# A versus B!

## Topological nonseparability and the Aharonov-Bohm effect

Tim Oliver Eynck[*], Holger Lyre[†], Nicolai von Rummell[†]

September 2001

**Abstract.** Since its discovery in 1959 the Aharonov-Bohm effect has continuously been the cause for controversial discussions of various topics in modern physics, e.g. the reality of gauge potentials, topological effects and nonlocalities. In the present paper we juxtapose the two rival interpretations of the Aharonov-Bohm effect. We show that the conception of nonlocality encountered in the Aharonov-Bohm effect is closely related to the nonseparability which is common in quantum mechanics albeit distinct from it due to its topological nature. We propose a third alternative interpretation based on the loop space of holonomies which serves to solve some of the problems and we trace back the topological nonlocality and thereby the Aharonov-Bohm effect to their quantum mechanical origin. All three discussed interpretations are, of course, empirically equivalent. In fact, they present us with an instructive case study for the thesis of theory underdetermination by empirical data.

## 1   Introduction

In 1959 Yakir Aharonov and David Bohm discovered an (unexpected) influence of the electromagnetic gauge potential on the quantum wave function. This effect is called the Aharonov-Bohm (AB) effect. For more than 40 years now, it has been at the focus of a continuing debate on a variety of challenging topics including the reality of gauge potentials, topological effects and nonlocalities in physics.[1]

It is probably fair to say that after years of intense discussion in the physics literature, the first three decades of which are most accessibly summarized and reviewed in the book by Peshkin and Tonomura (1989), the dust has more or less settled and agreement been achieved within the physics community at least as to the mathematical formulation of the AB effect and its experimental verification. Today, it is the broad textbook consensus that the AB effect shows

---

[*]NIKHEF, Postbus 41882, NL-1009 DB Amsterdam, The Netherlands, Email: teynck@nikhef.nl

[†]Institut für Philosophie, Ruhr-Universität Bochum, D-44780 Bochum, Germany,
Emails: holger.lyre@ruhr-uni-bochum.de, nicolai.rummell@ruhr-uni-bochum.de

[1]Aharonov and Bohm themselves paved the way for this debate not only with their 1959 publication, but with a series of papers (Aharonov and Bohm, 1961, 1962, 1963).

the fundamental role played by the gauge potentials in quantum theory as opposed to their auxiliary role in classical physics.[2] Since formulating (quantum) electrodynamics in terms of the electromagnetic field strength would require some sort of action at a distance to explain the AB effect, most physicists happily accepted the observable significance of gauge potentials. However, the gauge freedom of the potentials prevents one from promoting them to the same reality status as that enjoyed by for instance the electromagnetic field in the classical theory.

Despite the aforementioned debate on its physics aspects, philosophical interest in the AB effect, and more generally, in gauge theories, has only arisen fairly recently.[3] Such studies investigate the conceptual framework of the so eminently successful quantum field theories and often consider the AB effect as at least an important case study. To this end, they attempt to state more precisely the concepts and assumptions important already in the physics discussion but sometimes defined only implicitly or even glossed over there. Rather than mimicking a pseudo-mathematical style of first giving all relevant definitions and only then entering into the discussion of the AB effect we prefer to define these concepts as we go along. In this way we hope to motivate deviations from formulations closer to the working physicist's intuition, irrespective of whether stated explicitly elsewhere in the literature or merely implied. At the focus of our attention will be the different concepts of locality and questions pertaining to the ontological status of physical structures.[4] Our terminology will be to use local/locality as generic terms comprising particular types of locality such as separability, local action and point-like interaction.

We believe that the philosophical case of the AB effect is not expecting its final verdict for a long time yet but rather still summoning the witnesses. In section 2 of this paper we will give a brief description of the AB effect and review the, *prima facie*, two rival interpretations it allows for, namely the A-interpretation, in which the gauge potential $A_\mu$ is considered the basic entity, as opposed to the B-interpretation, in which only the magnetic field $\vec{B}$ is used (these correspond to the A and B in the title of this study). Besides an introductory summary of the effect itself and its interpretational history, included to make our statement comprehensible also for the more philosophically-minded reader with less background physics knowledge, we will then introduce a third alternative, the so-called *loop approach* to gauge theories and concentrate on its possible merits for the case at hand in section 3. This C-interpretation is based on gauge invariant quantities only. It allows for a complete, but nevertheless still economic ontology that could probably be accepted also by the working physicist. We explain how the—in this formulation particularly important—property of *nonseparability* stems from the interplay of a non-simply connected base space and the electromagnetic gauge group, and, consequently, refer to it as *topological nonseparability*. Albeit ourselves favouring the third (i.e. C), we aim to present the pros and cons of all three interpretations in a concise form so as to enable the reader to see clearly what is gained by accepting a particular one of them but also which price it comes at. We then argue why the triple ABC constitutes an example for theory underdetermination by

---

[2]The latter being due to the so-called *gauge freedom* of the potentials, i.e. in classical electrodynamics the transformation $A_\mu \to A_\mu + \partial_\mu \Lambda(x)$ with an arbitrary real-valued function $\Lambda(x)$ leaves the physics unaltered.

[3]To considerable extent triggered by Sunny Auyang's book (1995); cf. also Teller (1997), Redhead (1998).

[4]To be sure our study is not the first of this kind. See in particular the intriguing discussion of Healey's 1997 paper by Maudlin (1998), Leeds (1999), and again Healey (1999, 2001).
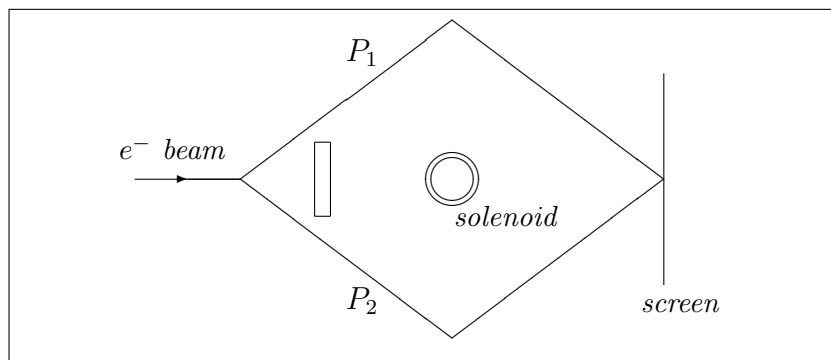
Figure 1: Schematic experimental configuration of the AB effect.

empricial data.

Finally, in section 4, we identify the topological nonseparability, which is at the heart of the AB effect and its nonlocality, as a genuine *quantum structure*. A bit whimsically, we shall call this the "Q-metainterpretation", i.e. it explains the AB effect at the meta level of identifying the theory to which it belongs as a quantum effect. It gives us an important insight into the twofold way quantum theory presents itself as nonlocal: first, in terms of its state space nonseparability, as we call it, responsible for the typical quantum mechanical correlations and, second, the newly introduced and truly quantum phenomenon of topological nonseparability in gauge theories.

## 2   Two *prima facie* interpretations of the AB effect

The AB effect is usually described in the context of the theory of a quantized non-relativistic electron coupled to classical (i.e. non-quantized) electrodynamics. Recall the typical experimental set-up of the AB effect (a schematic illustration is given in figure 1): An electron beam is split into two, which pass around a solenoid on paths $P_1$ and $P_2$, respectively, and are then brought to interference at a screen. Neither does the magnetic field $\vec{B}$ have a non-vanishing component in the region outside the solenoid (the magnetic field lines are confined to the solenoid, which for the sake of simplicity is thought of as infinitely long) nor does the electron penetrate the solenoid (which may be shielded). Thus, the set-up does not allow for any local (in a sense to be made precise below) coupling between the electron and the magnetic field. Nevertheless, a shift in the interference pattern will be observed upon alteration of the $\vec{B}$-field.

Before reviewing its two rivaling interpretations let us reformulate the crucial equations of the AB effect in the compact and most convenient calculus of differential forms. Recall the general form of Stokes' theorem $\oint_{\partial M} \omega = \int_M d\omega$ for some differential form $\omega$. Let $\mathcal{S}$ be the surface bounded by some closed curve $\mathcal{C} = \partial \mathcal{S}$ (paths $P_1$ and $P_2$ in the particular case at hand),

then with $dA = B$ we obtain[5]

$$\oint_{\mathcal{C}} A = \int_{\mathcal{S}} dA = \int_{\mathcal{S}} B. \tag{1}$$

Thus, the quantum phase picked up by the electron is given by

$$\Delta\alpha = \oint_{\mathcal{C}} A \tag{2}$$

and the magnetic flux is

$$\Phi_{mag} = \int_{\mathcal{S}} B. \tag{3}$$

The two interpretations to be summarized in the ensuing two subsections each focus on one side of expression (1). Somewhat reversing both the course of history and the running of the alphabet we present first the B-, and only afterwards the A-interpretation.

## 2.1 The B-interpretation

As was pointed out by DeWitt (1962), the AB effect allows for an understanding in terms of the magnetic field strength alone without making use of the potential. But this comes at a high price. From the standpoint of the B-interpretation the AB effect is a nonlocal effect with the nonlocality originating from a mysterious "nonlocal interaction" of $\psi$ (outside the solenoid) and $\vec{B}$ (inside the solenoid), a truly spooky "action at a distance".

Hence this interpretation of the AB effect is not in accordance with the notion of *local action*, i.e. the concept that influences are always mediated such that spatially separated objects can have no immediate effect on each other. But how is this to be understood? For the clarity of our argument we employ a rather weak notion of local action in the spirit of Healey (1997, p. 24). So if the B-interpretation violates this weak constraint it will automactically violate any stronger notion of local action.

**Local action**
*An interpretation of a theory meets our criterion of* local action *if all causes of an event propagate only via continuous physical processes.*

The locality condition figuring in the derivation of Bell's inequalities, namely that the causes of an event propagate continuously via some physical process not faster than at the speed of light, can be regarded as the relativistic formulation of our definition. In our opinion this notion of local action is rather close to the physicist's intuition of what a theory needs to be called local. To accomodate a given set of measuring results, i.e. explain what they were caused by, one must at least be able to tell a convincing causal story involving only continuous physical processes (but not necessarily observables). Note, however, that the concept of local action is different from that of a *point-like interaction* which mostly implies that *the interacting entities*

---

[5]Unless otherwise stated, we set $c$, $\hbar$, $e = 1$ throughout this paper.

*can be defined at a single point, couple to each other at that very point and are non-zero in overlapping space-time regions.*[6]

Coming back to our formulation of the principle of local action we see that this condition is not met by the B-interpretation of the AB effect. In an explanation of the effect that involves only the field strength, there will be no continuous physical transmission process from the change of the magnetic field inside the solenoid to the shift in the interference pattern, because the boundary conditions of vanishing $\psi$ and magnetic field on the cylinder produce a gap that cannot be bridged. However, a violation of the principle of local action does not immediately imply a violation of causality as can be seen in the case of the AB effect. Even in the relativistic case it is far from clear how a violation of local action would lead to causal paradoxes or superluminal signals. Although these fascinating questions do arise in connection with the AB effect we will not concentrate further on this matter in the present work.

But before concluding that the B-interpretation really does violate the principle of local action one must be able to refute a common objection often raised at this point. Since a quantum wave function is never identically zero, Strocchi and Wightman (1974) conclude that there is always a tail penetrating the solenoid regardless of the shielding employed. So by calculating the AB effect with the boundary condition of vanishing $\psi$ on the cylinder they claim that one is fundamentally mistaken and obtains only approximately the correct observable results—thereby losing the ability to appreciate the important conceptual differences. However, if the AB effect truly arose from a local interaction of $\psi$ and $\vec{B}$ inside the solenoid, this would have the immediate consequence that the clarity of the effect observed on the screen should somehow scale with the quality of the shield. But this has not been observed. It thus seems safe to exclude a local/point-like interaction of $\psi$ and $\vec{B}$ inside the solenoid. Another similar objection put forward is that since the magnetic field lines have to close somewhere the wavefunction $\psi$ can interact locally with the field in some far-out region. However, the wavefunction can be shielded from these remote regions, too. This is usually realized with torroidal magnets (Tonomura et al., 1983, Tonomura, 1998).

## 2.2 The A-interpretation

In this subsection we summarize the usual interpretation of the AB effect as a local interaction between $\psi$ and the gauge potential $A_\mu$ (therefore simply referred to as the A interpretation). In this form the AB effect is often presented in physics textbooks as teaching us that the field strength does not suffice to describe all observable effects in electromagnetism. It is then concluded that the wavefunction of the electron couples locally to the electromagnetic potential which unlike in classical physics no longer figures merely as a convenient calculation device but is to be considered a physically real field. This view has become popular through the Feynman lectures (cf. Feynman, 1963, vol. II, part 1, chap. 15-5). As was soon noticed the problem

---

[6]To be sure this definition is far from being precise: It presupposes the term "point" (or, to be mathematically correct, its infinitesimal neighbourhood) which is part of what should be defined. Moreover, the word "entities" is not intended to suggest any commitment to a particular ontology. Since it is not essential for the purpose of our paper to spell out a rigorous definition of the concept of point-like interaction, this informal description should be sufficient.

with this account is that since the electromagnetic potential is not gauge invariant it cannot without hesitation be considered physically real. On the contrary, in modern field theories gauge invariance is considered the *conditio sine qua non* for a quantity to qualify as an observable and thereby to be physically real. At least this is the working physicist's perspective on that matter. How then does this tie in with definitions introduced above?

The A interpretation is in agreement with the principle of local action. Clearly all causes for the shift in the interference pattern lie in its backward light cone. Moreover, the effect may be explained as the result of an interaction between $\psi$ and $A_\mu$ on the way to the screen in the sense that at least $\psi$ and $A_\mu$ are defined and non-zero in this region. Hence, as far as local action is concerned, the A-interpretation leads 1:0 over the B-interpretation.

But does this really render the AB effect local? The reader should bear in mind that locality implies more than just the idea that all causes propagate via continuous physical processes. Consider two spatially separated objects or systems. In case all observable properties, or all possible information of the compound system, can be associated to either one or the other sub-system, the description of the system is usually called local or separable. Conversely, the description is called nonseparable, if this procedure fails, i.e. information is lost or cannot be distributed between the constituents.[7] In other words, the principle of separability may roughly be formulated as follows: two spatiotemporally distant systems (objects) posses their own independent physical state. Here, we propose a more general formulation:

**Separability**
*Given a physical system S and a partition of its spatiotemporal support R into space-time regions the principle of* separability *states that it is always possible to decompose S, as induced by the partition of R, into subsystems with associated observable properties and to retrieve—to any desired accuracy—the properties of S from the properties of these subsystems.*

For sure there exist nonseparable phenomena already in classical physics in the form of correlations, but these do not pose any serious difficulties and can be explained by appealing to common causes. The most prominent example of nonseparability stems from the tensor product structure of the Hilbert space as the state space of quantum mechanics. We may therefore call it a *state space nonseparability*: The full state vector of a composite quantum mechanical system cannot always be factorized into a product of state vectors of the component systems, but rather is a linear combination thereof. These so called entangled quantum systems do not conform to the above given formulation of separability, since their state vectors represent the complete physical state of a quantum system. We should bear in mind that the principle of separability points to a certain conception of spacetime, whereas the dynamics of a quantum system evolves on a Hilbert space, which itself is not directly related to position space. This in the first place allows for the nonseparability of quantum mechanics. However, the AB effect is nonseparable in a different sense, as we shall argue in the following.

---

[7] From this perspective the notion of separability is closely related to holism. For a recent discussion see Redhead (1987), Healey (1997), Maudlin (1994) and Esfeld (2001).

A closer look at the phase integral (2), which is needed for an adequate description of the phenomenon, reveals that the A-interpretation is nonseparable with respect to our definition above. Although it is of course possible to decompose the path of the integral attributing properties (values of the gauge potential) to the elements of this partition and from this to retrieve all properties of the complex system (including the phase shift crucial to the AB effect), the A-interpretation still fails to meet all the criteria given in our definition. The values of the gauge potential associated to the space regions of the partition fail to be observable because of the gauge freedom still inherent in the potentials. Hence, a partition into *observable* properties fails.

Thus, accepting the A-interpretation does not render the AB effect local.[8] To be more precise, the A-interpretation gives a nonseparable—and, hence, nonlocal—account, since the AB effect is indeed a *topological effect*: its explanation requires knowledge about the underlying space (or spaces) as a whole (compare again footnote 7). The reader may note that *global* topological effects are nonlocal in the sense that the information cannot be distributed between spacetime regions. Hence, the AB effect is an example of what we like to call *topological nonseparability*, since the observable effect of the shift of the interference fringe cannot be reduced to observable properties associated to spacetime regions.

In more formal terms, the phase factor $exp(i\pi BR^2)$ picked up by an electron on its way from source to screen does not depend on the particular path we choose to compute it. It does, however, count the number of times the electron winds around the solenoid. For $n$ windings we obtain $exp(ni\pi BR^2)$. The phase factor is the image of a mapping from the electron's configuration space to the (electromagnetic) gauge group, i.e. $\mathbb{S}^1 \to \mathbb{S}^1$. These mappings fall into equivalence classes under homotopy and constitute the fundamental group of $U(1)$. Were it not for the non-trivial topology of *both* the base space and the gauge group, any two magnetic fields confined to the inside of the solenoid would necessarily have to have the same (null) effect on the interference pattern. Therefore, only the non-trivial topology of both spaces produces the AB effect and its peculiar type of nonlocality is best addressed as a topological nonseparability. It is obvious to ask for more examples falling into this category of nonseparability, or in other words for generalizations of the AB effect. From the above topological considerations it is clear that at least for the closest analogy, i.e. with the same base space, this requires the responsible gauge group to be non-simply connected. This is not the case for those typically considered in the standard model, the fundamental groups of the higher $SU(N)$ Lie groups are all trivial. Of course, other topological effects (solitons, instantons, non-trivial vacua etc.) on higher-dimensional base manifolds are investigated in the physics literature, but the analysis of their relation to our formulations of local action and separability is far beyond the scope of this study.

---

[8] Moreover, Michael Redhead (2002) has pointed out that the A-interpretation leads to an indeterminism at the level of what he calls "surplus structure", namely the additional mathematical structure encoded in the potential $A_\mu$ (i.e. gauge freedom) which is not governed by the field equations. This argument is basically the same as the famous "hole argument" (Earman and Norton, 1987). The B-interpretation also—but in a way more 'naturally'—takes into account the nonseparability of the AB effect as expressed by the phase integral (3). But, as opposed to the A-interpretation, the B-interpretation does not commit us to any surplus structure and, hence, the indeterminism problem can be avoided.

## 2.3   A versus B!

So far the score looks like a draw between A- and B-interpretations. While the A interpretation meets the criteria of local action and violates separability, exactly the opposite can be said about the B interpretation.

| locality type | — interpretation — | |
|---|---|---|
|  | A | B |
| local action | yes | no |
| separability | no | yes |

We have seen that the responsibility for the AB effect may be shifted from the magnetic fields to the gauge potentials and vice versa without any empirically discernible consequences. However, despite leading to the same empirical predictions these two interpretations present us with quite different perspectives on the concept of space or spacetime as a consequence of the different locality properties entailed. How is this possible? We will try to give a partial answer to this question by focusing on the connection and differences between these two types of locality. The first point to notice is that our definition of local action is about processes while separability is concerned with observable properties. Although we do not attempt exhaustive definitions of these two concepts, it seems fair to say that processes are related to the time-evolution of a system, whereas observable properties are measured at a particular instant of time. Thus, local action may be characterized as *diachronous*—in contrast to the *synchronous* notion of separability. The latter describes a feature of possible outcomes of measurements on complex systems, while local action restricts the set of convincing "stories" or explanations we are inclined to accept as describing the causal processes leading to the measured results.

If one finds spatially distant systems not to be independent of each other (i.e. observes some sort of correlation between the measuring results on individual constituent sub-systems), one is led to look for a process together with a "common cause" that links the two systems.[9] But here lies an important difference between classical and quantum physics. In quantum theory we loose our firm grip on what happens between two measurements, i.e. we lack a deterministic causal description of the underlying processes taking place in spacetime. This is why two interpretations of the same quantum effect may imply different underlying space or spacetime concepts. In quantum theory we have to rely on more indirect methods of investigation than in classical mechanics. Within the framework of the latter it would not be possible to switch from an

---

[9]As far as separability is concerned, classical and quantum mechanical correlations are on the same footing. The correlation between the colour of Bertlmann's socks (Bell, 1987, chap. 16) and the direction of polarization of electron spins in EPR type experiments both allow for only nonseparable accounts of the compound systems. But the two correlations are distinct insofar as the supposed processes of the common cause in case of Bertlmann's socks meet the criterion of relativistic locality that figures essentially in the derivation of Bell's inequalities while quantum mechanical correlations do not. Thus, supposing that quantum mechanical correlations have a common cause and demanding that this common cause obey the criterion of relativistic locality leads to disagreement with experimental observations, as expressed in Bell's inequalities and their empirical refutation (Bell, 1987, chaps. 1, 2).

interpretation that violates separability to an empirical equivalent interpretation that violates local action, since it is—at least in principle—possible to obtain knowledge about what happens at each spacetime point and tell a causal story involving observables only.

## 2.4   A case of theory underdetermination

Viewed in this light, the AB effect turns out as a good example for theory underdetermination by empirical data. Indeed, Stokes' formula (1) works as a "hinge" connecting the two interpretations and allowing us to switch back and forth between them. Although the connection between these two interpretations is obvious in terms of a formula ensuring their empirical equivalence and excluding the possibility of an *experimentum crucis*, there are contexts of the B-interpretation (the concept of local action, for instance) which are not translatable into the A-language and vice versa.

One might object that in case of the AB effect it is devious to speak of theory underdetermination by empirical data since the B-interpretation cannot be taken seriously for several reasons. One is that the Schrödinger equation takes a simple and, also, "beautiful" form only with the gauge potential $A_\mu$, whereas somehow coupling to the field strength $F_{\mu\nu}$ would seem "unnatural". Another reason might be that neither quantum mechanics nor classical electrodynamics are valid theories so that there exists the possibility to look beyond the model in which the AB effect is described.[10] At first glance, the success of the gauge-based standard model provides a practical argument for the A-interpretation.

But nevertheless it seems fruitful to us to study even cases of interpretations that can be considered rivaling only without taking into account guiding principles such as beauty or simplicity and without looking beyond the current model. Studying such notions of one interpretation that do not have counterparts in the other[11] or vice versa one gets to see most clearly the crucial theoretical concepts on which the two interpretations disagree. So by finding and analyzing examples of theory underdetermination the attention is drawn to concepts that are not yet fully understood and that point beyond the current model without appealing to simplicity or knowledge from more sophisticated models. This is why we think the AB effect is an example worth noting, because it sheds some light on the concept of locality.[12]

## 3   From AB to the C-interpretation

In section 2 we already mentioned the famous EPR correlations of quantum theory and its underlying state space nonseparability, so-called after the tensor product structure of the Hilbert

---

[10]Belot (1998) has made a similar point.

[11]These concepts are bound to be theoretical terms since the two interpretations are empirically equivalent.

[12]Cf. Lyre and Eynck (2001) for modern examples of theory underdetermination in gravity and a general discussion of the idea of "practical underdetermination of non-final theories" which aims to capture besides its interpretational aspects also its role as motivation for the practising scientist to focus attention and scientific effort.

space employed to describe it. Clearly, this concept of nonseparability is of a genuinely quantum origin. However, the nonseparability we observe in the AB effect has a different origin, namely the existence of a non-vanishing phase integral (2)—the *holonomy*. We will offer in this section a third interpretation of the AB effect suggested by a less well known formulation of electrodynamics. Based on the pioneering work of Yang (1974) and Wu and Yang (1975) it is possible today to represent pure gauge theories in terms of holonomies only. Working with equivalence classes of closed curves given by (2) the so-called *loop approach* (cf. Gambini and Pullin, 1996) is studied by some physicists mainly in the hope that it could provide solutions to problems occuring in the still elusive quantization of gravity. It is not our aim to in any way judge this non-mainstream programme from a physicist's point of view. We merely observe that it goes a step further than the more common A-interpretation in the sense that not only does it accept the holonomies as the most general electromagnetic observables but it also achieves a full description of the dynamics of the theory in terms of these quantities, including their mathematically consistant quantization. It seems therefore appropriate to introduce an alternative, third interpretation of the AB effect, namely the *C-interpretation* which is based on *closed curves*, i.e. holonomies

$$S(C) = exp\left(i \oint_{\mathcal{C}} A\right). \tag{4}$$

We would like to stress that this third option is of conceptual importance even if the loop approach fails to live up to the hopes of those researchers currently studying it. If, in the worst case scenario for those working on it, the loop approach turns out to be calculationally inconvenient and in all cases pragmatically inferior to the standard formulation of gauge theories, it will still provide a viable C-interpretation of the AB-effect.

As we have seen, the A-interpretation fails to give a separable account whereas the B-interpretation is not in accordance with local action. The C-interpretation decidedly leaves us with a nonseparable account of the AB effect (as the A-interpretation does), and, in the same manner, does respect local action. Thus, our total balance now looks like the following:

| | — interpretation — | | |
| --- | --- | --- | --- |
| | A | B | C |
| point-like interaction | yes | no | no |
| local action | yes | no | yes |
| separability | no | yes | no |
| observability | no | yes | yes |

Indeed, the AB effect provides an interesting example of theory underdetermination with respect to the three rivaling interpretations A, B and C. As we already mentioned in the discussion of the A- and B-interpretations in section 3, practising physicists usually favour the A-interpretation due to metatheoretical criteria. Now, as far as our central notions of nonlocality are concerned, the A- and C-interpretations have equal rights. But the A-interpretation is consistent with the concept of point-like interaction, whereas the C-interpretation is not. However, the A-interpretation, also, unavoidably introduces unobservable surplus structure (compare again footnote 8).

By way of contrast the C-interpretation has the particular advantage of meeting a simple but nevertheless sensible, straightforward criterion of physical reality, namely that *a theoretical entity may be considered physically real once any of its alterations has discernible (i.e. observable) effects.* By its very nature, the loop approach deals with gauge-invariant observables only. In this way, the C-interpretation avoids the introduction of mysterious surplus structure—and at the same time respects local action (as the B-interpretation does). Therefore, in our understanding, the C-interpretation clearly has to be favoured. As some of us have argued elsewhere (Drieschner, Eynck, and Lyre, 2002), the loop approach is mathematically equivalent to the restriction of equivalence classes of gauge potentials, which we called "prepotentials". Prepotentials are, indeed, topologically nonseparable entities. Ontologically speaking, they constitute the basic entities in gauge physics (just as holonomies do in the C-interpretation of the AB effect).

## 4 Quantum origin of the AB effect: the Q-metainterpretation

According to the C-interpretation, which in our view provides the proper ontological description of the AB effect, prepotentials—or holonomies—are to be considered the true basic entities in gauge physics. Does this observation allow one to categorize the AB effect as either classical or quantum? One could argue that $U(1)$, the gauge group of QED, already occurs in classical electrodynamics. But although the classical theory certainly admits a transformation of the form $A_\mu \to A_\mu + \partial_\mu \Lambda(x)$ which leaves the physics unaltered, the $\Lambda(x)$ featuring in it could be seen as the generator of either $U(1)$ or $\mathbb{R}$. In other words, given $\mathbb{R}$ as a Lie algebra, the Lie group is underdetermined at the classical level ($G = \mathbb{R}$ or $G = U(1)$ ). It is only in a quantum theory that one could refer to the gauge principle which, involving the transformation of a spinor, certainly singles out the unitary group from the two possible candidates.

However, the gauge principle has recently been questioned by Brown (1999), Healey (2001), Teller (2000) and Lyre (2000, 2001)[13]. These critiques doubt that the usual textbook treatment of the gauge principle suffices to motivate the introduction of the electromagnetic field and its coupling to fermionic matter, or in general any non-abelian gauge field. It is argued in different ways that the gauge principle merely expresses some freedom to choose internal coordinates and afterwards "repair" a non-trivial choice through the introduction of a potential or connection, but does not necessitate the introduction of a real physical field and corresponding finite field strength. In order to make connection with observation some extra empirical input is called for. We do not at all intend to dispute the great practical and phenomenological success of gauge theories, but we would like to stress that if one shares the doubts expressed in the aforementioned references, then one might want to reverse the usual argumentation, in which the AB effect appears as a consequence of a complete formulation of electrodynamics, and instead consider the effect as at least one empirical bridge leading from an internal coordinate transformations to a consistent descripton of electromagnetic phenomena[14].

---

[13]Of these Lyre goes as far as to propose a generalised equivalence principle intended to provide the true empirical input for gauge theories.

[14]In the parlance of Lyre (2000, 2001) the AB effect would then as a kind of *experimentum crucis* proof the equivalence not of the conceptually different charges in electromagnetism and the gauge transformation of a

Brushing these worries aside, one can conclude that the true knowledge about the nature of the gauge group comes from quantum theory. Unitary groups such as $U(1)$ appear in physics at a fundamental level as symmetry groups of quantum Hilbert spaces. This allows us to trace the topological nonseparability back to its truly quantum roots: the AB effect is indeed a quantum effect. Hence, we end with a *Q-metainterpretation* of the AB effect providing us with a new type of "quantum structure": topological nonseparability as a second type of nonseparability—or, more general, nonlocality—in quantum physics.

## 5 Conclusion

In this paper we have presented in a hopefully concise form the possible interpretations of the AB effect. They differ as to which elements of the mathematical formalism are considered an element of physical reality and thus primary, namely A (the electromagnetic potential), B (the field strengths) or C (closed curves, holonomies). We have tried to define and highlight the concepts of locality and separability etc. and attempted to explain in detail which of these feature in a particular interpretation. Although ourselves clearly favouring the C-interpretation, we have intended to act as impartial arbitrators in presenting the score and hope that the reader holding a different point of view on the matter will at least be willing to agree with our catalogue of consequences entailed by subscribing to any one of the three interpretations. We consider this an important endeavour since although the working scientist can often make do with an incomplete explication of basic assumptions and traditional, sometimes deep-rooted prejudices occasionally even promote the advance of physics, the philosopher of science as a spectator *a posteriori* should nevertheless take a critical stand and pay careful attention to these underlying foundations.

This is particularly true in the field of quantum gauge field theory. What Bell (1987, p. 28) once said about the interpretations and implications of quantum mechanics is doubly true of gauge theory and its ramifications: one is mistaken to believe that all questions have long been answered and that the answers could be fully understood in just 20 minutes.

### Acknowledgements

---

spinor, respectively, but of the equivalence of the transformation groups involved. These could *a priori* have been either $\mathbb{R}$ or $U(1)$ for electromagnetism, but only the latter of these two allows for an identification (no matter whether empirically motivated or by means of some generalized equivalence principle) with the term occuring in the spinorial gauge transformation.

# References

Aharonov, Y., and Bohm, D. (1959). Significance of electromagnetic potentials in the quantum theory. *Physical Review*, *115*(3), 485-491.

Aharonov, Y., and Bohm, D. (1961). Further considerations on electromagnetic potentials in the quantum theory. *Physical Review*, *123*(4), 1511-1524.

Aharonov, Y., and Bohm, D. (1962). Remarks on the possibility of quantum electrodynamics without potentials. *Physical Review*, *125*(6), 2192-2193.

Aharonov, Y., and Bohm, D. (1963). Further considerations of the role of electromagnetic potentials in the quantum theory. *Physical Review*, *130*(4), 1625-1632.

Auyang, S. Y. (1995). *How is quantum field theory possible?* New York: Oxford University Press.

Bell, J. S. (1987). *Speakable and unspeakable in quantum mechanics.* Cambridge: Cambridge University Press.

Belot, G. (1998). Understanding electromagnetism. *The British Journal for the Philosophy of Science*, *49*, 531-555.

Brown, H. R. (1999). Aspects of objectivity in quantum mechanics. In Butterfield, J., and Pagonis, C. (Eds.). *From physics to philosophy.* Cambridge: Cambridge University Press.

DeWitt, B. S. (1962). Quantum theory without electromagnetic potentials. *Physical Review*, *125*(6), 2189-2191.

Drieschner, M., Eynck, T. O., and Lyre, H. (2002). Comment on Redhead: The interpretation of gauge symmetry. In Kuhlmann, Lyre, and Wayne.

Earman, J., and Norton, J. (1987). What price spacetime substantivalism? The hole story. *The British Journal for the Philosophy of Science*, *83*, 515-525.

Esfeld, M. (2001). *Holism in philosophy of mind and philosophy of physics.* Dordrecht: Kluwer.

Feynman, R. P. (1963). *The Feynman lectures on physics.* Reading, Mass.: Addison-Wesley.

Gambini, R., and Pullin, J. (1996). *Loops, knots, gauge theories and quantum gravity.* Cambridge: Cambridge University Press.

Healey, R. (1997). Nonlocality and the Aharonov-Bohm effect. *Philosophy of Science*, *64*, 18-41.

Healey, R. (1999). Quantum analogies: A reply to Maudlin. *Philosophy of Science*, *66*, 440-447.

Healey, R. (2001). *On the reality of gauge potentials.* (E-print PITT-PHIL-SCI00000328)

Kuhlmann, M., Lyre, H., and Wayne, A. (Eds.). (2002). *Ontological Aspects of Quantum Field Theory.*, Singapore: World Scientific. (In preparation)

Leeds, S. (1999). Gauges: Aharonov, Bohm, Yang, Healey. *Philosophy of Science*, *66*, 606-627.

Lyre, H. (2000). A generalized equivalence principle. *International Journal of Modern Physics D*, *9*(6), 633-647. (E-print arXiv:gr-qc/0004054)

Lyre, H. (2001). The principles of gauging. *Philosophy of Science*, *68*(3, Supplement). (E-prints arXiv:quant-ph/0101047, PITT-PHIL-SCI00000113)

Lyre, H., and Eynck, T. O. (2001). *How to gauge gravity? Underdetermination in gravitational theories.* (Preprint.)

Maudlin, T. (1994). *Quantum non-locality and relativity.* Oxford: Blackwell.

Maudlin, T. (1998). Discussion: Healey on the Aharonov-Bohm effect. *Philosophy of Science*, *65*, 361-368.

Peshkin, M. A., and Tonomura, A. (1989). *The Aharonov-Bohm effect.* Berlin: Springer.

Redhead, M. (1987). *Incompleteness, nonlocality, and realism.* Oxford: Clarendon Press.

Redhead, M. (1998). Review: S. Y. Auyang, How is quantum field theory possible? (New York, Oxford University Press, 1995). *The British Journal for the Philosophy of Science*, *49*, 499-507.

Redhead, M. (2002). The interpretation of gauge symmetry. In Kuhlmann, Lyre, and Wayne.

Strocchi, F., and Wightman, A. S. (1974). Proof of the charge superselection rule in local relativistic quantum field theory. *Journal of Mathematical Physics*, *15*, 2198-2224.

Teller, P. (1997). A metaphysics for contemporary field theories–Essay review: S. Y. Auyang, How is quantum field theory possible? (New York, Oxford University Press, 1995). *Studies in History and Philosophy of Modern Physics*, *28*, 507-522.

Teller, P. (2000). The gauge argument. *Philosophy of Science*, *67*(3, Supplement), S466-S481.

Tonomura, A. (1998). *The quantum world unveiled by electron waves.* Singapore: World Scientific.

Tonomura, A., Umezaki, H., Matsuda, T., Osakabe, N., Endo, J., and Sugita, Y. (1983). Is magnetic flux quantized in a toroidal ferromagnet? *Physical Review Letters*, *51*, 331-334.

Wu, T. T., and Yang, C. N. (1975). Concept of nonintegrable phase factors and global formulation of gauge fields. *Physical Review D*, *12*(12), 3845-3857.

Yang, C. N. (1974). Integral formalism for gauge fields. *Physical Review Letters*, *33*(7), 445-447.