 522, 531–535, 559, 590, 615, 616, 624, 631
- unique trajectory, 495–499, 505–512, 521–525, 541–546, 563, 616, 617
- uniqueness, 499–502, 513, 514, 526, 527, 533–537, 550–560, 568, 591, 615–625, 630, 631, 640, 641
- uniqueness of historical facts, *see* uniqueness of measurement results, quantum measurement
- unobservable aspects, 22, 23, 490, 498–513, 531–536, 566–568, 585, 590–595, 617, 623–630
- unobservable causes, 617–623
- unobservable hidden variables, 501
- unobservable history, 513
- unobservable state, 526, 594
- unobservable trajectory, 495–499, 506, 523, 529, 544, 563, 625–630
- unobserved path, 499, 522–525, 530–537, 547, 559, 617–624, 631
- unobserved physical attribute, 498–513, 518–526, 533–536, 545–554, 591–594, 616–624, 630
- unobserved portion of history, 533–536, 591, 617–619, 625
- unpredictability, 470, 501, 559, 619–621, 631
- vacuum fluctuations, 495, 496
- virtual processes, 503n, 630
- visualization, 505–510, 525
- wave function, 482, 487n, 489, 493n, 497–500, 509–519, 531, 543–548, 559, 560, 619, 620
- wave function phase continuity, 597, 598
- wave function phase invariance, 598, 629
- wave packet, 495
- wave-particle duality, 497, 498
- wavelike nature, 497
- quantum theory, 8–19, 388, 393, 429, 432, 480

- absolute determinism, 535, 546
- alternative interpretations, 483–489, 502, 509, 518, 526, 535, 536
- basic principles, 12
- central problem of interpretation, 568
- classical hidden-variables theory, 485, 500–502, 507–511, 518, 519, 535, 546, 615, 620, 621
- complementarity principle, 505, 530
- consistent interpretation, 492–501, 507–514, 519–527, 541, 549, 565, 570, 583, 590, 623, 639, 640
- consistent-histories interpretation, 498, 513, 530, 531n, 548n, 556, 561–563, 570, 592, 601
- constant and uniform spacetime background, 623–625, 630, 640
- conventional interpretation, *see* orthodox interpretation, quantum theory
- conventional quantum mechanics, 503, 612
- Copenhagen interpretation, 501–505, 563, 641
- cosmological theory, 583
- counter-intuitive aspect, 491–493, 587, 618
- currently favored interpretation, 509, 568
- domain of validity, 482
- epistemological viewpoint, 493
- equivalent processes, 316
- existing interpretations, 456, 492–494, 528
- flat invariant spacetime, 623
- foundations, 631
- fundamental time asymmetry, 484–489
- hidden-variables theory, 470
- idealistic position, 548
- incomplete interpretations, 498, 548, 568
- incompleteness, 570, 615, 616
- inconsistent interpretations, 492–494, 510–518, 622
- instrumentalist interpretation, 511–514
- interpretation problem, 10–13, 19–23, 455, 456, 491–494, 503–508, 521, 535n, 538, 556, 582, 631–633
- irreversibility, 383
- known interpretations, 569
- lack of intelligibility, 491, 506
- many-worlds interpretation, 471–476, 509, 527, 552–562, 567, 599, 622, 641
- mathematical framework, 455, 456, 487–502, 513, 529, 556–560, 569, 570, 592, 644
- mathematical structure, 483, 518, 519
- measurement problem, 455–458, 472, 486, 498, 509, 538, 550–556, 562–570, 592, 601, 640, 641
- naive realist interpretations, 620
- non-realist interpretation, *see* instrumentalist interpretation, quantum theory
- objectification problem, 527, 631
- ontological viewpoint, 493
- orthodox interpretation, 10, 455,

- 486, 491, 498–503, 508–511, 516–518, 530, 538–541, 548–555, 568, 600, 621
- probabilistic inferences, 563
- probabilistic nature, 491, 622
- QBism, 543
- quantum fields in curved spacetime, 612
- quantum logic, 494
- quantum reality problem, 455, 456
- quantum scale of action, 627
- realist interpretation, 10, 22, 455–459, 488–490, 498–512, 518–521, 526, 527, 533, 538, 546–551, 559–563, 584n, 617, 631, 641
- realist time-symmetric interpretation, 516, 527, 528, 533–540, 535n, 546, 547, 563–566, 590, 601, 616–621
- reformulation, 616, 621–623
- relational interpretation, 517, 568
- revised interpretation, 483, 491, 492
- satisfactory interpretation, 562, 583
- selection principle, 606
- standard theory, 486–488, 519–522
- state-of-the-art interpretation, 565
- statistical predictions, 482–485, 490, 493n
- time-symmetric, 468, 483–490, 508–512, 519, 520, 526, 535, 536, 588
- time-symmetric equations, 384
- time-symmetric interpretation, 526, 590, 599, 623
- traditional interpretation, *see* orthodox interpretation, quantum theory
- transactional interpretation, 485–487, 535
- uncertainty principle, 492, 497, 507, 513, 623
- unsettling viewpoint, 525
- quantum unitarity, 377
- quark, 245, 308, 314–316
- real probability, 534
- reformulated quantum theory
 - causally-determined geodesics of a random spacetime, 630
 - classical gravitational field fluctuations, 624–631
 - classical spacetime continuum approximation, 626, 627
 - dependence of uncertainty on energy magnitude, 628
 - dependence of uncertainty on spatial scale, 629
 - energy exchanges with the fluctuating gravitational field, 628
 - fluctuating equilibrium of gravitational forces, 626–628
 - fluctuating spacetime background, 623
 - massless particles with curved trajectories, 630
 - massless particles with non- c velocities, 630
 - random local space curvature, 624–627
 - relevance of particle masses, 628
 - statistically equivalent theory, 617, 630, 640
 - straight trajectories in a locally curved spacetime, 630

- temporary violations of energy conservation, 628
- Reichenbach, Hans, 84, 468, 584, 599
- relativistic frame dragging, 78, 85, 86, 244
- relativity of particle existence, 628
- repulsive force field
 - energy sign of, 116, 117
- rest mass, 247
- Rovelli, Carlo, 248n
- Rutherford atom model, 130

- Schrödinger, Erwin, 491
- scientific realism, 10, 22, 23, 499, 516–518, 526, 548, 564, 616, 640
- scientific research, 15–19, 639
- second law of thermodynamics, 32, 275, 438, 458, 581
 - adjustment of initial conditions, *see* microscopic state preparation, second law of thermodynamics
 - anti-thermodynamic evolution, 376–380, 415, 422, 436, 476
 - branch systems, *see* isolated systems, second law of thermodynamics
 - Clausius' definition of entropy change, 394
 - coarse-graining, 260, 375–377, 384–389, 638
 - constraint on process description, 189
 - contraction of space, 414–418, 435–438
 - degradation of energy, 131–133, 145
 - energy, 270, 271
 - entropy, 131, 149–151, 189–195, 218, 259, 269, 270, 375–394, 405–409, 429, 430, 436, 461–467, 473–477, 482, 489, 490, 516–519, 543–547, 556–558, 579–581, 586–595, 600, 607–610
 - entropy-decreasing fluctuation, 377–382, 415, 438, 474–477, 587, 594, 609
 - expansion of space, 376, 385, 402, 432–438
 - gravitational entropy, 9, 236, 258–267, 276, 375, 385–397, 404–422, 428–438, 450, 473, 480, 490, 496, 548, 558, 576, 581, 589–595, 600–603, 608n, 638
 - growth of inhomogeneities, 334, 385, 388, 408–411, 417–420
 - heat death, 438
 - homogeneous final state, 420
 - inhomogeneity of the matter distribution, 411
 - inhomogeneous final state, 418
 - inhomogeneous initial state, 419–421
 - isolated systems, 377–381, 433, 473
 - macroscopic parameters, 375, 376, 386
 - macroscopic state, 430
 - matter disintegration, 131, 132
 - matter entropy, 385, 391, 420, 430
 - microscopic configuration, *see* microscopic state, second law of thermodynamics
 - microscopic degrees of freedom, 131, 375–379, 386–392, 595, 603
 - microscopic parameters, 375
 - microscopic state, 336, 375–379,

- 386–388, 396, 408–411, 429, 430
- microscopic state preparation, 195, 379–383
- negative entropy, 404
- non-equilibrium state, 375–378, 433
- objective entropy growth, 231, 260, 261, 375–377, 385–393, 638
- particle motion reversal, 379
- Poincaré return or recurrence, *see* entropy-decreasing fluctuation, second law of thermodynamics
- reduction of inhomogeneities, 418
- smoothness of matter distribution, 266
- static equilibrium, 267
- subjective entropy growth, 375, 388, 578n
- temperature, 269, 270
- thermal equilibrium, 131, 263–268, 376–381, 394–397, 408, 430, 438, 475
- unique coarse-graining, 260, 386
- unlimited entropy growth, 600
- violation, 148–151, 195, 248, 473, 474, 586, 609, 623
- sensibility to initial conditions, *see* classical instability
- sign of energy, 538, 596
- Smolin, Lee, 625
- solipsism, 22, 499, 564, 565, 577, 606, 641
- space-like separated events, 520, 539
- relative time order, 464
- spatial curvature parameter, 356–363
- special relativity, 186, 464
- causal time, 466
- space-like interval, 426
- time-like interval, 466, 613
- specific densities, 293–296, 351–355, 372, 401, 449
- specific expansion rates, 298, 299, 324, 332, 352–372, 400, 401, 447–449, 636
- spin angular momentum, 315
- sign of action, 256
- spin density of matter, 256, 400n
- spreading wave front, 603
- Stükelberg, Ernst, 49, 464
- state vector, 199
- static force field, 393n
- statistical mechanics, 15–19, 263–267, 374–377, 383, 501, 512, 550
- Brownian motion, 625
- equilibrium thermodynamics, 267, 272
- multiple and near simultaneous interactions, 625, 626
- near-equilibrium thermodynamics, 267, 268, 627
- non-equilibrium thermodynamics, 267, 268, 383
- Planck's definition of entropy, 480
- quantum field theory, 455
- thermal equilibrium, 634
- strong nuclear interaction, 140
- structure formation, 274–276, 302, 303, 341–350, 365, 409, 444, 636
- bottom-up process, 342
- conventional models, 348, 636
- earliest galaxies, 303, 341, 346, 636
- elliptical galaxies, 341
- gravitational instability, 342–347, 410, 417

- kinetic energy reduction, 342
- radiation emission, 342
- rate acceleration, 344n, 346, 347, 445, 636
- supernovae, 282
- Susskind, Leonard, 451
- symmetry-breaking phase transitions, 441
- ‘t Hooft, Gerard, 235
- teleological problem of time, 467
- thermal time, 578
- thermodynamic time asymmetry, *see* time irreversibility
- thermodynamics, 23, 374, 578
- tidal effect, 244
- time as a reference system, 577
- time asymmetry, *see* time irreversibility
- time directionality, 16–19, 455, 460–462, 577, 583
 - local topological order, 469
- time irreversibility, 148, 419, 438, 458–465, 463n, 489, 515, 550, 556–569, 574–577, 603–606, 608n, 611–614, 624, 639
 - absence of records of future, 610
 - astronomical processes, 463n
 - backward teleology, 432
 - backward-in-time propagation, 32, 192, 644
 - Boltzmann’s solution, 381, 382
 - boundary conditions, 462, 482
 - conditions on current state, 433, 434, 479
 - cosmic evolution, 418–420, 435–438, 482
 - cosmological arrow of time, 376, 411, 436
 - derived property, 434
 - dissipation, 463n, 557, 578, 602–606
 - electromagnetic arrow of time, 481
 - favorable conditions, 409, 421
 - formation of records, 32, 190, 411, 436, 463, 473–475, 502–506, 523–525, 557, 558n, 602–614
 - forward-in-time viewpoint, *see* unidirectional time, time irreversibility
 - from state preparation, 433
 - fundamental irreversibility, 231, 383, 384, 434, 463–467, 475, 482, 556, 584, 600
 - information flow, 32, 461–466, 473–481, 488, 516, 543, 619
 - interaction with radiation, 480
 - irreducible time asymmetry, *see* fundamental irreversibility, time irreversibility
 - irreversible processes, 267, 268, 376, 384, 585
 - long-lasting history, 382
 - objective notion, 261, 375, 376, 384
 - origin, 9–12, 19, 25, 195, 276, 375, 381–388, 408–411, 421, 422, 432–437, 456, 473, 490, 638
 - preferred direction of time, 421–437
 - psychological arrow of time, 381, 382, 411
 - reversal of the arrow of time, 420
 - singular status of position space, 603–606
 - temporal parallelism, 433, 434, 463n,

- 473
- thermodynamic arrow of time, 9–12, 19, 25, 26, 131, 190, 195, 206–210, 275, 320, 376–384, 408–413, 423, 424, 434–438, 447, 452, 460–475, 490, 516, 550, 574, 581, 593–595, 609, 633, 638
- thermodynamic time, 189, 481, 644
- unidirectional time, 183, 189–224, 241–244, 252–254, 305–311, 316, 375, 414–419, 424, 454, 459–482, 519, 528, 578–581, 586–600, 605, 614, 644
- unidirectionality, 383, 384, 427
- wave retardation, 480
- wavefront propagation, 381, 479–482
- wavelike processes, 480
- time travel, 115, 146, 147, 474–481, 528, 581, 586
- causality violation, 147–151, 470–475
- continuous information flow, 587
- contradictory accounts, 471–476
- knowledge paradox, 474, 586
- paradox, 470–477, 609
- time-direction degree of freedom, 9, 31–43, 189–195, 577, 633, 644
- antiparticles, 190, 210–213, 224, 306–310, 321, 461–467, 473, 482, 519, 529, 538, 587, 644
- backward-in-time-propagating particle, 12, 26, 31–42, 49, 124, 190–194, 210, 224, 306–311, 316, 321–324, 644
- backward-propagating particle, 460–482, 519, 529, 587, 596
- bidirectional time, 189–219, 225, 242, 254, 311–316, 436–438, 454, 480, 593–596, 610, 621
- chronological order, 187, 192
- condition of continuity of the flow of time, 41, 124–126, 305–317, 324, 325, 586n, 637
- direction of propagation in time, 9, 31–40, 124, 129, 157, 182–211, 218–224, 243, 305–316, 321–325, 337–340, 354, 460–468, 522, 538, 596, 642–644
- Feynman’s interpretation, 33, 49, 460–464, 644
- particle world-line, 124, 306–315, 469
- relativity of the sign of charges, 31–41
- relativity of the sign of energy, 31–37, 58, 642
- reversal of action, 124, 157, 181, 305
- reversal of energy, 124, 129, 157, 305
- time-direction-dependent property, 191–194
- time-symmetric viewpoint, *see* bidirectional time, time-direction degree of freedom
- uniquely ordered sequence of events, 577
- time-symmetric physical laws, 378–384, 409, 429–432
- time-symmetric stochastic gravitational field theory, 627
- Tollaksen, Jeff, 537
- topological defects, 441
- cosmic strings, 441

- magnetic monopoles, 441
- transition probability, 497, 505n, 508, 522, 523, 529–534, 558, 567, 608, 612
- unidirectional variable, 244
- unification scale, 243n, 252
- unified theory of interactions, 247n, 338, 339, 440
- uniform matter distribution
 - underdensity, 273
- universal causal chain, 459, 581, 586, 587, 611–613, 624–627, 640
 - bifurcation point, 593–596
 - closed causal chain, 587–601, 606–615, 620, 629, 640
 - closure requirement, 459, 593–616, 630, 639, 640
 - constraint of non-divergence of superspace trajectories, 605, 606
 - end of time, 593, 594, 600
 - interrupted trajectory, 593
 - local topological ordering properties of spacetime, 584
 - non-decreasing probability of closure, 600
 - parallel stretching of superspace trajectories, 605
 - relativity of direction of propagation in time, 593
 - relativity of past and future, 593
 - relevance to quantum cosmology, 579–582
 - sequential order of events, 576, 583
- universal force, 82–84
- universe
 - absence of causal relationships, 423, 424
 - backward-in-time evolution, 416, 465
 - Big Crunch, 414–418, 428, 435
 - Big Crunch singularity, 414–424, 449
 - boundary conditions, 411
 - causal horizon, 572–574
 - causal relationships, 21, 423–433, 440, 448–454, 472, 473, 482, 501, 510–514, 558, 572–581, 588–590, 610, 638
 - causal self-determination, 588–590, 607–609
 - causal structure, 586
 - closed, 285, 352, 358, 364–369, 416
 - co-moving volume, 402
 - collapsing, 401
 - conservation of energy, 356, 357, 594
 - cosmic time, 426
 - creation out of nothing, 353, 354, 359, 369, 420, 421
 - deceleration of expansion, 290–292
 - extended vacuum state, 420–422, 435
 - final singularity, *see* Big Crunch singularity, universe
 - flat space, 365–373, 366n, 635
 - four-dimensional, 572–575
 - geometry, 361, 635
 - global space curvature, 282, 373
 - global time symmetry, 420, 438
 - gravitational energy, 354–374, 412, 423, 635
 - gravitational energy of curvature, 356–363
 - gravitational momentum, 354, 366n
 - gravitational potential energy, 352–

- 357, 365–368
 gravitational potential energy of
 matter, 356–374, 449
 history, 514, 587–592, 597
 incomplete instance of reality, 618
 initial high-density state, 421
 initial low-density state, 411, 420–
 422
 invariant total energy, 570–572, 584
 isolated system, 570
 isotropic expansion, 439, 440
 kinetic energy of expansion, 135,
 352–374, 423, 439, 449
 low-gravitational-entropy Big Crunch,
 548
 matter energy, 290, 332, 354–374,
 423, 439–443, 450, 635
 matter momentum, 354
 negative curvature, 352, 357, 358,
 364
 non-zero curvature of space, 365n
 open, 364–369
 positive curvature, 352, 357, 358,
 364, 372
 radiation energy, 126, 290, 322,
 347, 372, 443
 radius of curvature, 372, 450, 451
 scale factor, 105, 285–299, 327,
 355–357, 367, 401, 440, 634
 static, 411
 thermal energy of matter, 362
 total energy, 353–356, 365n
 unique future, 464–467, 477, 576
 unique past, 464, 465, 477, 519,
 576
 vacuum-dominated era, 288–292
 variation of expansion rate, 368
 wave function, 571, 576, 585, 597,
 598
 zero-angular-momentum condition,
 570n
 zero-energy condition, 88, 298, 352–
 374, 423, 439, 447–451, 570n,
 572n, 635
 zero-momentum condition, 88, 354,
 355, 366n, 570n, 572n
 Unruh effect, 628
 vacuum decay problem, 30, 113, 128–
 133, 179, 633
 vacuum energy, 18, 273–282
 absence of microscopic structure,
 403
 accelerated expansion, 273, 281,
 282, 288–293, 298, 370
 accelerating observer, 496
 average density, *see* cosmological
 constant, vacuum energy
 bi-dimensional universe analogy,
 285, 294
 Casimir effect, 45
 concentration, 329–334
 cosmological constant, 8–11, 100,
 105, 121, 167, 175–181, 228,
 273–282, 287–299, 326–332, 355–
 357, 365–372, 444, 449, 450,
 496, 634–636
 decelerated expansion, 288, 370–
 372
 electrically neutral vacuum, 120,
 153, 331, 637
 equation of state, 372
 equilibrium state, 101, 171
 from quantum fluctuations, 99, 171
 gravitational influence, 278
 gravitational potential energy, 365n,

- 372
- ground state, 113, 132
- homogeneity, 402–406
- initial magnitude, 281
- interactions with matter, 121, 136–138, 145–155, 177–181
- local absence of absence, 111
- local variations, 175, 296, 326–336, 349, 372, 637
- maximum contributions, 171–176, 171n, 181, 282–296, 326–331, 351, 637
- maximum density, 277, 365, 373, 496
- natural zero density, 105, 281, 282, 299, 441, 635, 636
- negative densities, 45, 99, 101, 113, 128, 133, 633
- negative-energy observers, 105, 278–284, 635
- observer-dependent volumes, 284–287, 296, 326, 327, 449
- persistent microscopic structure, 558n
- positive and negative contributions, 99, 105, 106, 170–174, 181, 277–297, 326–331, 637
- quantum-gravitational cut-off, 277
- self-amplifying cosmological constant, 288, 299, 327, 636
- self-reducing cosmological constant, 288–293, 298, 365–369
- smaller past value, 291, 292
- uniform distribution, 398, 496
- variable cosmological constant, 289–292, 299, 324, 355, 367–372, 400–403, 444, 636
- virtual particles, 105, 120, 141, 281, 331, 495–497, 628
- virtual processes, 99, 100, 141
- zero-point fluctuations, 11, 99, 105, 106, 137, 153, 173–176, 273–283, 292–300, 326–331, 342, 496, 633–635
- voids in a matter distribution, 8, 91–111, 119, 120, 153, 154, 326, 336, 346, 394
- Birkhoff’s theorem, 92, 93
- effects on expansion of space, 91–93
- entropy growth, 396, 407
- from explosive processes, 350
- gravitational attraction, 101, 110, 111, 300, 301, 397, 404, 636
- gravitational dynamics, 91, 92, 121, 155
- gravitational field, 403–407
- gravitational forces, 407, 626
- hollow sphere analogy, 93–96
- inertial mass, 91
- largest voids, 341–346
- spherical voids, 92
- surrounding overdense shell, 97, 107, 108
- unexpectedly large gravitational repulsion, 346
- unexpectedly large voids, 344, 345
- vacuum dark matter, 333, 334, 407
- voids in negative vacuum energy
- mutual interactions, 109
- positive-energy matter, 101–109, 120, 121, 126, 136, 153–155, 167–173, 292–297, 331, 351, 403, 496, 637
- voids in vacuum charge distribution
- positive-energy matter, 331, 637

- Von Neumann equation, 554
Von Neumann, John, 486, 512, 569, 607
- wave equations
 absorber theory, 479–485
 advanced solutions, 478–489, 535, 596
 boundary conditions, 486, 487
 constructive interference, 479, 480, 523, 612, 629
 destructive interference, 479, 480, 486, 523, 532, 533, 540, 598, 612, 629
 periodic boundary conditions, 598
 probability wave, 507
 relativistically invariant, 479
 retarded solutions, 478, 479, 485–489, 535
 stationary wave, 598
 wave functions propagating in opposite directions of time, 490
- wavelength, 246
- weak anthropic principle, 105, 298, 299, 367–370, 421, 450, 451, 554, 574, 578n, 580, 635, 641
- Weinberg, Steven, 298
- Wheeler, John, 41, 479, 480
- white hole, 249, 414, 415, 436–438
- wormhole, 146, 147
 exotic matter, 146, 147
 instability, 147
 throat, 146
 traversable, 146, 147
- Zeh, Heinz Dieter, 554
- zero energy level, 131, 168