1
2 **Assessing climate model projections: state of the art and philosophical reflections**
5
6
7 Joel Katzav
8 The Department of Philosophy and Ethics, Eindhoven University of Technology, the
9 Netherlands
10
11 Henk A. Dijkstra
12 Institute for Marine and Atmospheric Research, Utrecht University, the Netherlands
13
14 A. T. J. (Jos) de Laat
15 The Royal Netherlands Meteorological Institute, the Netherlands
16
17
18
19
20
21
22
23
24
25
26
27
28
29
30
31
32
33
34
35
36
37
38
39
40
41
42
43
44
45
46
47
48
49
50

51

52                                    **Abstract**

53    The present paper draws on climate science and the philosophy of science in order to

54    evaluate climate-model-based approaches to assessing climate projections. We

55    analyze the difficulties that arise in such assessment and outline criteria of adequacy

56    for approaches to it. In addition, we offer a critical overview of the approaches used in

57    the IPCC working group one fourth report, including the confidence building,

58    Bayesian and likelihood approaches. Finally, we consider approaches that do not

59    feature in the IPCC reports, including three approaches drawn from the philosophy of

60    science. We find that all available approaches face substantial challenges, with IPCC

61    approaches having as a primary source of difficulty their goal of providing

62    probabilistic assessments.

63

64

65

66

67

68

69

70
71
72
73
74
75
76
77
78
79
80
81
82

83
84 **1. Introduction**

85 The climate system is the system of processes that underlie the behavior of

86 atmospheric, oceanic and cryospheric phenomena such as atmospheric temperature,

87 precipitation, sea-ice extent and ocean salinity.  Climate models are designed to

88 simulate the seasonal and longer term behavior of the climate system. They are

89 mathematical, computer implemented representations that comprise two kinds of

90 elements. They comprise basic physical theory – e.g., conservation principles such as

91 conservation of momentum and heat – that is used explicitly to describe the evolution

92 of some physical quantities – e.g., temperature, wind velocity and properties of water

93 vapor. Climate models also comprise parameterizations. Parameterizations are

94 substitutes for explicit representations of physical processes, substitutes that are used

95 where lack of knowledge and/or limitations in computational resources make explicit

96 representation impossible. Individual cloud formation, for example, typically occurs

97 on a scale that is much smaller than global climate model (GCM) resolution and thus

98 cannot be explicitly resolved. Instead, parameterizations capturing assumed

99 relationships between model grid-average quantities and cloud properties are used.

100     The basic theory of a climate model can be formulated using equations for the

101 time derivatives of the model's state vector variables, $x_i$, $i = 1, ..., n,$ as is

102 schematically represented by

103 $$\frac{\partial x_i}{\partial t} = F_i(x_1...x_n, y_1,..., y_n, t) + G_i(t) \tag{1}$$

104     In Eqt. (1), $t$ denotes time, the functions $G_i$ represent external forcing factors

105 and how these function together to change the state vector quantities, and the $F_i$

106 represent the many physical, chemical and biological factors in the climate system and

107 how these function together to change the state vector quantities. External forcing

108    factors – e.g., greenhouse gas concentrations, solar irradiance strength, anthropogenic

109    aerosol concentrations and volcanic aerosol optical depth – are factors that might

110    affect the climate system but that are, or are treated as being, external to this system.

111          The $x_i$ represent those quantities the evolution of which is explicitly described

112    by basic theory, that is the evolution of which is captured by partial time derivatives.

113    The $y_i$ represent quantities that are not explicitly described by basic theory. So these

114    variables must be treated as functions of the $x_i$, i.e., the $y_i$ must be parameterized. In

115    this case, the parameterizations are schematically represented in Eqt. (2).

116    $$y_i = H_i(x_1,...,x_n) \qquad\qquad (2)$$

117    Given initial conditions $x_i(t_0)$ at time $t = t_0$ and boundary conditions, the climate

118    model calculates values of the state vector at a later time $t = t_1$ in accordance with

119    Eqt. (1).

120          Climate models play an essential role in identifying the causes of climate

121    change and in generating projections. Projections are conditional predictions of

122    climatic quantities. Each projection tells us how one or more such quantities would

123    evolve were external forcing to be at certain levels in the future. Some approaches to

124    assessing projections derive projections, and assess their quality, at least partly

125    independently of climate models. They might, for example, use observations to decide

126    how to extend simulations of present climate into the future (Stott et al., 2006) or

127    derive projections from, and assess them on the basis of, observations (Bentley, 2010;

128    Siddall et al., 2010). We focus on climate-model-based assessment. Such assessment

129    is of the projections of one or more climate models and is assessment in which how

130    good models are in some respect or another is used to determine projection quality. A

131    climate model projection (CMP) quality is a qualitative or quantitative measure, such

132    as a probability, that is indicative of what we should suppose about CMP accuracy.

133       It is well recognized within the climate science community that climate-

134    model-based assessment of projection quality needs to take into account the effects of

135    climate model limitations on projection accuracy (Randall et al., 2007; Smith, 2006;

136    Stainforth et al., 2007a). Following Smith (2006) and Stainforth (2007a), we

137    distinguish between the following main types of climate model limitations:

138        (a) External forcing inaccuracy – inaccuracy in a model's representation of
139            external forcing, that is in the $G_i$ in Eqt. (1).
140

141        (b) Initial condition inaccuracy – inaccuracy in the data used to initialize
142            climate model simulations, that is in the $x_i(t_0)$.
143

144        (c) Model imperfection – limitations in a model's representation of the climate
145            system or in our knowledge of how to construct this representation,
146            including:
147

148            1. Model parameterization limitations – limitations in our knowledge of
149            what the optimal or the appropriate parameter values and parameterization
150            schemes for a model are. This amounts, in the special case where
151            parameterizations are captured by Eqt. (2), to limitations in our knowledge
152            of which functions $H_i$ one should include from among available
153            alternatives.
154

155            2. Structural inadequacy – inaccuracy in how a model represents the
156            climate system which cannot be compensated for by resetting model
157            parameters or replacing model parameterizations with other available
158            parameterization schemes. Structural inaccuracy in Eqt. (1) is manifested
159            in an insufficient number of variables $x_i$ and $y_i$ as well as in the need for
160            new functions of these variables.
161

162    Parameterization limitations are illustrated by the enduring uncertainty about climate

163    sensitivity and associated model parameters and parameterization schemes. A

164    relatively recent review of climate sensitivity estimates underscores the limited ability

165    to determine its upper bound as well as the persistent difficulty in narrowing its likely

166    range beyond 2 to 4.5 °C (Knutti and Hegerl, 2008). The 21 GCMs used by Working

167    Group One of the IPCC fourth report (WG1 AR4) illustrate structural inadequacy.

168    These sophisticated models are the models of the World Climate Research

169    Programme's Coupled Model Intercomparison Project phase 3 (CMIP3) (Meehl et al.,

170  2007a). Some important sub-grid and larger than grid phenomena that are relevant to

171  the evolution of the climate system are not accurately represented by these models,

172  some are only represented by a few of the models and some are not represented at all.

173  Parameterization of cloud formation, for example, is such that even the best available

174  parameterizations suffer from substantial limitations (Randall et al., 2003). None of

175  the models represent the carbon cycle, only some represent the indirect aerosol effect

176  and only two represent stratospheric chemistry (CMIP3, 2007). The models also omit

177  many of the important effects of land use change (Mahmood et al., 2010; Pielke,

178  2005). Many of their limitations, e.g., the limited ability to represent surface heat

179  fluxes as well as sea ice distribution and seasonal changes, are the result of a

180  combination of structural inadequacy and parameterization limitations (Randall et al.,

181  2007, p. 616). CMIP3 simulations illustrate initial condition inaccuracy. Due to

182  constraints of computational power and to limited observations, these simulations start

183  from selected points of control integrations rather than from actual observations of

184  historical climate (Hurrell et al., 2009).

185      The most ambitious assessments of projection quality, and these are primarily

186  climate-model-based assessments, are those of WG1. The first three WG1 reports rely

187  primarily on the climate-model-based approach that we will call the confidence

188  building approach. This is an informal approach that aims to establish confidence in

189  models, and thereby in their projections, by appealing to models' physical basis and

190  success at representing observed and past climate. In the first two reports, however,

191  no uniform view about what confidence in models teaches about CMP quality is

192  adopted (IPCC 1990; IPCC 1996). The summary for policymakers in the WG1

193  contribution to the IPCC first assessment report, for example, qualifies projections

194  using diverse phrases such as 'we predict that', 'confidence is low that' and 'it is likely

195    that' (IPCC 1990). A more systematic view is found in WG1's contribution to the

196    third IPCC assessment report (WG1 TAR). It made use of a guidance note to authors

197    which recommends that main results be qualified by degrees of confidence that are

198    calibrated to probability ranges (Moss and Schneider, 2000). The summary for

199    policymakers provided by WG1 TAR does assign projections such degrees of

200    confidence. It expresses degrees of confidence as degrees of likelihood and takes, e.g.,

201    'very likely' to mean having a chance between 90 and 99 %, and 'likely' to mean

202    having a chance between 66 % and 90 %. The chapter on projections of future climate

203    change, however, defines degrees of confidence in terms of agreement between

204    models. A very likely projection, for example, is defined (roughly) as one that is

205    physically plausible and is agreed upon by all models used (IPCC 2001).

206         WG1 AR4's assessment of projection quality has two stages. First, confidence

207    in models is established as in previous reports. This is mostly achieved in Chapter 8 –

208    which describes, among other things, successful simulations of natural variability

209    (Randall et al., 2007) – and in chapter 9 – which focuses on identifying the causes of

210    climate change, but also characterizes model successes at simulating $20^{th}$ century

211    climate change (Hegerl et al., 2007). The second stage is carried out in Chapter 10 –

212    which provides WG1 AR4's global projections (Meehl et al., 2007b) – and Chapter 11

213    – which focuses on regional projections (Christensen et al., 2007). In these chapters,

214    expert judgment is used to assign qualities to projections given established confidence

215    in models and the results of formal, probabilistic projection assessment (Meehl et al.,

216    2007b). WG1 AR4 is the first WG1 report that makes extensive use of formal

217    assessment, though it recognizes that such approaches are in their infancy

218    (Christensen et al., 2007; Randall et al., 2007). Both climate-model-based and partly

219    climate-model-independent formal approaches are used.

220       Although WG1 AR4 assesses models using degrees of confidence, it does not

221    assess projections in these terms. Nor does it equate projection likelihoods with

222    degrees of agreement among models. It does, however, implement the advice to

223    provide probabilistically calibrated likelihoods of projections (IPCC 2005). For

224    example, unlike WG1 TAR, WG1 AR4 provides explicit likelihood estimates for

225    projected ranges of global mean surface temperature (GMST) changes. It estimates

226    that the increase in GMST by the end of the century is likely to fall within -40 to +60

227    % of the average GCM warming simulated for each emission scenario and provides

228    broader uncertainty margins than the GCM ensemble in particular because GCMs do

229    not capture uncertainty in the carbon cycle (Fig. 2).

230       The sophistication of WG1 AR4's assessments was enabled by the increasing

231    ability to use multi-GCM and perturbed physics GCM ensembles. Thus, while WG1's

232    first two reports relied on simple models to produce long term GMST projections,

233    WG1 TAR and WG1 AR4 relied primarily on state-of-the-art GCM ensembles to

234    assess these and other projections. WG1 AR4 nevertheless still relied on simpler

235    models, including intermediate complexity and energy balance models (Randall et al.,

236    2007).

237       In this review, we provide a critical discussion of the (climate-model-based)

238    approaches to assessing projection quality relied on in WG1 AR4 and more recent

239    work by climate scientists. In doing so, we build on the substantial climate science

240    literature, including WG1 AR4 itself. We, however, extend this literature using the

241    perspective of the philosophy of science. Our discussion does focus more than climate

242    scientists themselves tend to on precisely why assessing projection quality is difficult,

243    on what is required of an adequate approach to such assessment and on the limitations

244    of existing approaches. We, nevertheless, also address some of the practical concerns

245    of climate scientists. We outline three views of how to assess scientific claims that are

246    drawn from the philosophy of science and consider how they might further assist in

247    assessing projection quality. Important issues that space does not allow us to address

248    are the special difficulties that assessment of regional projection quality raises. An

249    issue that deserves more attention than we have given it is that of how uncertainty

250    about data complicates assessing projection quality.

251        We begin (Section 2) by considering what kinds of qualities should be

252    assigned to projections, especially whether probabilistic qualities should be assigned.

253    We then (Section 3) discuss why assessing projection quality is difficult and outline

254    criteria for adequate approaches to doing so. Using these criteria, we proceed to

255    discuss (Sections 4–7) the approaches that were used in WG1 AR4, namely the

256    confidence building, the subjective Bayesian and the likelihood approaches. Finally

257    (Section 8), we discuss approaches that are not used, or are not prominent in, WG1

258    AR4, including the possibilist and three philosophy-of-science-based approaches.

259

260    **2. Probabilistic and non-probabilistic assessment**

261    Probabilistic assessment of projection quality will here be taken to include assigning

262    probabilities or informative probability ranges to projections or projection ranges.

263    Such assessment has been argued for on the ground that it is better suited to handling

264    the inevitable uncertainty about projections than deterministic assessments are

265    (Raisanen and Palmer, 2001). But philosophers of science, computer scientists and

266    others point out that probabilities fail to represent uncertainty when ignorance is deep

267    enough (Halpern, 2003; Norton, 2011). Assigning a probability to a prediction

268    involves, given standard probability frameworks, specifying the space of possible

269    outcomes as well as the chances that the predicted outcomes will obtain. These,

270    however, are things we may well be uncertain about given sufficient ignorance. For

271    example, we might be trying to assess the probability that a die will land on '6' when

272    our information about the kind and bias of the die is limited. We might have the

273    information that it can exhibit the numerals '1', '6' and '8' as well as the symbol '*', but

274    not have any information about what other symbols might be exhibited or, beyond the

275    information that '6' has a greater chance of occurring than the other known symbols,

276    the chances of symbols being exhibited. The die need not be a six sided die. In such

277    circumstances, it appears that assigning a probability to the outcome '6' will

278    misrepresent our uncertainty.

279         Assigning probability ranges and probabilities to ranges can face the same

280    difficulties as assigning probabilities to single predictions. In the above example,

281    uncertainty about the space of possibilities is such that it would be inappropriate to

282    assign the outcome '6' a range that is more informative than the unhelpful 'somewhere

283    between 0 and 1'. The same is true about assigning the range of outcomes '1', '6' and

284    '8' a probability.

285         One might suggest that, at least when the possible states of a system are

286    known, we should apply the principle of indifference. According to this principle,

287    where knowledge does not suffice to decide between possibilities in an outcome

288    space, they should be assigned equal probabilities. Some work in climate science

289    acknowledges that this principle is problematic, but suggests that it can be applied

290    with suitable caution (Frame et al., 2005). Most philosophers argue that the principle

291    should be rejected (Strevens, 2006a). We cannot know that the principle of

292    indifference will yield reliable predictions when properly applied (North, 2010). If,

293    for example, we aim to represent complete ignorance of what value climate sensitivity

294    has within the range 2 to 4.5 °C, it is natural to assign equal probabilities to values in

295   this range. Yet whether doing so is reliable across scenarios in which greenhouse

296   gasses double depends on what climate sensitivity actually tends to be across such

297   scenarios and it is knowledge of this tendency that is, given the assumed ignorance,

298   lacking. Further, we can only define a probability distribution given a description of

299   an outcome space and there is no non-arbitrary way of describing such a space under

300   ignorance (Norton, 2008; Strevens, 2006a). What probability should we assign to

301   climate sensitivity's being between 2 and 4 °C, given complete ignorance within the

302   range 2 to 6 °C? 50 % is the answer, when the outcome space is taken to be the given

303   climate sensitivity range and outcomes are treated as equiprobable. But other answers

304   are correct if alternative outcome spaces are selected, say if the outcome space is

305   taken to be a function not just of climate sensitivity but also of feedbacks upon which

306   climate sensitivity depends. And in the supposed state of ignorance about climate

307   sensitivity, we will not have a principled way of selecting a single outcome space.

308        Although the case of the die is artificial, our knowledge in it does share some

309   features with our knowledge of the climate system. We are, for example, uncertain

310   about what possible states the climate system might exhibit, as already stated in the

311   case of climate sensitivity. A central question in what follows is to what extent our

312   ignorance of the climate system is such that probabilistic assessment of projection

313   quality is inappropriate.

314        Acknowledging that probabilistic assessment is inappropriate in some case is

315   by no means then to give up on assessment. Assigning non-probabilistic qualities can

316   commit us to less than assigning probabilities or probability ranges and thus can better

317   represent uncertainty. Judging that it is a real possibility that climate sensitivity is 2

318   °C does not require taking a position on the full range of climate sensitivity. Nor need

319   rankings of climate sensitivities according to plausibility do so. Other non-

320  probabilistic qualities the assignment of which is less demanding than that of

321  probabilities or probability ranges are sets of probability ranges and the degree to

322  which claims have withstood severe tests (see Halpern (2003) for a discussion, and

323  formal treatment, of a variety of non-probabilistic qualities. We discuss severe-test-

324  based and real-possibility-based assessments in sections 8.4 and 8.1 respectively).

325

326  **3. Why is assessing projection quality difficult?**

327  Projections, recall, are predictions that are conditional on assumptions about external

328  forcing. So errors in assumptions about external forcing are not relevant to assessing

329  projection quality. Such assessment need only take into account the effects of initial

330  condition inaccuracy and model imperfection. In the present section, we consider why

331  these kinds of limitations make assessing projection quality difficult. This question is

332  not answered just by noting that climate models have limitations. Scientific models

333  are in general limited, but it is not generally true that assessing their predictions is a

334  serious problem. Consider standard Newtonian models of the Earth-Sun system. Such

335  models suffer from structural inadequacy. They represent the Earth and the Sun as

336  point masses. Moreover, they tell us that the Earth and the Sun exert gravitational

337  forces on each other, something that general relativity assures us is not strictly true.

338  Still, assessing to what extent we can trust the predictions these models are used to

339  generate is something we typically know how to do.

340

341  **3.1 Initial condition inaccuracy and its impact on assessing projections**

342  We begin by considering the difficulties associated with initial condition error. Work

343  in climate science emphasizes the highly nonlinear nature of the climate system (Le

344  Treut et al., 2007; Rial et al., 2004), a nature that is reflected in the typically nonlinear

345      form of the $F_i$ in Eqt. (1). Nonlinear systems are systems in which slight changes to

346      initial conditions can give rise to non-proportional changes of quantities over time

347      (Lorenz, 1963). This high sensitivity can make accurate prediction inherently difficult.

348      Any errors in simulations of highly nonlinear systems, including even minor errors in

349      initial condition settings, might be multiplied over time quickly. The high sensitivity

350      to initial conditions also, as climate scientists note, threatens to make assessing

351      prediction quality difficult. The way in which error grows over time in such systems

352      cannot be assumed to be linear and might depend on how the system itself develops

353      (Palmer, 2000; Palmer et al., 2005).

354      However, how serious a problem sensitivity to initial conditions is for

355      assessing projection quality is not a straightforward matter. The known inaccuracy in

356      model initial condition settings means that high sensitivity of the evolution of climatic

357      quantities to initial conditions might be important. Yet, a climatic quantity the

358      evolution of which is going to be highly nonlinear at one temporal scale may continue

359      to exhibit approximately linear evolution on another such scale. Greenland ice volume

360      may, for example, evolve linearly in time over the coming few decades but

361      nonlinearly over more than three centuries (Lenton et al., 2008). If this is so,

362      nonlinearity will only be a limited obstacle to assessing projections of Greenland ice

363      volume. More generally, whether, and to what extent, a climatic process is nonlinear

364      will depend on the desired projection accuracy, the quantity of interest, the actual

365      period and region of interest and the temporal and spatial scale of interest (IPCC

366      2001). Thus, whether the highly nonlinear behavior of the climate system is a problem

367      for assessing projection quality will have to be determined on a case by case basis.

368

369      **3.2 Tuning and its impact on assessing projections**

370     Further features of climate modeling complicate determining the impact of model

371     imperfection on CMP quality. The first of these features is tuning. Tuning is the

372     modification of parameterization scheme parameters so as to accommodate – create

373     agreement with – old data. A prominent instance is the setting of parameters

374     associated with the small-scale mixing processes in the ocean. Tuning to current day

375     conditions is hard to avoid given the limited available data about the climate system.

376     Moreover, climate scientists worry that when model success results from

377     accommodation, it provides less confirmation of model abilities than success that

378     results from out-of-sample prediction, that is from prediction that is made prior to the

379     availability of the data but that nevertheless accurately captures the data (Knutti,

380     2008; Smith, 2006; Stainforth et al., 2007a). Prominently, there is the suspicion that

381     accommodation threatens to guarantee success irrespective of whether models

382     correctly capture those underlying processes within the climate system that are

383     relevant to its long term evolution (Schwartz et al., 2007). This impacts assessing

384     projection quality. Difficulty in assessing the extent to which a model's basic

385     assumptions hold will give rise to difficulty in assessing its projections.

386        Work in the philosophy of science, however, shows that whether, and under

387     what conditions, the accommodation of data provides reduced confirmation is an

388     unresolved one (Barrett and Stanford, 2006). On the one hand, some philosophers do

389     worry that accommodation raises the threat of generating empirical success

390     irrespective of whether one's theoretical assumptions are correct (Worrall, 2010). On

391     the other hand, if we prioritize out-of-sample prediction over accommodation,

392     evidence might be good evidence of the suitability of model *A* for generating a set of

393     projections *R* for the late 21$^{st}$ century and not so good evidence for the suitability of

394     model *B* for this purpose even though the models are intrinsically identical. This

395 might occur because the developers of model *B* happen to learn, while those of *A* do

396 not learn, of relevant evidence at the stage of model development. In such

397 circumstances, the developers of *B* might end up accommodating the evidence while

398 the developers of *A* successfully predict it. Resulting differing degrees of confidence

399 in the models would, paradoxically, have to be maintained even if it were recognized

400 that the models are intrinsically identical. If accommodated evidence as such is poor

401 evidence, what determines whether evidence is good evidence for a model is the

402 model's history and not just its intrinsic characteristics (see, e.g., Hudson (2007) for

403 worries about the value of out-of-sample prediction).

404     Unfortunately, while the philosophy of science literature tells us that tuning

405 might not be so bad, it still leaves open the possibility that it is problematic. So how

406 tuning affects CMP accuracy still needs to be addressed.

407     Of course, different approaches to parameterization affect CMP quality

408 differently. For example, stochastic parameterizations, i.e., parameterizations that

409 introduce small but random variations in certain model parameters or variables, are

410 arguably sometimes better than standard deterministic parameterizations (Palmer et

411 al., 2005). The worries about tuning, however, arise for all available parameterization

412 techniques.

413

414 **3.3 The long term nature of projections and its impact on assessing projections**

415 A second factor that, according to some climate scientists, complicates determining

416 the impact of model imperfection is the fact that climate models cannot be tested

417 repeatedly across relevant temporal domains (Frame et al., 2007; Knutti, 2008). We

418 can repeatedly compare weather model forecasts with observations. Success

419 frequencies can then be used to provide probabilistic estimates of model fitness for the

420  purpose of generating accurate forecasts. Recently, some old CMPs have been directly

421  assessed (Hargreaves, 2010). But many CMPs have fulfillment conditions that are

422  never realized and, anyway, CMPs are generally too long term to allow repeated

423  direct testing. Thus, it has been argued, it is hard to take the impact of many model

424  implemented assumptions about long term climate into account in assessing model

425  suitability for generating projections.

426  But the fact that we cannot test our models' predictions over the time scales of

427  the predictions is not itself a difficulty. Consider predictions of Earth orbit variation

428  induced changes in solar radiation at the top of atmosphere over the next million

429  years. Here, predictions are generated using model implemented theory about orbital

430  physics, including Newtonian mechanics and an understanding of its limitations

431  (Laskar et al., 2004). This theory is what grounds confidence in the predictions,

432  though the theory and the models based upon it are only tested against relatively

433  short-term data. As the general views we will discuss about how scientific claims are

434  assessed illustrate, there is no need to assume that estimates of a model's ability must

435  be, or are, made on the basis of numerous observations of how well the model has

436  done in the past.

437

438  **3.4 Basic theory, recognized model imperfection and assessing projections**

439  There are nevertheless two more factors other than tuning that complicate taking into

440  account the effects of model imperfection in assessing projection quality. The first,

441  which is not explicitly discussed in the climate science literature but which climate

442  scientists no doubt recognize, is the combination of known model imperfection with

443  the fact that the background knowledge used in constructing models provides a

444  limited constraint on model construction.

445     Philosophers of science observe that theory provides essential information

446     about model reliability (Humphreys, 2004). Newtonian physics, general relativity and

447     other theories provide essential information about when, and to what extent, we can

448     neglect aspects of the solar system in applying Newtonian theory to model the orbit of

449     the Earth. The same, we have noted, is true of models of how changes in the Earth's

450     orbit affect top of the atmosphere solar radiation. In the case of climate modeling,

451     however, the extent to which theory can guide climate model construction and

452     projection quality assessment is limited. After all, parameterization is introduced

453     precisely because of a limited ability to apply explicit theory in model construction.

454     We do not, for example, have a quantitative theory of the main mechanisms of

455     the stratospheric circulation. As a result, while our partial understanding of these

456     mechanisms can be used in arguing that  CMIP3 GCMs' limited ability to represent

457     the stratosphere adversely affects their simulations of tropospheric climate change, the

458     way and extent to which it does so will remain a matter of ongoing investigation (as

459     in, e.g., Dall' Amico (2010)).

460     A limited ability to apply theory in model construction will even make it

461     difficult to decide what we can learn about CMP accuracy from whatever success

462     models have. For easy, relatively theory neutral, ways of drawing conclusions from

463     model successes are hard to come by given model imperfection.

464     Model imperfection implies that models will only have limited empirical

465     success, as indeed is found in the case of climate models. The strongest claim reported

466     by WG1 AR4 on behalf of simulated GCM multi-model annual mean surface

467     temperatures is that, outside of data poor regions such as the polar regions, simulated

468     temperatures were usually within 2 °C of observed temperatures. For most latitudes,

469     the error in simulated zonally averaged outgoing shortwave radiation was about 6%.

470 Simulation of the strength of the Atlantic Meridional Overturning Circulation (MOC)

471 suffers from substantial inaccuracies (Fig. 3). And the same is true of simulation of

472 precipitation patterns, especially on regional scales (Randall et al., 2007). Such

473 inaccuracies short-circuit a simple argument for assigning a high quality to CMPs,

474 namely one that assigns them such a quality on the ground that they were generated

475 by models which simulate data well across the board. Indeed, there is reason to think

476 that increased ability to simulate the current mean climate state across large sets of

477 climate variables is a limited constraint on CMP accuracy (Abe et al., 2009; Knutti et

478 al., 2010). For example, it has been shown (Knutti et al., 2010) that the range of

479 CMPs of precipitation trends is not substantially affected by whether it is produced by

480 all the CMIP3 models or by a subset of high performing models. Assessment of a

481 projection's quality requires correctly identifying which, if any, aspects of model

482 performance are relevant to the projection's accuracy.

483       Further difficulty in figuring out what to infer from what model success there

484 is arises from the well recognized interdependency of climatic processes. Changes in

485 some climatic processes inevitably give rise to changes in others. Changes in cloud

486 cover, land usage, soil hydrology, boundary layer structure and aerosols will, for

487 example, affect surface temperature trends and vice versa. Thus, an accurate

488 simulation of some quantity $x$ will require an appropriate simulation of related

489 quantities upon which $x$ depends. And our assessment of the quality of a projection of

490 $x$ will have to take into account both the accuracy with which $x$ has been simulated

491 and the accuracy with which related quantities have been simulated. One cannot

492 simply argue that since some models simulate a certain climatic quantity well, their

493 projections of this quantity are good (Parker, 2009).

494    Easy, relatively theory neutral ways of assessing what to infer from limited

495    model successes might also be hampered by structural instability, which is, like high

496    sensitivity to changes in initial conditions, a feature of nonlinear systems. A system is

497    structurally unstable when slight changes to its underlying dynamics would give rise

498    to qualitatively different system evolutions. Components of the climate system do

499    exhibit structural instability (Ghil et al., 2008; McWilliams, 2007). This means that

500    minor observed errors in simulating current climate might, given model imperfection,

501    lead to substantial errors in CMPs.

502

**3.5 Unrecognized model imperfection and assessing projections**

503

504    The final source of difficulty for assessing projection quality in light of model

505    imperfection is the possibility, worried about by scientists from all fields, that our

506    models are wrong in unrecognized ways. Empirically successful theories and models

507    have often turned out to rest on mistaken assumptions about which theoretical – that is

508    not directly observable – processes and entities explain observable phenomena

509    (Laudan, 1981). This is true of theories and models of the climate system. Prior to the

510    1990s, for example, climate models that were used to provide spatial simulations of

511    global surface temperatures did not include a representation of the role of aerosols in

512    the climate system and this turned out to be a surprisingly substantial incompleteness

513    in the simulations (Wigley, 1994). Moreover, current candidates for substantially

514    underestimated forcing, feedbacks and internal variability exist (e.g., terrestrial

515    biogeochemical feedbacks (Arneth et al., 2010) and feedbacks amplifying the effects

516    of solar luminosity (Kirkby, 2007)).

517    Some philosophers have concluded, largely on the basis of the history of

518    successful but superseded theories and models, that a theory or model's predictive

519    success should not be used to justify belief in what the theory or model tells us about

520    theoretical entities and processes (see, e.g., Stanford (2006)). On their view, theories

521    and models should be taken to be no more than tools for predicting observable

522    phenomena. The sad truth, however, is that it is currently unclear what we are entitled

523    to assume about how complete empirically successful theories and models are (see

524    Saatsi (2005) and Psillos (1999) for two of many further alternative perspectives on

525    this unresolved issue). In particular, it is unclear what we are entitled to assume about

526    how complete climate models and our knowledge of the climate system are, including

527    about how complete our knowledge of climatic factors that are materially relevant to

528    CMP accuracy is. This complicates assessment. For example, difficulty in estimating

529    the completeness of GCMs' representations of the effects of solar luminosity

530    fluctuations means difficulty in assessing projections of GMST trends.

531

532    **3.6 Criteria of adequacy for approaches to assessing projections**

533    Our discussion of why assessing projection quality is difficult helps to spell out

534    criteria of adequacy for approaches to such assessment. Adequate approaches will,

535    given initial condition inaccuracy, have to assess projection quality in light of the

536    possible path dependent nature of error propagation. Given the inevitable use of

537    parameterization, they will have to take the possible effects of tuning into account.

538    They will also have to take the impact of model imperfection into account. Doing so

539    involves paying attention to climate models' limited ability to simulate climate, to the

540    difficulty in determining which aspects of model empirical success are relevant to

541    assessing which projections, to the interdependence of the evolution of climatic

542    quantities along with the effect of this interdependence on error propagation and to

543    possible structural instability. Doing so also requires attending to the history induced

544    lack of clarity about unrecognized model imperfection. If the claim is that we are

545    entitled to ignore the history of successful but superseded models and thus to cease

546    worrying about unrecognized model imperfection, we need to be told why. Otherwise,

547    the impact of unrecognized climate model limitations on the accuracy of their

548    projections needs to be taken into account.

549         Since we know that only some of the projections of climate models will be

550    accurate, an adequate approach to assessing projection quality will have to provide

551    projection (or class of projections) specific assessments (Gleckler et al., 2008; Parker,

552    2009). It should judge the quality of a CMP on the basis of how fit the model or

553    models which generated it are for the purpose of doing so, i.e., for the purpose of

554    correctly answering the question the CMP answers.

555

556    **4. The confidence building approach**

557    We now discuss the confidence building approach to assessing projection quality.

558    This approach, recall, focuses on what model agreement with physical theory as well

559    as model simulation accuracy confirm. Better grounding in physical theory and

560    increased accuracy in simulation of observed and past climate is used to increase

561    confidence in models and hence in CMPs. Given the emphasis on grounding in

562    physical theory, the reliance here is primarily on GCMs.

563         In the uncertainty assessment guidance note for WG1 AR4 lead authors,

564    degrees of confidence in models are interpreted probabilistically. Specifically, they

565    are calibrated to chance ranges, e.g., very high confidence in a model is interpreted as

566    its having an at least 9 in 10 chance of being correct (IPCC 2005). The chance that a

567    model is correct can be thought of as the model's propensity to yield correct results

568    with a certain frequency, but neither the guidance note nor the report itself indicate

569 how chances should be interpreted. Indeed, they do not indicate how the talk of

570 chances of models' being correct relates to the talk of CMP likelihoods, and the report

571 does not go beyond establishing increased confidence in models in order to assign

572 them specific degrees of confidence. This last fact makes it unclear how the report's

573 use of 'increased confidence' relates to the explication of degrees of confidence in

574 terms of chances. Better grounding in physical theory is illustrated by the, at least

575 partly theoretically motivated, inclusion in some GCMs of interactive aerosol modules

576 (Randall et al., 2007). Illustrations of improved simulation accuracy are given below.

577

578 **4.1 Initial condition inaccuracy and the confidence building approach**

579 WG1 AR4 states that many climatic quantities of interest, including those relating to

580 anthropogenic climate change, are much less prone to nonlinear sensitivity to initial

581 conditions than weather related quantities and are thus more amenable to prediction

582 (Le Treut et al., 2007). This relative insensitivity to initial conditions is argued for

583 primarily on the basis of GCM simulations in which initial conditions are varied.

584 Notably, CMIP3 multi-model simulations of $20^{th}$ century GMST, in which ranges

585 reflect different initial condition runs of participating models, suggest little internal

586 variability in GMST over periods of decades and almost none over the whole century

587 (See Fig. 1 and (Hawkins and Sutton, 2009)).

588        WG1 AR4 acknowledges that confidence in simulations of response to

589 changes in initial conditions depends on resolving worries about the effects of

590 relevant model imperfection (Meehl et al., 2007b). But the claim is that these worries

591 can be mitigated by examining how well GCMs simulate important sources of the

592 climate system's nonlinear responses, e.g., the El Niño – Southern Oscillation (ENSO)

593 and the MOC. Thus, the ability of GCMs to simulate observed nonlinear change in the

594     Atlantic MOC in response to fresh water influx has been used to argue that they can

595     produce reliable projections of aspects of 21[st] century MOC behavior but that

596     confidence in projections beyond the 21[st] century is very limited (Pitman and Stouffer,

597     2006).

598          Computational resources, however, only allowed a very limited range of initial

599     conditions to be explored by CMIP3 GCMs (CMIP3, 2007). As to the question of the

600     extent to which GCM ability to simulate (in)sensitivity to initial conditions does help

601     with assessment in light of model imperfection and tuning, it is addressed in the

602     following sections. Here we only note that the need to address this question has been

603     made pressing since WG1 AR4. Recent work suggests that GCMs do not adequately

604     capture the structure of the climate system prior to abrupt changes in the past and are,

605     in some circumstances, insufficiently sensitive to initial conditions. They can, for

606     example, only simulate the cessation of the MOC under about 10 times of the best

607     estimate of actual fresh water influx that has brought it about in the past (Valdes,

608     2011). There is, in addition, a spate of studies according to which CMIP3 GCMs

609     substantially underestimate the extent to which 20[th] century GMST anomalies are due

610     to internal variability, including initial condition variability, on multidecadal scales

611     (Semenov et al., 2010; Swanson et al., 2009; Wu et al., 2011). Some work suggests

612     that the underestimates extend to periods of 50 to 80 years in length (Wyatt et al.,

613     2011).

614          Recognizing the potential significance of initial conditions to improving

615     multidecadal CMPs, some recent work aims to take on the challenge of limited

616     available data in order to initialize simulation runs to actual observed initial

617     conditions (Hurrell et al., 2009). More extensive exploration of the impact of varying

618     GCM simulation initial condition settings is also being carried out (Branstator and

619     Teng, 2010).

620

621     **4.2 Parameterization, tuning and the confidence building approach**

622     WG1 AR4 addresses the difficulty of assessing projection quality in light of tuning by

623     taking increased simulation accuracy to increase confidence in models only when this

624     accuracy is not a result of direct tuning, i.e., only when it is not the result of tuning a

625     parameter for a certain quantity to observations of that quantity (Randall et al., 2007,

626     p. 596). But tuning can be indirect. GCMs do not possess parameters for GMST

627     trends, and thus cannot be directly tuned to observations of these trends. Nevertheless,

628     there is (CCSP, 2009) substantial uncertainty about radiative forcings, and especially

629     about aerosol forcing, allowing forcing parameters to be tuned to yield close

630     agreement between simulated and observed $20^{th}$ century mean GMST trends (Fig. 1).

631     That this tuning occurs is, as is widely recognized within the climate science

632     community, suggested by the observation that different models achieve such

633     agreement by substantially different combinations of estimates of climate sensitivity

634     and radiative forcing [CCSP, 2009; Knutti, 2008b].

635        The difficulty in assessing projection quality in light of parameterization

636     limitations is partly, if implicitly, addressed by noting improvements in

637     parameterization schemes since the publication of WG1 TAR. As schemes that

638     incorporate a better understanding of the climate system and show better agreement

639     with data become available, we acquire a better understanding of the limitations of

640     older schemes and increase trust in model performance. Such improvement, however,

641     leaves open the question of how to handle worries about tuning. Moreover, increased

642     quality of parameterizations does not indicate how to assess the impact of the

643    inevitable remaining underdetermination in parameterization choice on projection

644    quality. Thus, it remains unclear how accurate CMPs actually are.

645         Another strategy that is not explicitly discussed in WG1 AR4, but which is

646    consistent with the confidence building approach, is suggested by the idea that

647    grounding in basic theory increases confidence in models. Perhaps, in some cases, the

648    role of basic theory in generating CMPs is sufficient so as to eliminate, or

649    substantially reduce, worries arising from the use of parameterizations. It has been

650    argued that while simulating the feedback effect of increased water vapor inevitably

651    makes use of parameterizations, this effect is dominated by processes that are

652    represented by the equations of fluid dynamics and thus will continue to be accurately

653    simulated by climate models (Dessler and Sherwood, 2009). It has also been

654    suggested that, since GCMs use the equations of fluid dynamics, our ability to predict

655    nonlinear MOC evolution that results from its fundamental properties is beginning to

656    mature, unlike our ability to predict nonlinear evolution it might exhibit as a result of

657    terrestrial ecosystems (Pitman and Stouffer, 2006).

658         One difficulty here is how to determine that properties represented by basic

659    physical theory largely determine the evolution of projected quantities. Insofar as

660    estimates that this is so rely on – as, e.g., Dessler and Sherwood (2009) rely on –

661    climate model results, it is assumed that available parameterizations are adequate and

662    the reliance on parameterization is not bypassed. Further, even if we have managed to

663    isolate properties that are represented by basic theory and determine the evolution of a

664    projected quantity, we cannot escape worries relating to the use of parameterization.

665    Parameterization always plays an essential role even in descriptions of subsystems of

666    the climate for which we possess basic equations. Basic equation discretization in

667    GCMs brings with it grid-scale dependent parameterization, e.g., grid-scale dependent

668  convection parameterization, of subgrid processes. How this discretization and

669  associated parameterization affects CMP accuracy, especially in light of how it affects

670  model ability to simulate highly nonlinear dynamics, needs adequate treatment.

671

672  **4.3 Structural inadequacy and the confidence building approach**

673  Increased model grounding in basic physical theory and increased accuracy in

674  simulation results across a range of such results does indicate increased structural

675  adequacy. Moreover, confidence building exercises do typically acknowledge a wide

676  variety of model limitations. What we need, however, are arguments connecting

677  increased success with the quality of specific classes of CMPs. This includes

678  arguments addressing the issue of how total remaining inadequacy affects CMP

679  quality.

680  Thus, for example, WG1 AR4 offers information such as that more state-of-

681  the-art models no longer use flux adjustments, that resolution in the best models is

682  improving, that more physical processes are now represented in models and that more

683  such processes are explicitly represented (Randall et al., 2007). But we need

684  arguments that connect these successes to an overall estimate of remaining structural

685  inadequacy and tell us what this inadequacy means for the quality of specific classes

686  of CMPs. It is one thing to be shown that simulated multi-model mean surface

687  temperatures are, outside of data poor regions, usually within 2 °C of observed

688  temperatures, another to be shown how this information bears on the quality of CMPs

689  of mean surface temperature trends and yet another to be shown how it bears on the

690  quality CMPs of mean precipitation trends.

691  While the needed arguments can be further developed, it remains to be seen

692  how far they can be developed. Further, it is likely that these arguments will, to a

693    substantial extent, be based on theory and expert judgment, thus limiting the extent to

694    which the confidence building approach is model based.

695

696    **4.4 The appeal to paleoclimate**

697    An important distinction needs to be made between model ability to simulate 20[th]

698    century climate and model ability to simulate paleoclimate. The latter provides

699    opportunities for out-of-sample testing, as WG1 AR4 notes (Jansen et al., 2007, p.

700    440). Such testing is of particular significance as it has the potential to help in

701    addressing the question of the extent to which tuning to current climate is a problem.

702    Indeed, there is growing recognition of the importance of palaeodata, including of its

703    importance for model assessment (Caseldine et al., 2010). In this context, there is an

704    ongoing debate about whether to conclude that GCMs lack representations of crucial

705    mechanisms/feedbacks because these models have difficulties in accurately

706    simulating past warm, equable climates with a weak equator-to-pole temperature

707    gradient (Huber and Caballero, 2011; Spicer et al., 2008).

708        Although this may change in the future, the burden of assessing models in

709    light of data nevertheless currently rests firmly on the ability of models to simulate

710    recent climate. This is so for at least three reasons. First, simulation experiments with

711    paleodata are still limited. WG1 AR4's appeal to such simulations is confined

712    primarily to two instances. WG1 AR4 uses model ability to simulate aspects of the

713    climate system during the Last Glacial Maximum (LGM) in order further to support

714    the claim that models have captured the primary feedbacks operating in the climate

715    system at the time (Jansen et al., 2007, p. 452). WG1 AR4 also uses model ability to

716    simulate climate responses to orbital forcing during the mid-Holocene in order to

717    improve confidence in model ability to simulate responses to such forcing (Jansen et

718    al., 2007, p. 459). Second, most of the models WG1 AR4 relies on in generating

719    projections are not among the models it relies on in discussing paleoclimate

720    simulations (Schmidt, 2010). And when the same models are relied on in both

721    contexts, model resolution usually varies across the contexts  (Braconnot et al., 2007).

722    Practical constraints mean lower resolution models have to be used to simulate

723    paleoclimate. Thus it is unclear what the paleoclimate simulation successes allow us

724    to conclude about model fitness for the purpose of generating projections. Third, there

725    are substantial, unresolved issues about how uncertain paleoclimate reconstructions

726    are, and thus about what we can learn from them (Snyder, 2010; Wunsch, 2010).

727

728    **4.5 Inter-model results, robust projections and the confidence building approach**

729    The confidence building approach is strengthened, both in WG1 AR4 and elsewhere,

730    by noting that state-of-the-art GCMs provide a robust and unambiguous picture of the

731    evolution of some large scale features of climate. Such multi-model results are

732    supposed to increase confidence in projections. For example, state-of-the-art GCMs

733    predict that GMST evolution will be roughly linear over much of this century, thus

734    supposedly reducing worries about the sensitivity of such evolution to initial condition

735    changes and to minor variations in model structure (Knutti, 2008).

736        How does the appeal to multi-model results help in assessing projection

737    quality, as opposed to improving projection accuracy? We outline two views about

738    how it does so and then critically discuss these views.

739        A common assumption in formal analyses of multi-model ensemble results,

740    and to some extent in applications of the confidence building approach, is that model

741    errors are independent of each other and thus tend to cancel out in calculations of

742    multi-model means (Meehl et al., 2007b; Palmer et al., 2005; Tebaldi and Knutti,

743    2007). Indeed, there is empirical evidence that multi-model means are more accurate

744    than are the results of individual models (see Gleckler et al. (2008) as well as, for

745    further references, Knutti et al. (2010)). Given the assumptions of error independence

746    and of error cancellation, one could argue that we can expect a reduction of error in

747    ensemble means with increased model numbers and thus can take the number of

748    models used in generating means to be an indicator of CMP quality (Tebaldi and

749    Knutti, 2007).

750    In addition, or alternatively, one can assume that ensemble models are to some

751    extent independent of each other in that they explore alternative model structures and

752    parameterizations that are consistent with our knowledge of the climate system

753    (Murphy et al., 2007). Ensemble projection ranges can then be viewed as at least

754    partial explorations of our uncertainty about the climate system and can thus be used

755    to tell us something about projection quality. One might suggest, in particular, that the

756    greater the extent to which the range of uncertainty is explored by an ensemble, the

757    greater the extent to which the projections/projection ranges it produces are robust or

758    insensitive to uncertain assumptions and thus the more probable these results are

759    (Weisberg (2006) describes the general logic behind appeals to robustness). Multi-

760    model ensemble projection ranges are sometimes interpreted probabilistically, e.g.,

761    the range of generated projections is supposed to span the range of possibilities and

762    each projection is assigned a probability equal to the fraction of models that generate

763    it (as in Räisanen and Palmer (2001) and, to some extent, in WG1 TAR (IPCC 2001)).

764    The appeal to multi-model results does not, and is not intended to, address the

765    issue of tuning or the difficulty of figuring out what to infer about the quality of

766    specific CMPs from the partial empirical successes of models. Further, worries about

767 the use of multi-model ensembles have been raised both within and without climate

768 science.

769      Philosophers have pointed out that individual model error can only cancel out

770 to a limited extent because limited knowledge and limited computational resources

771 mean that where one model's error is not repeated by another model, the other model

772 will probably have to introduce a different error (Odenbaugh and Alexandrova, 2011).

773 Limited knowledge and limited computational resources also mean that substantial

774 model imperfection will inevitably be shared across models in ensembles (Odenbaugh

775 and Alexandrova, 2011). Multi-model ensembles in all fields of research accordingly

776 inevitably leave us with substantial error the impact of which on results is not

777 estimated. So, while coming to rely on multi-model ensembles might entitle us to be

778 more confident in projections than we would have been otherwise, it does not appear

779 to allow us to assign qualities that, like probabilies and informative probability

780 ranges, involve specifying the full range of possible evolutions of projected quantities.

781      Climate scientists' examination of GCM ensemble results confirms that such

782 ensembles only provide limited improvement in agreement with empirical data and

783 that much of the remaining disagreement arises from biases that are systematic across

784 ensemble members (Knutti et al., 2010). For present day temperature, for example,

785 half of the bias exhibited by the ensemble of models used by CMIP3 would remain

786 even if the ensemble were enlarged to include an indefinite number of models of

787 similar quality (Fig. 4). The observation that models share model imperfections is also

788 acknowledged in climate science research, including in WG1 AR4. Climate modelers

789 tend to aim at constructing the best models they can for their shared purposes and in

790 doing so inevitably use shared knowledge and similar technology. As a result, climate

791 models tend to be similar, sharing many of the same imperfections (Allen and Ingram,

792   2002; Knutti, 2010; Meehl et al., 2007b; Stainforth et al., 2007a; Tebaldi and Knutti,

793   2007).

794        A related problem is that, although model limitations are extensively examined

795   in the literature, discussion of the extent to which models in specific multi-model

796   ensembles differ in ways that are relevant to assessing projections is limited (Knutti et

797   al., 2010).

798        Recognizing the limited extent to which model error cancels out, some climate

799   scientists have suggested that we should not assume that the larger the ensemble the

800   closer means are to representing reality. Instead, they suggest, one should assume that

801   the correct climate and the climates simulated by models in an ensemble are drawn

802   from the same distribution, e.g., from the standard normal (Gaussian) distribution.

803   Under this new assumption, the failure of an increase in ensemble size to improve

804   simulation results is no longer interpreted as indicating systematic bias. One can then,

805   the suggestion is, assume that when a proportion $r$ of an ensemble yield a given

806   projection, $r$ is the probability of that projection (Annan and Hargreaves, 2010). But

807   the assumption that model probability distributions coincide with the real climate

808   distribution cannot be made in general, as is illustrated in the case of the already

809   mentioned GCM inability realistically to simulate historical Atlantic MOC collapse.

810   Indeed, structural inadequacy that is known to be shared by ensemble models means

811   that we know that the correct climate *cannot* be represented by current models.

812        Let us now look at the second argument for appealing to inter-model results in

813   assessing projection quality, the one according to which multi-model ensembles allow

814   us to explore our uncertainty. Since existing climate models share many uncertain

815   assumptions, the projections/projection ranges multi-model ensembles produce do not

816   reflect full explorations of our uncertainty (Parker, 2011; Pirtle et al., 2010).

817    Moreover, once again, such ensembles do not allow assigning projection qualities the

818    assignment of which involves estimating the full range of possible evolutions of

819    projected quantities.

820         The GCMs used by WG1 AR4 only sample some of the recognized range of

821    uncertainty about aerosol forcing, perhaps because of the already mentioned tuning

822    relating to this forcing. As a result, the spread of estimated temperature anomalies

823    these models provide (Fig. 1) substantially underestimates the uncertainty about this

824    anomaly and, accordingly, would be misleading as a guide to projection quality

825    (Schwartz et al., 2007). So too, if we take the range of natural variability covered by

826    the simulations represented in Fig. 1 to reflect our uncertainty about natural variability

827    over the next three decades, we will assign a very low probability to the prediction

828    that natural variability will substantially affect GMST trends over this period.

829    Keeping in mind, however, that these models may well similarly and substantially

830    underestimate internal variability over the next 30 years would lead us to reduce our

831    confidence in this prediction. Worse, if we cannot estimate the probability that the

832    ensemble is wrong (something the ensemble cannot help us with!) about internal

833    variability here, we are not in a position to assign the prediction a probability.

834         A number of suggestions have been made within the climate science

835    community about how partially to address the above worries about the use of multi-

836    model ensembles. Assessments that are explicit about the extent to which climate

837    models in any multi-model ensemble differ in ways that are relevant to assessing

838    projection quality should be offered (IPCC 2010; Knutti et al., 2010). If, for example,

839    internal variability in the MOC is an important source of uncertainty for projections of

840    mean sea surface temperatures over the next 30 years and our ensemble is in the

841    business of making such projections, it should be clear to what extent the simulations

842    produced by the ensemble differ from each other in ways that explore how internal

843    variability in the MOC might occur. Assessing projection quality relevant differences

844    in models is a substantial task, one that goes well beyond the standard multi-model

845    exercise.

846          In addition, while limited knowledge and resources, e.g., restrictions to certain

847    grid resolutions, mean that there is no question of exploring all of existing uncertainty,

848    provision of second and third best guess modeling attempts could provide a clearer

849    picture of our uncertainty and its impact on CMP quality (Knutti et al., 2010; Smith,

850    2006).

851          A difficulty to keep in mind is that of determining how a model component

852    that is shared by complex models that differ in complex ways affects CMP quality.

853    Assessment of model components and their impact on model performance is a

854    challenge that is – because of the need to evaluate models in light of background

855    knowledge – part and parcel of assessing models fitness for purpose. This challenge is

856    complicated when the projection is generated by complex models that implement

857    common components but differ in other complex ways. For the same component may,

858    as a result, function in different ways in different models (Lenhard and Winsberg,

859    2010). Examining how a parameterization of cloud microphysics affects CMPs may,

860    for example, be hampered if the parameterization scheme is embedded in models that

861    substantially differ in other parameterizations and/or basic theory.

862          The comparison of substantially differing models will also exacerbate existing

863    challenges for synthesizing the results of multi-model ensembles. Climate scientists

864    have noted that synthesizing the results of different models using a multi-model mean

865    can be misleading even when, as in the case of the CMIP3 models, the models

866    incorporate only, and only standard, representations of atmosphere, ocean, sea ice and

867    land [Knutti et al., 2010]. For example, the CMIP3 multi-model mean of projected

868    local precipitation changes over the next century is 50 % smaller than that which

869    would be expected if we were to assume that at least one, we know not which, of the

870    CMIP3 models is correct. So it seems that using a mean in this case is misleading

871    about what the models describe (Knutti et al., 2010). Synthesizing the results of

872    different models may be even more misleading where models differ substantially in

873    how they represent processes or in which processes they represent, e.g., if some of the

874    models do and some do not include representations of biogeochemical cycles (Tebaldi

875    and Knutti, 2007). In such circumstances, for example, a mean produced by two

876    models may well be a state that is impossible according to both models.

877

878    **5. The subjective Bayesian approach**

879    Perhaps the main approach to supplement the confidence building approach in WG1

880    AR4 is the subjective Bayesian approach. We first consider this formal,

881    supplementary approach as it is used to assess projection quality in light of difficulties

882    in parameter choice (Hegerl et al., 2006; Murphy et al., 2004). We then consider how

883    it has been extended.

884

885    **5.1 The subjective Bayesian approach to parameter estimation**

886    A simple, but representative, application of the standard version of the Bayesian

887    approach to parameter, including projection parameter, estimation involves

888    calculating the posterior probability distribution function P($F$ | data, $M$) using Bayes'

889    theorem, as in Eqt. (3) (Frame et al., 2007). P($F$ | data, $M$) specifies the probabilities

890    of values of a parameter, $F$, given data and a model $M$. P(data | $F$, $M$) is the likelihood

891    of $F$ and captures, as a function of values of $F$, the probability that the data would be

892 simulated by *M*. In the Bayesian context, 'the likelihood of *F*' refers to a probability

893 function for data rather than, as it would on the WG1 AR4 use of 'likelihood', to a

894 probability range for *F*. The prior probability distribution function P(*F* | M) is the

895 probability distribution function of *F* given only *M* and thus prior to consideration of

896 the data. P(data) is a normalizing constant required to ensure that the probabilities

897 sum up to 1.

898 $$P(F \mid \text{data}, M) = P(\text{data} \mid F, M)P(F \mid M)/P(\text{data}) \tag{3}$$
899

900 The probabilities in Eqt. (3) are, on the subjective Bayesian approach, to be

901 interpreted as precise, quantitative measures of strength of belief, so called 'degrees of

902 belief'. What makes the subjective Bayesian approach subjective is that unconstrained

903 expert opinion – the beliefs of certain subjects irrespective of whether they meet

904 objective criteria of rationality such as being well grounded in empirical evidence – is

905 used as a central source for selecting prior probability distributions. Still, the

906 subjective Bayesian approach often uses uniform assignments of priors. In doing so, it

907 borrows from what is usually called 'objective Bayesianism' (see Strevens (2006b) for

908 a discussion of the different forms of Bayesian approaches to science).

909 Bayes' theorem allows us to take existing estimates of parameter uncertainty –

910 here captured by P(*F* | M) – and to constrain these using information from perturbed

911 physics experiments about how well a model simulates data as a function of parameter

912 settings – information here captured by the likelihood function P(data | *F*, M).

913 Assume experts provide prior probability distributions for parameters relating to total

914 radiative and present-day indirect aerosol forcing and that we calculate the probability

915 that a model gives, as a function of the parameters' values, to observed oceanic and

916 atmospheric temperature change. Bayes' rule can then yield posterior probability

917  distributions for the parameters (Fig. 5). Bayesian parameter estimation has tended to

918  rely on models of intermediate complexity and on energy balance models.

919  The Bayesian hope is that the constraints provided by simulation success on

920  parameter estimates will increase the objectivity of such estimates. Moreover, Bayes'

921  theorem provides, what the confidence building approach does not provide, a clear

922  mechanism that relates simulation accuracy to conclusions about CMP quality, thus

923  helping to address the problem of what to infer from available simulation accuracy

924  given the existence of model imperfection.

925  Nevertheless, the standard version of the Bayesian approach to parameter

926  estimation faces substantial problems. The standard interpretation of the probability

927  distributions $P(F \mid M)$ and $P(F \mid \text{data}, M)$ is that they are probability distributions for $F$

928  that are conditional on the correctness of a version of $M$. In the present context, what

929  is being assumed to be correct is a model version in which one or more parameters are

930  unspecified within a certain range. For the goal is to select parameter values from

931  within a range of such values. Now, it is on the basis of the standard interpretation of

932  $P(F \mid M)$ and $P(F \mid \text{data}, M)$ that standard justifications, using so-called Dutch Book

933  arguments, for updating beliefs in accord with Bayes' theorem proceed. Dutch Book

934  arguments generally assume that the, typically statistical, model versions upon which

935  probabilities are conditional are correct. It is argued that, given this assumption, the

936  believer would end up with beliefs that are not as true as they might have been, or

937  would incur a financial loss, if his or her beliefs were not updated in accord with

938  Bayes' theorem (see Jeffrey (1990) and Vineberg (2011) for examples). But if, as in

939  the cases we are concerned with, the model version upon which distributions are

940  conditional is not correct, applying Bayes' theorem may offer no advantage and may

941  be a disadvantage.

942       Assume that our subject relies on a CMIP3 GCM to determine whether a

943    specified fresh water influx will lead to a collapse in the MOC and that the specified

944    influx is a tenth of that needed to get the model to simulate collapse. Assume also that

945    some exploration of plausible parameter settings in the GCM does not alter results

946    substantially. Applying Bayes's theorem on the assumption that the model is, up to

947    plausible parameter modification, correct means that the probability we assign the

948    outcome 'collapse' is 0. The modeler acquiesces to the theorem. Unfortunately, as we

949    now know, the model's results are misleading here. In this case, not applying Bayes'

950    theorem may lead to more realistic judgments.

951       Thus, the standard use of Bayes' theorem in parameter estimation requires an

952    alternative to the standard interpretation of its conditional probabilities. We will also

953    need an alternative to the standard justifications for applying Bayes' theorem.

954       Even if we have settled on some interpretation of the conditional posterior

955    probabilities produced by Eqt. (3), there remains the question of what we can infer

956    about reality from these probabilities. There remains, in other words, the question of

957    what distribution of probabilities for $F$, $P(F)$, we should adopt given the conditional

958    distribution $P(F \mid \text{data}, M)$. We might have a probability distribution for climate

959    sensitivity that is conditional on the data and a model. But what should we infer from

960    this about actual climate sensitivity? We cannot properly answer such questions until

961    we have gone beyond assessing how parameter choice affects projection quality and

962    have also assessed how structural inadequacy, parameterization scheme choice and

963    initial condition inaccuracy do so (Rougier, 2007).

964       Rougier provides a non-standard version of the Bayesian approach to

965    parameter estimation that has the substantial advantage of allowing us to factor in

966    estimates of structural inadequacy into subjective Bayesian parameter estimates

(Rougier, 2007). Nevertheless, his work takes estimates of structural inadequacy as given and thus does not, by itself, tell us how more comprehensive assessments of projection quality are to be produced.

Additional difficulties for the Bayesian approach relate to the usage of prior probabilities. We rehearse two familiar worries about this usage. First, estimates of $P(F \mid M)$ are usually made after data that bears on the estimates is in hand and it is hard to estimate what probability distribution would be assigned to $F$ independently of knowledge of this data. Failure properly to estimate $P(F \mid M)$ may lead to counting the same data twice, once in estimating priors and once in estimating likelihoods (Frame et al., 2007).

Second, while some climate scientists have argued that the explicit setting out of subjective priors by experts is desirable because it makes subjective judgments explicit (Hargreaves, 2010), philosophers of science have pointed out that it leaves open the question of the extent to which experts' views are evidence based and thus puts reliable and unreliable priors on a par (Sober, 2002). This issue becomes particularly worrying in the context of climate modeling. We know that prior selection may be based on results involving tuning and be required even when data underdetermines parameter value choice. So there is a risk that assigning a prior to a parameter value will beg the question against alternative choices and thus yield estimates of climatic variables we are by no means obliged to accept. The worry of question begging is exacerbated by arguments to the effect that the influence of likelihoods, and thus of data, on the shape and width of prior distributions is often minor (Frame et al., 2005).

A common way of trying to minimize the impact of the appeal to expert opinion is to represent the state of ignorance that existed prior to the consideration of

likelihoods using uniform prior distributions within expert specified ranges. We have already seen that uniform distributions are not suitable for representing ignorance. Moreover, to assume a uniform prior distribution will often be to ignore knowledge we have of the relative plausibility of various prior assignments (Annan and Hargreaves, 2011; Rougier, 2007). So too, a uniform assignment of priors for one parameter will sometimes, because of the non-linear relationship between some model variables, provide a non-uniform prior assignment to another (Frame et al., 2005). It has been suggested that best practice given the worries about prior selection is to provide readers with posteriors as well as likelihoods. This would somewhat clarify the role data actually have had in determining posteriors (Frame et al., 2007).

Another way in which the influence of priors might be minimized is by repeated updating of posteriors in response to new evidence over time. As already noted, however, evidence with which to test models is mostly limited to familiar 20th century datasets. There is thus currently limited scope for successive updating of priors.

As to the idea that the appeal to likelihoods in deriving posterior probabilities will provide an objective constraint on parameter selection, it also has problems. Likelihoods measure agreement with data, irrespective of whether such agreement results from tuning (Katzav, 2011). In addition, we have seen that an adequate assessment of projection quality needs to take into account not only agreement with data, but also how error for each simulated quantity develops over projection scenarios as a function of error associated with other such quantities. Finally, there are various likelihood metrics or ways of measuring agreement with data. Choice between these and how such choice affects posteriors is only beginning to be explored (see, e.g., Ishizaki et al. (2010)).

1017

1018

1019 **5.2 The subjective Bayesian approach and multi-model ensembles**

1020 The subjective Bayesian approach has been extended to assessing multi-GCM

1021 ensemble output. This extension, which will be called the subjective Bayesian MM

1022 approach, involves taking an ensemble and producing a statistical model of its

1023 simulation results. Comparing the statistical model and available data yields a

1024 likelihood function that captures the probability the ensemble gives to the data. Bayes'

1025 theorem can then be used, in conjunction with the likelihood function and estimates of

1026 prior probability distributions for the statistical model's parameters, in order to

1027 produce a posterior probability distribution for these parameters (Furrer et al., 2007a;

1028 Furrer et al., 2007b; Leith and Chandler, 2010; Tebaldi et al., 2005; Tebaldi and

1029 Knutti, 2007).

1030 Some variants of the subjective Bayesian MM approach give each ensemble

1031 model equal weight in calculating ensemble posterior probability distributions (Leith

1032 and Chandler, 2010). Other variants weight the contribution of each ensemble model

1033 to posteriors as a function of how well the model simulates aspects of the climate

1034 system (Tebaldi et al., 2005).

1035 Many analyses, e.g., those in WG1 TAR  and some of those in WG1 AR4, of

1036 multi-model ensemble results produce projections that are just averages of individual

1037 model results and that have uncertainty ranges which reflect inter-model variability.

1038 This does not yield probabilistic estimates of multi-model ensemble results. The

1039 subjective Bayesian MM approach does yield such estimates. The hope is that doing

1040 so helps to take into account structural inadequacy and limited knowledge of how to

1041 select parameterization schemes. The subjective Bayesian MM approach does not

1042 explicitly tackle the issue of how initial condition inaccuracy affects CMP quality.

1043          The subjective Bayesian MM approach suffers from many of the problems of

1044          the subjective Bayesian approach to parameter estimation. The subjective Bayesian

1045          MM approach faces the problems that arise from the use of prior probabilities. It also

1046          suffers from the problems relating to the choice of likelihood metrics and the failure

1047          to take into account how error for each simulated quantity develops as a function of

1048          error associated with other such quantities. Even weighting models in assessing

1049          projection quality is not a clear advantage given that the data used to do so may well

1050          have already been used in model construction.

1051          Finally, there remain the issues of how to interpret the conditional

1052          probabilities used in Bayes' theorem given model imperfection and of how the

1053          conditional probabilities produced by Bayes' theorem relate to unconditional

1054          probabilities. On the subjective Bayesian MM approach, one updates priors on the

1055          assumption that the statistical model of multi-model ensemble results is correct.

1056          However, given that we know that multi-model ensemble results are biased, this

1057          assumption is false. And any inference from probabilities that are conditional upon

1058          data and an ensemble to unconditional probabilities can only be made given a full

1059          assessment of the effects of initial condition error and model imperfection on CMP

1060          accuracy. We have seen, however, that multi-model ensembles do not provide such an

1061          assessment.

1062

1063          **6. The likelihood approach**

1064          One response to the subjective Bayesian approach's difficulties with subjective prior

1065          probabilities is to try to avoid the use of priors all together. This is what the likelihood

1066          approach does using GCMs. It aims to produce probability distributions for

1067          parameters solely in light of how well models simulate data as a function of parameter

1068     settings, that is solely in light of likelihood functions such as P(data | $F$, $M$) (Allen et

1069     al., 2006). Doing so requires not discounting any parameter settings prior to

1070     simulation and thus providing likelihood functions that span a much broader range of

1071     parameter values than is usual. This has become possible, though usually only in

1072     experiments that perturb the parameters of a single model structure, with the

1073     distributed computing techniques used by climateprediction.net (Frame et al., 2007).

1074     The results of such attempts are distributions that are less biased due to those

1075     parameters that are perturbed, but that are far broader than those otherwise produced.

1076     An application of the likelihood approach is as follows: we take the climate

1077     sensitivities of each of a multi-thousand climateprediction.net ensemble of GCM

1078     variants and estimate the true climate sensitivity to be a weighted sum of these

1079     sensitivities. The weight of each sensitivity is determined by the probability the

1080     variant it belongs to gives to observations of a number of climatic quantities,

1081     including mean sea level temperature, precipitation and surface heat fluxes (Piani et

1082     al., 2005).

1083     The likelihood approach can also be used to minimize the impact of structural

1084     inadequacy and uncertainty about choice of parameterization scheme on CMP

1085     accuracy. It can do so by producing assessments that are only based on the best

1086     simulations available for specific parameter settings (Sanderson et al., 2008). But

1087     focusing on best results does not take into account how they are affected by initial

1088     condition inaccuracy, tuning or aspects of model imperfection other than parameter

1089     choice uncertainty. The same is true of what might be called the multi-model

1090     likelihood approach. This approach uses correlations between GCMs' predictions of

1091     trends for a quantity and related observations formally to select the best predictions

1092     (Boe et al., 2009; Shukla et al., 2006).

1093

1094 **7. Putting it all together**

1095 As we have noted, WG1 AR4 often uses expert judgment that takes the results of the

1096 approaches we have been discussing, as well as partly model-independent approaches,

1097 into consideration in assigning final projection qualities. Insofar as final assignments

1098 are model based, however, the shared limitations of the approaches we have been

1099 discussing remain untouched. In particular, insofar as final assessments are model

1100 based, they face serious challenges when it comes to assessing projection quality in

1101 light of structural inadequacy, tuning and initial condition inaccuracy. Moreover, they

1102 continue to be challenged by the task of assigning probabilities and informative

1103 probability ranges to projections.

1104

1105 **8. Assessing projections: what else can be done?**

1106 We now examine approaches that differ from those that play center stage in WG1

1107 AR4. The first approach, the possibilist approach, is described in the climate science

1108 literature but is primarily non-probabilistic. The remaining approaches are

1109 philosophy-of-science-based approaches. There are currently four main, but not

1110 necessarily mutually exclusive, philosophical approaches to assessing scientific

1111 claims. One of these is the already discussed subjective Bayesian approach. The other

1112 three are those that are discussed below.

1113

1114 **8.1 The possibilist approach**

1115 On the possibilist approach, we should present the range of alternative projections

1116 provided by models as is, insisting that they are no more than possibilities to be taken

1117 into account by researchers and decision makers and that they provide only a lower

1118    bound to the maximal range of uncertainty (Stainforth et al., 2007a; Stainforth et al.,

1119    2007b). Climate model results should, accordingly, be presented using plots of the

1120    actual frequencies with which models have produced specific projections (as in Fig.

1121    6). At the same time, one can supplement projected ranges with informal, though

1122    sometimes probabilistic, assessments of confidence in projections that appeal, as the

1123    confidence building approach appeals, to inter-model agreement and agreement with

1124    physical theory (Stainforth et al., 2007a).

1125    Informal approaches to assessing projection quality must address the same

1126    central challenges that quantitative approaches must address. So, insofar as the

1127    possibilist position allows informal probabilistic assessments of projection quality, it

1128    must address the difficulties that all probabilistic approaches face. However, one

1129    could easily purge the possibilist approach of all probabilistic elements and assess

1130    projections solely in terms of their being possibilities. Moreover, there are obvious

1131    ways to develop purely possibilistic assessment further. Purely possibilistic

1132    assessment can, in particular, be used to rank projections. Possibilities can, for

1133    example, be ranked in terms of how remote they are.

1134    The purged possibilist approach would still face challenges. Presenting CMPs

1135    as possibilities worthy of consideration involves taking a stance on how CMPs relate

1136    to reality. For example, if we are presented with an extreme climate sensitivity range

1137    of 2 to 11 K (Fig. 6) and are told that these are possibilities that should not have been

1138    neglected by AR3 WG1's headline uncertainty ranges (Stainforth et al., 2005), a claim

1139    is implicitly being made about which climate behavior is a real possibility. It is

1140    implied that these possibilities are unlike, for example, the possibility that the United

1141    States will more than halve its budget deficit by 2015. Thus a possibilist assessment of

1142    projection quality needs to be accompanied by an examination of whether the

1143    projections are indeed real possibilities. The same considerations apply to 'worst case

1144    scenarios' when these are put forward as worthy of discussion in policy settings or

1145    research. The threat that arises when we do not make sure that possibilities being

1146    considered are real possibilities is that, just as we sometimes underestimate our

1147    certainty, we will sometimes exaggerate our uncertainty.

1148          Nevertheless, the challenges facing purely possibilistic assessment are

1149    substantially more manageable than those facing probabilistic assessment. To say that

1150    something is a real possibility at some time $t$ is, roughly, to say that it is consistent

1151    with the overall way things have been up until $t$ and that nothing known excludes it

1152    (see Deutsch (1990) for a similar definition). A case for a projection's being a real

1153    possibility can, accordingly, be made just by arguing that we have an understanding of

1154    the overall way relevant aspects of the climate system are, showing that the

1155    projection's correctness is consistent with this understanding and showing that we do

1156    not know that there is something that ensures that the projection is wrong. There is, as

1157    observed in discussing probabilistic representations of ignorance, no need to specify a

1158    full range of alternatives to the projection here. Further, state-of-the-art GCMs can

1159    sometimes play an important role in establishing that their projections are real

1160    possibilites. State-of-the-art GCMs' projections of GMST are, for example and given

1161    the extent to which GCMs capture our knowledge of the climate system, real

1162    possibilities.

1163

1164    **8.2 The critical approach**

1165    The first philosophy-of-science-based approach that is not discussed in the IPCC

1166    reports and that will be discussed here is the critical approach (Freedman, 2009;

1167    Longino, 1990). According to this approach, scientific claims are rational to the extent

1168    that they result from open, critical discussion. Longino offers a prominent view of

1169    what such discussion involves. She holds that open critical discussion occurs within a

1170    community to the degree that the community has recognized avenues for criticism of

1171    evidence, methods, assumptions and reasoning; the community's members share

1172    standards of criticism; the community is responsive to criticism and intellectual

1173    authority is shared equally among qualified members (Longino, 1990). Petersen offers

1174    what can be thought of as a version of the critical approach, one that is designed to

1175    assist in, among other things, assessing CMP quality. He provides procedures, and a

1176    classification of types of uncertainty, that are supposed to help systematizing

1177    qualitative assessments of model assumptions and thus to facilitate open, critical

1178    discussion of the quality of model-based-claims (Petersen, 2012).

1179       The motivation for the critical approach is twofold. On the one hand,

1180    according to its proponents, critical discussion allows overcoming individual

1181    subjective bias. On the other hand, there are no available standards beyond our current

1182    standards by which scientific claims can be judged. So, it is argued, rationality cannot

1183    amount to more than the application of available standards of critical discussion and

1184    the acceptance of the deliverances of these standards.

1185       The critical approach is not really an alternative to the approaches used in

1186    WG1 AR4. Rather it is a framework that tells us in what conditions the deliverances

1187    of these approaches are acceptable. Petersen's framework could, for example, be used

1188    to guide applying the confidence building approach.

1189       Further, according to the critical approach, we can recognize that an

1190    assessment of the quality of a projection is limited while nevertheless accepting the

1191    projection. For, on this approach, where acceptance of models' fitness for the purpose

1192    generating projections is a result of open, critical discussion, accepting the models'

1193    projections is reasonable even if the discussion in question has substantial limitations,

1194    e.g., if the impact of unknown structural inadequacy on the projections has not been

1195    taken into account. The critical approach would thus, for example, warrant trust in

1196    state-of-the-art GCMs for the purpose of generating the GMST projections presented

1197    in Fig. 2, subject to expert correction in light of known GCM limitations and provided

1198    that the trust results from open, critical discussion.

1199        Acceptance of models' fitness for purpose can, however and as Longino's

1200    criteria for such criticism state, only be the result of open, critical discussion if there

1201    are shared standards for assessing fitness for purpose. In the absence of shared

1202    standards, agreement will be the result of the arbitrary preference of some standards

1203    over others rather than the uptake and assessment of relevant alternatives. In the case

1204    of CMP assessment, what we need for acceptance of model fitness for purpose to be

1205    the result of open, critical discussion is agreement about issues such as whether

1206    assessment should be probabilistic, whether it should be formal and so on. The present

1207    paper makes it clear, however, that it would be premature to agree on these issues and,

1208    indeed, that there is no such agreement.

1209        A more general worry about the critical approach is that, by itself, it leaves

1210    unaddressed the question of when the results of open, critical discussion are reliable

1211    (Crasnow, 1993). Unless we have an assessment of how reliable current critical

1212    discussion of model fitness for purpose is, it is unclear why we should accept the

1213    results of such discussion.

1214

1215    **8.3 Inference to the best explanation and climate model evaluation**

1216    The next philosophy based approach to assessing projection quality is the inference to

1217    the best explanation (IBE) approach (Lipton, 2004). In discussing the confidence

1218  building approach we saw model confidence being increased on the basis of

1219  improvement in model virtues such as agreement with background knowledge

1220  (including grounding in basic theory), increased realism, agreement with observations

1221  and model scope – that is, roughly, the number of distinct classes of facts a model

1222  simulates. An additional model virtue that is appealed to in climate modeling

1223  (Shackley, 1997) but is not explicitly discussed in WG1 AR4 is simplicity – which is

1224  directly tied to the number and complexity of model assumptions. Yet WG1 AR4

1225  does not, recall, tell us how to map combinations of model virtues onto non-

1226  comparative assessments of model confidence. It tells us when confidence should be

1227  increased on the basis of model virtues but not when confidence should be high. The

1228  IBE approach does and does so in a way that aims to take structural inadequacy into

1229  account.

1230      Theories and models explain phenomena in the sense that they provide

1231  derivations or simulations that show how phenomena are caused or fit into broader

1232  patterns of phenomena (Bokulich, 2011). Thus, GCMs can be said to explain GMST

1233  trends and rising sea levels because the simulations they provide show how these

1234  phenomena causally depend on anthropogenic greenhouse gas trends. How good the

1235  explanations of a model or theory are depends on what combination of virtues the

1236  model or theory has. How good a climate model's explanations are, for example,

1237  depends on how accurate its simulations are, how detailed its descriptions of climatic

1238  mechanisms are, the extent to which it can simulate climate in different periods and so

1239  on. This allows proponents of the IBE approach to propose that how confident we

1240  should be in a theory or model depends on how good the explanations it provides are,

1241  and thus on how good its virtues make its explanations (Lipton, 2004; Thagard, 1978).

1242    That is, it allows the proposal that IBE determines how confident we should be in our

1243    explanations. IBE, as applied to models, is just that form of inference which involves:

1244    (i)    the possession of alternative explanations of a body of data, where

1245           each alternative explanation rests on a model that explains the data;

1246    (ii)   a determination of which of the available alternative models that

1247           explain the data provides the best available explanation of the data, i.e.,

1248           of which of these models has the best combination of explanatory

1249           virtues;

1250    (iii)  an inference to the approximate truth of that model which provides the

1251           best available explanation, provided that the model explains the data

1252           well enough (this standard presentation of IBE has been adapted from

1253           Katzav (2012)).

1254        Since very successful theories do turn out to suffer from unexpected

1255    imperfections, even the most optimistic proponents of the IBE approach only allow

1256    that the very best explanations are good enough. Explanations that are good enough

1257    are usually identified with explanations that are not only empirically successful,

1258    simple, of wide scope and well grounded in background knowledge but that also

1259    provide confirmed novel predictions, that is confirmed predictions of phenomena that

1260    were out-of-sample when they were made *and* unexpected at the time. The idea

1261    behind this stringent definition is that, while it is true that the history of science

1262    provides examples of successful theories and models that have turned out to be

1263    fundamentally wrong, those theories or models which generate confirmed novel

1264    predictions arguably tend to survive, at least as approximations, in later theories (see

1265    Psillos (1999, pp. 101-111) for a standard discussion). Newtonian mechanics is one of

1266    the most successful theories ever, and it lead to its share of novel and confirmed

1267 predictions. Of course, like the already mentioned Newtonian Earth-Sun models,,

1268 Newtonian mechanics appears to be fundamentally wrong in many ways. But

1269 Newtonian mechanics can still be argued to be approximately true. After all, general

1270 relativity does show that we can recover the central equations of Newtonian

1271 mechanics given the right approximations.

1272       Unfortunately, IBE does not provide a way of assessing the quality of specific

1273 classes of CMPs from climate model successes. The IBE approach, like the

1274 confidence building approach in WG1 AR4, provides a way of establishing

1275 confidence in models as wholes (Katzav, 2012).

1276       Further, how accurate a climate model is depends not only on how good its

1277 explanations are but also on how well its parameterization schemes have been

1278 engineered to compensate for our limited ability to model climate. So confidence in a

1279 climate model, or in its fitness for some purpose, should not depend solely on the

1280 quality of its explanations (Katzav, 2012). As to the question whether, in any case,

1281 climate models' explanations are good enough to warrant inferring their approximate

1282 correctness, it is too complex to be addressed here.

1283       We also need to note the dispute about whether IBE should be relied on. When

1284 asked why we should think that IBE allows us to infer the approximate correctness of

1285 models when the future might provide us with surprises about model imperfection,

1286 proponents of IBE answer that we can only explain the success of our models by

1287 supposing that they are approximately true. The success of models would, otherwise,

1288 be a miracle (see, e.g., Musgrave (1988) and Worrall (2010)). Winsberg, however,

1289 provides examples of highly successful principles that do not appear to be

1290 approximately true (Winsberg, 2006). Opponents of IBE point out, further, that the

1291 justification of IBE is itself a kind of IBE and thus begs the question of whether IBE

1292    is acceptable (Laudan, 1981). The justification aims to get us to trust IBE on the

1293    grounds that the best explanation for the successes of a model is its approximate truth.

1294    Some, partly in light of the circular justification of IBE, aim to eschew IBE all

1295    together. Others, accepting that IBE cannot future proof our estimates of how good

1296    our models are, weaken IBE so that it is a form of inference that allows us to rank

1297    models according to explanatory capacity but that leaves open the question of how our

1298    best models relate to the truth. Yet others insist that IBE is fine roughly as it is,

1299    arguing that it is impossible, on pain of an infinite regress, to provide non-circular

1300    justification of all basic inferential principles and that IBE is a good candidate

1301    fundamental principle for justifying models and theories (see Psillos (1999) for a

1302    discussion of some of these views).

1303

1304    **8.4 Severe testing, climate models and climate model projections**

1305    The remaining approach to assessing scientific claims that we will discuss is the

1306    severe testing approach. The idea behind the severe testing approach is that the

1307    deliberate search for error is the way to get to the truth. Thus, on this approach, we

1308    should assess scientific claims on the basis of how well they have withstood severe

1309    testing or probing of their weaknesses (Mayo, 1996; Popper, 2005; Rowbottom,

1310    2011). There are a variety of definitions of 'severe test'. One prominent definition is

1311    Mayo's (Mayo, 1996; Parker, 2008). It, however, requires that for a test of a claim to

1312    be severe it must be very unlikely that the claim would pass the test if the claim were

1313    false, a requirement that very few tests of climate model fitness for purpose fulfill and

1314    thus which would render the severe testing approach largely unhelpful here. We,

1315    accordingly, explore the usefulness of the main alternative definition, which is

1316    Popper's.

1317        According to Popper, an empirical test of a theory or model is severe to the

1318 extent that background knowledge tells us that it is improbable that the theory or

1319 model will pass the test. Background knowledge consists in established theories or

1320 models other than those being tested (Popper, 2002, p. 150). Popper offers the 1919

1321 test of general relativity's prediction of the precise bending of light in the Sun's

1322 gravitational field as an example of a severe test. The observed bending was

1323 improbable and indeed inexplicable in light of background knowledge at the time,

1324 which basically consisted in Newtonian mechanics. For similar reasons, the precise

1325 precession of Mercury also provided a severe test of general relativity.

1326        A crucial difference between the severe testing approach and the approaches

1327 pursued by WG1 AR4 is that the severe testing approach never allows mere

1328 agreement, or increased agreement, with observations to count in favor of a claim.

1329 That simulation of observed phenomena has been successful does not tell us how

1330 unexpected the data are and thus how severely the data have tested our claims. If, for

1331 example, the successful simulation is the result of tuning, then the success is not

1332 improbable, no severe test has been carried out and no increased confidence in model

1333 fitness for purpose is warranted. Notice, however, that the fact that claims are tested

1334 against in-sample data is not itself supposed to be problematic as long as the data does

1335 severely test the claims [Mayo, 1996]. This is illustrated by the prediction of the

1336 precession of Mercury. The prediction was not novel or even out-of-sample. It was

1337 well measured by Le Verrier in 1859 and was known by Einstein when he constructed

1338 his theory (Earman and Glymour, 1978). Another crucial difference between the

1339 severe testing approach and those pursued by WG1 AR4 is that the severe testing

1340 approach is not probabilistic. The degree to which a set of claims have withstood

1341 severe tests, what Popper calls their degree of corroboration, is not a probability.

1342        How might one apply a (Popperian) severe testing approach to assessing

1343   projection quality? What we need, from a severe testing perspective, is a framework

1344   that assigns a degree of corroboration to a CMP, $p$, as a function of how well the

1345   model (or ensemble of models), $m$, which generated $p$ has withstood severe tests of its

1346   fitness for the purpose of doing so. Such severe tests would consist in examining the

1347   performance of some of those of $m$'s predictions the successes of which would be both

1348   relevant to assessing $m$'s fitness for the purpose of generating $p$ and improbable in

1349   light of background knowledge. Assessing, for example, a GCM's projection of 21$^{st}$

1350   century GMST would involve assessing how well the GCM performs at severe tests

1351   of relevant predictions of 20$^{th}$ century climate and/or paleoclimate. That is it would

1352   involve assessing how well the GCM performs at simulating relevant features of the

1353   climate system that we expect will seriously challenge its abilities. A relevant

1354   prediction will be one the accuracy of which is indicative of the accuracy of the

1355   projection of 21$^{st}$ century GMST. Examples of relevant features of the climate the

1356   accurate simulation of which will be a challenge to IPCC-AR5 models are the effects

1357   of strong ENSO events on the GMST, effects of Atlantic sea surface temperature

1358   variations (associated with the MOC) on the GMST and special aspects of the GMST

1359   such as its late 30s and early 40s positive trends. That these data will challenge IPCC-

1360   AR5 models is suggested by the difficulty CMIP3 models have in adequately

1361   simulating them (Katzav, 2011).

1362        The above ideas about applying the severe testing approach will, as a step

1363   towards their operationalization, be elaborated on somewhat and put more formally. $p$

1364   is corroborated by data just in case the data are probable in light of $p$ but improbable

1365   in light of background knowledge, $B$. Symbolically, $p$ is corroborated by data just in

1366   case $P(data \mid B) < 0.5$ and $C(p \mid data, B)$ satisfies

1367 $$C(p \mid \text{data}, B) \propto P(data \mid p, B) - P(data \mid B) > 0 \qquad (4)$$

1368 Here $P(data \mid p, B)$ is the probability of the data in light of $p$ and $B$, and $P(data \mid B)$ is

1369 the probability of the data in light of $B$ alone. $C(p \mid \text{data}, B)$ itself results when the

1370 right hand side of (1) is normalized so as to yield a result that is between 1 and -1,

1371 where 1 signifies the highest degree of corroboration and -1 signifies the highest

1372 degree of falsification (Popper, 1983).

1373 Now, we want to assign a degree of corroboration to $p$ as a function of the

1374 fitness of $m$ for the purpose of generating $p$. So one could identify $P(data \mid p, B)$ with

1375 the probability that $m$ gives to data which are relevant to testing $m$'s fitness for the

1376 purpose of generating $p$, that is with $P(data \mid q, m)$, where $q$ is $m$'s prediction about the

1377 relevant data. One could also identify $P(data \mid B)$ with the probability given to the

1378 relevant data by an established rival, $m1$, to $m$, that is with $P(data \mid q1, m1)$, where $q1$

1379 is $m1$'s prediction for the data. Thus, in the context of assessing $m$'s suitability for

1380 generating $p$, (4) could be interpreted as:

1381 $$C(p \mid \text{data}, m, m1) \propto P(data \mid q, m) - P(data \mid q1, m1) > 0 \qquad (5)$$

1382 If one's focus is on assessing IPCC-AR5 projections of 21$^{st}$ century GMST, it is

1383 natural to identify the probability background knowledge gives to data with the

1384 probability the CMIP3 ensemble gives to them. Accordingly, one could, for example,

1385 calculate the degree of corroboration of the projection of GMST of a particular AR5

1386 GCM for the 21$^{st}$ century in light of the model's simulation of data relating to ENSO

1387 strength by calculating the difference between the probability the model gives to these

1388 data – $P(data \mid q, m)$ in (5) – and the probability the CMIP3 ensemble gives to them –

1389 $P(data \mid q1, m1)$ in (5).

1390 How might the severe testing approach help us with the difficulties involved in

1391 assessing projection quality? The severe testing approach allows us to bypass any

worries we might have about tuning since it only counts success that does not result from tuning, success that surely does exist, in favor of CMPs (Katzav, 2011). The severe testing approach can thus, at least, be used as a check on the results of approaches that do not take tuning into account. If, for example, the subjective Bayesian approach assigns a high probability to a projection and the severe testing approach gives the projection a high degree of corroboration, we can at least have some assurance that the probabilistic result is not undermined by tuning.

Underdetermination in choice between parameters/available parameterization schemes might also be addressed by the severe testing approach. Substituting different parameterization schemes into a model might result in varying degrees of corroboration, as might perturbing the model's parameter settings. Where such variations exist, they allow ranking model fitness for purpose as a function of parameter settings/parameterization schemes. Similarly, degrees of corroboration can be used to rank fitness for purpose of models with different structures. The resulting assessment has, like assessment in terms of real possibilities, the advantage that it is less demanding than probabilistic assessment or assessment that is in terms of truth or approximate truth. Ranking two CMPs as to their degrees of corroboration, for example, only requires comparing the two CMPs. It does not require specifying the full range of alternatives to the CMPs. Nor does it require that we take some stand on how close the CMPs are to the truth, and thus that we take a stand on the effects of unknown structural inadequacy on CMP accuracy. Popper's view is that a ranking in terms of degrees of corroboration only provides us with a ranking of our conjectures about the truth. The most highly corroborated claim would thus, on this suggestion, be our best conjecture about the truth. Being our best conjecture about the truth is, in principle, compatible with being far from the truth.

1417    Consider now some of the limitations of the severe testing approach. To begin

1418    with, while the fact that the severe testing approach is, in some respects, less

1419    demanding than other approaches has its advantages, it also have its disadvantages.

1420    Suppose we rank a claim according to degree of corroboration. What does this imply

1421    for the usability of the claim in research and in decision making? Popper's suggestion

1422    that the most highly corroborated claim is our best conjecture about the truth suggests

1423    a role for corroboration in the context of research. But when is our best conjecture

1424    close enough to the truth to be relevant to practice, e.g., to decision making (Salmon,

1425    1981)? Popper's response is not straightforward (Miller, 2005). However, one can

1426    make use of Popper's idea that claims should be assessed by severe tests without

1427    buying into the rest of his views about science. The beginnings of an alternative

1428    response is as follows: the overall degree of corroboration of a claim depends not just

1429    on how the claim has done at this or that single test, but also on how broadly it has

1430    been tested. A claim's degree of corroboration is thus correlated with the extent to

1431    which the claim is consistent with the overall way things are and, therefore, with the

1432    extent to which the claim is a real possibility. A high enough degree of corroboration

1433    will, accordingly, allow us to conclude that a claim is a real possibility and that it

1434    should be used in decision making.

1435    Another basic worry is that our description of the severe testing approach

1436    presupposes that we are able to determine, prior to using the severe testing approach,

1437    whether data are relevant to assessing fitness for purpose. This includes sometimes

1438    being able to determine, independently of the severe testing approach, that inaccuracy

1439    in simulating a quantity is not substantially relevant to the accuracy of projections of

1440    other quantities. But being able to provide such determinations is something we

1441    required of adequate approaches to assessing projection quality.

1442

1443 **9. Conclusion**

1444 There remain substantial difficulties for WG1 AR4's (climate-model-based)

1445 approaches to assessing projection quality, particularly because they aim at

1446 probabilistic assessment. Indeed, worries about probabilistic assessment of projection

1447 quality are increasingly being raised by those working on projection quality

1448 assessment (Parker, 2010; Smith, 2006; Stainforth et al., 2007a).

1449 The commonly used versions of the subjective Bayesian approach leave us,

1450 because of their limited ability to represent known climate model imperfection, with a

1451 puzzle about why Bayesian updating should be used. Rougier's version does allow a

1452 more complete representation of model imperfection, though it does not actually

1453 provide us with a way of assessing such imperfection. The likelihood approach was

1454 only briefly discussed. It is limited to assessment that takes uncertainty about

1455 parameter choice into account. The confidence building approach has the advantage

1456 of flexibility. It can, since confidence need not be expressed probabilistically, provide

1457 non-probabilistic assessments. So too, the argumentation it uses can in principle be

1458 extended to providing assessments of fitness for purpose, though it currently tends to

1459 stop at assessing models as such.

1460 In examining approaches not used in WG1 AR4, we saw that the similarity

1461 between the confidence building and IBE approaches suggests that IBE might be used

1462 to extend the confidence building approach. The many who do not share in the

1463 skepticism about IBE will be tempted to use the criterion of explanatory goodness in

1464 order to establish the approximate correctness of climate models. At the same time,

1465 we saw that the IBE approach does not help us to select which CMPs we are entitled

1466 to be confident in. We also saw that considering explanatory quality alone is not the

1467      appropriate way of assessing climate model performance. The critical approach

1468      provides not so much a way of assessing projection quality as one of systematizing

1469      such assessments and legitimizing its results. The legitimization it would provide,

1470      however, is problematic because of the lack of agreement about how to assess

1471      projection quality and because of the need to address the question of when consensus

1472      is a guide to truth.

1473      The possibilist and severe testing approaches are promising in that they

1474      propose specific non-probabilistic measures of CMP quality. The severe testing

1475      approach has the additional advantage that it provides a way of trying to get a handle

1476      on the effects of tuning on CMP accuracy. As we have noted, however, both

1477      possibilist and severe testing approaches face problems.

1478      Some of the difficulties that arise in assessing projection quality are

1479      difficulties that would arise irrespective of actual projection accuracy. Tuning may

1480      well not affect the ability of models reliably to generate some important class of

1481      projections. But our uncertainty about the very practice of tuning means that, even if

1482      the projections in question are accurate and reliably generated, we will find it difficult

1483      to decide whether they are accurate. Similarly, the non-linear nature of the climate

1484      system may well not adversely affect the accuracy of some class of projections. But

1485      our uncertainty about whether non-linearity is pertinent to the projections will mean

1486      that we will find it difficult to decide whether they are accurate. This is frustrating, but

1487      does not alter the predicament we find ourselves in with respect to developing

1488      adequate approaches to assessing projection quality.

1489
1490                             References
1491
1492 Abe, M., H. Shiogama, J. C. Hargreaves, J. D. Annan, T. Nozawa and S. Emori,
1493          Correlation between Inter-Model Similarities in Spatial Pattern for Present and
1494          Projected Future Mean Climate, *Sola, 5*, 133-136, 2009.

1495   Allen, M. R., J. Kettleborough and D. A. Stainforth, Model error in weather and
1496        climate forecasting, in Predictability of weather and climate, edited by T.
1497        Palmer and R. Hagedorn, pp. 391-427, Cambridge University Press,
1498        Cambridge, 2006.

1499   Allen, M. R. and W. J. Ingram, Constraints on future changes in climate and the
1500        hydrologic cycle, *Nature, 419*(6903), 224-232, 2002.

1501   Annan, J. D. and J. C. Hargreaves, Reliability of the CMIP3 ensemble, *Geophys. Res.*
1502        *Lett., 37*(2), L02703, 2010.

1503   Annan, J. D. and J. C. Hargreaves, On the generation and interpretation of
1504        probabilistic estimates of climate sensitivity, *Climatic Change, 104*(3-4), 423-
1505        436, 2011.

1506   Arneth, A., S. P. Harrison, S. Zaehle, K. Tsigaridis, S. Menon, P. J. Bartlein, J.
1507        Feichter, A. Korhola, M. Kulmala, D. O'Donnell, G. Schurgers, S. Sorvari and
1508        T. Vesala, Terrestrial biogeochemical feedbacks in the climate system, *Nature*
1509        *Geosci, 3*(8), 525-532, 2010.

1510   Barrett, J. and P. K. Stanford, Prediction, in The Philosophy of Science: An
1511        Encyclopedia, edited by J. Pfeifer and S. Sarkar, pp. 589-599, Routledge, New
1512        York, 2006.

1513   Bentley, M. J., The Antarctic palaeo record and its role in improving predictions of
1514        future Antarctic Ice Sheet change, *Journal of Quaternary Science, 25*(1), 5-18,
1515        2010.

1516   Boe, J., A. Hall and X. Qu, September sea-ice cover in the Arctic Ocean projected to
1517        vanish by 2100, *Nature Geosci, 2*(5), 341-343, 2009.

1518   Bokulich, A., How scientific models can explain, *Synthese, 180*(1), 33-45, 2011.

1519   Braconnot, P., B. Otto-Bliesner, S. Harrison, S. Joussaume, J. Y. Peterchmitt, A. be-
1520        Ouchi, M. Crucifix, E. Driesschaert, T. Fichefet, C. D. Hewitt, M. Kageyama,
1521        A. Kitoh, A. Laine, M. F. Loutre, O. Marti, U. Merkel, G. Ramstein, P.
1522        Valdes, S. L. Weber, Y. Yu and Y. Zhao, Results of PMIP2 coupled
1523        simulations of the Mid-Holocene and Last Glacial Maximum - Part 1:
1524        experiments and large-scale features, *Climate of the Past, 3*(2), 261-277, 2007.

1525   Branstator, G. and H. Teng, Two Limits of Initial-Value Decadal Predictability in a
1526        CGCM, *J. Climate, 23*(23), 6292-6311, 2010.

1527   Caseldine, C. J., C. Turney and A. J. Long, IPCC and palaeoclimate - an evolving
1528        story?, *Journal of Quaternary Science, 25*(1), 1-4, 2010.

1529   CCSP, Atmospheric Aerosol Properties and Climate Impacts. A Report by the U.S.
1530        Climate Change Science Program and the Subcommittee on Global Change
1531        Research. [Mian Chin, Ralph A. Kahn, and Stephen E. Schwartz (eds.)].
1532        National Aeronautics and Space Administration, Washington, D.C., USA.,
1533        2009.

Christensen, J. H., B. Hewitson, A. Busuioc, A. Chen, X. Gao, I. Held, R. Jones, R. K. Kolli, W. T. Kwon, R. Laprise, V. Magaña Rueda, L. Mearns, C. G. Menéndez, J. Räisanen, A. Rinke, A. Sarr and P. Whetton, Regional Climate Projections, in Climate Change 2007: The Physical Science Basis. Contribution of Working Group I to the Fourth Assessment Report of the Intergovernmental Panel on Climate Change, edited by S. Solomon, D. Qin, M. Manning, Z. Chen, M. Marquis, K. M. Averyt, M. Tignor and H. L. Miller, Cambridge University Press, Cambridge, United Kingdom and New York, NY, USA, 2007.

CMIP3, World Climate Research Programme's Coupled Model Intercomparison Project phase 3 multi-model dataset (on line), *http://www-pcmdi. llnl. gov/ipcc/about_ipcc. php,* 2007.

Crasnow, S. L., Can Science be Objective? Longino's Science as Social Knowledge, *Hypatia, 8,* 194-201, 1993.

Dall' Amico, M., P. Stott, A. Scaife, L. Gray, K. Rosenlof and A. Karpechko, Impact of stratospheric variability on tropospheric climate change, *Climate Dynamics, 34*(2), 399-417, 2010.

Dessler, A. E. and S. C. Sherwood, Atmospheric Science: A Matter of Humidity, *Science, 323*(5917), 1020-1021, 2009.

Deutsch, H., Real Possibility, *Nous, 24*(5), 751-755, 1990.

Earman, J. and C. Glymour, Einstein and Hilbert - 2 Months in History of General Relativity, *Archive for History of Exact Sciences, 19*(3), 291-308, 1978.

Frame, D. J., B. B. B. Booth, J. A. Kettleborough, D. A. Stainforth, J. M. Gregory, M. Collins and M. R. Allen, Constraining climate forecasts: The role of prior assumptions, *Geophysical Research Letters, 32*(9), L09702, 2005.

Frame, D. J., N. E. Faull, M. M. Joshi and M. R. Allen, Probabilistic climate forecasts and inductive problems, *Philosophical Transactions of the Royal Society A-Mathematical Physical and Engineering Sciences, 365*(1857), 1971-1992, 2007.

Freedman, K., Diversity and the Fate of Objectivity, *Social Epistemology, 23*(1), 45-56, 2009.

Furrer, R., R. Knutti, S. R. Sain, D. W. Nychka and G. A. Meehl, Spatial patterns of probabilistic temperature change projections from a multivariate Bayesian analysis, *Geophysical Research Letters, 34*(6), L06711, 2007a.

Furrer, R., S. R. Sain, D. Nychka and G. A. Meehl, Multivariate Bayesian analysis of atmosphere - Ocean general circulation models, *Environmental and Ecological Statistics, 14*(3), 249-266, 2007b.

Ghil, M., M. D. Chekroun and E. Simonnet, Climate dynamics and fluid mechanics: Natural variability and related uncertainties, *Physica D: Nonlinear Phenomena, 237*(14-17), 2111-2126, 2008.

1574    Gleckler, P. J., K. E. Taylor and C. Doutriaux, Performance metrics for climate
1575          models, *Journal of Geophysical Research-Atmospheres, 113*(6), D06104,
1576          2008.

1577    Halpern, J. Y., Reasoning about uncertainty, MIT Press, London and Cambridge
1578          Massachusetts, 2003.

1579    Hargreaves, J., Skill and uncertainty in climate models, *Wiley Interdisciplinary*
1580          *Reviews: Climate Change,* 556-564, 2010.

1581    Hawkins, E. and R. Sutton, The Potential to Narrow Uncertainty in Regional Climate
1582          Predictions, *Bull. Amer. Meteor. Soc., 90*(8), 1095-1107, 2009.

1583    Hegerl, G. C., T. J. Crowley, W. T. Hyde and D. J. Frame, Climate sensitivity
1584          constrained by temperature reconstructions over the past seven centuries,
1585          *Nature, 440*(7087), 1029-1032, 2006.

1586    Hegerl, G. C., F. W. Zwiers, P. Braconnot, N. P. Gillett, Y. Yuo, J. A. Marengo
1587          Orsini, N. Nicholls, J. E. Penner and P. A. Stott, Understanding and
1588          Attributing Climate Change, Cambridge University Press, Cambridge, United
1589          Kingdom and New York, NY, USA, 2007.

1590    Huber, M. and R. Caballero, The early Eocene equable climate problem revisited,
1591          *Clim. Past, 7*(2), 603-633, 2011.

1592    Hudson, R., What is Really at Issue with Novel Predictions?, *Synthese, 155*(1), 1-20,
1593          2007.

1594    Humphreys, P., Extending ourselves: computational science, empiricism and
1595          scientific method, Oxford University Press, Oxford, 2004.

1596    Hurrell, J., G. Meehl, D. Bader, T. Delworth, B. Kirtman and B. Wielicki, A Unified
1597          Modeling Approach to Climate System Prediction, *Bull. Amer. Meteor. Soc.,*
1598          *90*(12), 1819-1832, 2009.

1599    IPCC 1990, Climate Change: The IPCC Scientific Assessment, Cambridge University
1600          Press, Cambridge, UK and New York, 1990.

1601    IPCC 1996, Climate Change 1995: The Science of Climate Change. Contribution of
1602          Working Group I to the Second Assessment Report of the Intergovernmental
1603          Panel on Climate Change, Cambridge University Press, Cambridge, UK and
1604          New York, 1996.

1605    IPCC 2001, Climate change 2001: The scientific basis. Contribution of Working
1606          Group I to the Third Assessment Report of the Intergovernmental Panel on
1607          Climate Change, Cambridge University Press, Cambridge, 2001.

1608    IPCC 2005, Guidance Notes for Lead Authors of the IPCC Fourth Assessment
1609          Report on Addressing Uncertainties (on line), *http://www. ipcc.*
1610          *ch/meetings/ar4-workshops-express-meetings/uncertainty-guidance-note. pdf,*
1611          2005.

1612    IPCC 2007, Climate Change 2007: Synthesis Report. Contribution of Working
1613         Groups I, II and III to the Fourth Assessment Report of the Intergovernmental
1614         Panel on Climate Change [Core Writing Team, Pachauri, R.K and Reisinger,
1615         A. (eds.)]. IPCC, Geneva, Switzerland, 104 pp., 2007.

1616    IPCC 2010, Report of the Intergovernmental Panel on Climate Change Expert
1617         Meeting on Assessing and Combining Multi Model Climate Projections, IPCC
1618         Working Group I Technical Support Unit, University of Bern, Bern, 2010.

1619    Ishizaki, Y., T. Nakaegawa and I. Takayabu,  Comparison of Three Bayesian
1620         Approaches to Project Surface Air Temperature Changes over Japan Due to
1621         Global Warming, *Sola, 6*, 21-24, 2010.

1622    Jansen, E., J.Overpeck, K.R.Briffa, J.-C.Duplessy, F.Joos, V.Masson-Delmotte,
1623         D.Olago, B.Otto-Bliesner, W.R.Peltier, S.Rahmstorf, R.Ramesh, D.Raynaud,
1624         D.Rind, O.Solomina, R.Villalba and D.Zhang, Palaeoclimate, in *Climate
1625         Change* 2007: *The Physical Science Basis. Contribution of Working Group I
1626         to the Fourth Assessment Report of the Intergovernmental Panel on Climate
1627         Change*, edited by S. Solomon, D. Qin, M. Manning, Z. Chen, M. Marquis, K.
1628         M. Averyt, M. Tignor and H. L. Miller, Cambridge University Press,
1629         Cambridge, UK and New York, 2007.

1630    Jeffrey, R. C., The logic of decision, The University of Chicago Press, Chicago and
1631         London, 1990.

1632    Katzav, J.,  Should we assess climate model predictions in light of severe tests?, *Eos,
1633         Transactions American Geophysical Union, 92(23)*, 195, 2011.

1634    Katzav, J.,  Hybrid models, climate models and inference to the best explanation,
1635         *British Journal for the Philosophy of Science, doi: 10.1093/bjps/axs002*, 2012.

1636    Kirkby, J.,  Cosmic Rays and Climate, *Surveys in Geophysics, 28*(5), 333-375, 2007.

1637    Knutti, R.,  Should we believe model predictions of future climate change?,
1638         *Philosophical Transactions of the Royal Society A-Mathematical Physical and
1639         Engineering Sciences, 366*(1885), 4647-4664, 2008.

1640    Knutti, R.,  The end of model democracy?, *Climatic Change, 102*, 395-404, 2010.

1641    Knutti, R., R. Furrer, C. Tebaldi, J. Cermak and G. A. Meehl,  Challenges in
1642         Combining Projections from Multiple Climate Models, *J. Climate, 23*(10),
1643         2739-2758, 2010.

1644    Knutti, R., T. F. Stocker, F. Joos and G. K. Plattner,  Constraints on radiative forcing
1645         and future climate change from observations and climate model ensembles,
1646         *Nature, 416*(6882), 719-723, 2002.

1647    Knutti, R. and G. C. Hegerl,  The equilibrium sensitivity of the Earth's temperature to
1648         radiation changes, *Nature Geosci, 1*(11), 735-743, 2008.

1649    Laskar, J., P. Robutel, F. Joutel, M. Gastineau, A. C. M. Correia and B. Levrard,  A
1650            long-term numerical solution for the insolation quantities of the Earth,
1651            *Astronomy & Astrophysics, 428*(1), 261-285, 2004.

1652    Laudan, L.,  A Confutation of Convergent Realism, *Philosophy of Science, 48*(1), 19-
1653            49, 1981.

1654    Le Treut, H., R. Somerville, U. Cubasch, Y. Ding, C. Mauritzen, A. Mokssit, T.
1655            Peterson and M. Prather, Historical overview of climate change, in *Climate*
1656            *Change* 2007: *The Physical Science Basis. Contribution of Working Group I*
1657            *to the Fourth Assessment Report of the Intergovernmental Panel on Climate*
1658            *Change*, edited by S. Solomon, D. Qin, M. Manning, Z. Chen, M. Marquis, K.
1659            M. Averyt, M. Tignor and H. L. Miller, Cambridge University Press, New
1660            York, 2007.

1661    Leith, N. A. and R. E. Chandler,  A framework for interpreting climate model outputs,
1662            *Journal of the Royal Statistical Society Series C-Applied Statistics, 59*, 279-
1663            296, 2010.

1664    Lenhard, J. and E. Winsberg,  Holism, entrenchment, and the future of climate model
1665            pluralism, *Studies In History and Philosophy of Science Part B: Studies In*
1666            *History and Philosophy of Modern Physics, 41(3)*, 253-262, 2010.

1667    Lenton, T. M., H. Held, E. Kriegler, J. W. Hall, W. Lucht, S. Rahmstorf and H. J.
1668            Schellnhuber,  Tipping elements in the Earth's climate system, *Proceedings of*
1669            *the National Academy of Sciences of the United States of America, 105*(6),
1670            1786-1793, 2008.

1671    Lipton, P., *Inference to the best explanation*, Routledge, London and New York,
1672            2004.

1673    Longino, H. E., Science as Social Knowledge: Values and Objectivity in Scientific
1674            Inquiry, Princeton University Press, Princeton, 1990.

1675    Lorenz, E.,  Deterministic Nonperiodic Flow, *Journal of Atmospheric Science, 20*,
1676            103-141, 1963.

1677    Mahmood, R., R. A. Pielke, K. G. Hubbard, D. Niyogi, G. Bonan, P. Lawrence, R.
1678            McNider, C. McAlpine, A. Etter, S. Gameda, B. D. Qian, A. Carleton, A.
1679            Beltran-Przekurat, T. Chase, A. I. Quintanar, J. O. Adegoke, S.
1680            Vezhapparambu, G. Conner, S. Asefi, E. Sertel, D. R. Legates, Y. L. Wu, R.
1681            Hale, O. W. Frauenfeld, A. Watts, M. Shepherd, C. Mitra, V. G. Anantharaj,
1682            S. Fall, R. Lund, A. Trevino, P. Blanken, J. Y. Du, H. I. Chang, R. E. Leeper,
1683            U. S. Nair, S. Dobler, R. Deo and J. Syktus,  Impacts of Land Use/Land Cover
1684            Change on Climate and Future Research Priorities, *Bull. Amer. Meteor. Soc.,*
1685            *91*(1), 37-46, 2010.

1686    Mayo, D. G., Error and the Growth of Experimental Knowledge, The University of
1687            Chicago Press, Chicago and London, 1996.

1688    McWilliams, J. C.,  Irreducible imprecision in atmospheric and oceanic simulations,
1689            *Proceedings of the National Academy of Sciences, 104*(21), 8709-8713, 2007.

1690     Meehl, G. A., C. Covey, T. Delworth, M. Latif, B. McAvaney, J. F. B. Mitchell, R. J.
1691           Stouffer and K. E. Taylor, The WCRP CMIP3 multimodel dataset - A new era
1692           in climate change research, *Bull. Amer. Meteor. Soc., 88*(9), 1383-1394,
1693           2007a.

1694     Meehl, G. A., T. F. Stocker, W. D. Collins, P. Friedlingstein, A. P. Gaye, J. M.
1695           Gregory, A. Kitoh, R. Knutti, J. M. Murphy, A. Noda, S. C. B. Raper, I. G.
1696           Watterson, A. J. Weaver and Z.-C. Zhao, Global Climate Projections, in
1697           *Climate Change* 2007: *The Physical Science Basis. Contribution of Working*
1698           *Group I to the Fourth Assessment Report of the Intergovernmental Panel on*
1699           *Climate Change*, edited by S. Solomon, D. Qin, M. Manning, Z. Chen, M.
1700           Marquis, K. M. Averyt, M. Tignor and H. L. Miller, Cambridge University
1701           Press, Cambridge, United Kingdom and New York, NY, USA, 2007b.

1702     Miller, D., Out of error: further essays on critical rationalism, Ashgate Publishing
1703           Limited, Aldershot, 2005.

1704     Moss, R. and S. H. Schneider, Uncertainties in the IPCC TAR: Recommendations to
1705           lead authors for more consistent assessment and reporting, in Guidance papers
1706           on the cross cutting issues of the third assessment report of the IPCC.
1707           Technical report., pp. 33-51, World Meteorological Organization, Geneva,
1708           2000.

1709     Murphy, J. M., B. B. B. Booth, M. Collins, G. R. Harris, D. M. H. Sexton and M. J.
1710           Webb, A Methodology for Probabilistic Predictions of Regional Climate
1711           Change from Perturbed Physics Ensembles, *Philosophical Transactions:*
1712           *Mathematical, Physical and Engineering Sciences, 365*(1857), 1993-2028,
1713           2007.

1714     Murphy, J. M., D. M. H. Sexton, D. N. Barnett, G. S. Jones, M. J. Webb and M.
1715           Collins, Quantification of modelling uncertainties in a large ensemble of
1716           climate change simulations, *Nature, 430*(7001), 768-772, 2004.

1717     Musgrave, A., The ultimate argument for scientific realism, in Relativism and realism
1718           in science, edited by R. Nola, pp. 229-252, Kluwer, Dordrecht, 1988.

1719     North, J., An empirical approach to symmetry and probability, *Studies in History and*
1720           *Philosophy of Modern Physics, 41*(1), 27-40, 2010.

1721     Norton, J. D., Ignorance and indifference, *Philosophy of Science, 75*(1), 45-68, 2008.

1722     Norton, J. D., Challenges to Bayesian Confirmation Theory, in Philosophy of
1723           Statistics, vol. 7, edited by P. S. Bandyopadhyay and M. Forster, pp. 391-440,
1724           Elsevier, New York, 2011.

1725     Odenbaugh, J. and A. Alexandrova, Buyer beware: robustness analyses in economics
1726           and biology, *Biology and Philosophy, 26*(5), 757-771, 2011.

1727     Palmer, T. N., Predicting uncertainty in forecasts of weather and climate, *Reports on*
1728           *Progress in Physics, 63*(2), 71-116, 2000.

1729    Palmer, T. N., G. J. Shutts, R. Hagedorn, E. Doblas-Reyes, T. Jung and M.
1730          Leutbecher,  Representing model uncertainty in weather and climate
1731          prediction, *Annual Review of Earth and Planetary Sciences, 33*, 163-193,
1732          2005.

1733    Parker, W. S.,  Computer simulation through an error-statistical lens, *Synthese,*
1734          *163*(3), 371-384, 2008.

1735    Parker, W. S.,  Confirmation and Adequacy-for-Purpose in Climate Modelling,
1736          *Aristotelian Society Supplementary Volume, 83*(1), 233-249, 2009.

1737    Parker, W. S.,  When Climate Models Agree: The Significance of Robust Model
1738          Predictions, *Philosophy of Science, 78*(4), 579-600, 2011.

1739    Parker, W. S.,  Predicting weather and climate: Uncertainty, ensembles and
1740          probability, *Studies In History and Philosophy of Science Part B: Studies In*
1741          *History and Philosophy of Modern Physics, 41(3),* 263-272, 2010.

1742    Petersen, A. C.,  Simulating Nature: A Philosophical Study of Computer Simulation
1743          Uncertainties and their Role in Climate Science and Policy Advice, CRC
1744          Press, Boca Raton, FL, 2012.

1745    Piani, C., D. J. Frame, D. A. Stainforth and M. R. Allen,  Constraints on climate
1746          change from a multi-thousand member ensemble of simulations, *Geophys.*
1747          *Res. Lett., 32*(23), L23825, 2005.

1748    Pielke, R. A.,  Land use and climate change, *Science, 310*(5754), 1625-1626, 2005.

1749    Pirtle, Z., R. Meyer and A. Hamilton,  What does it mean when climate models agree?
1750          A case for assessing independence among general circulation models,
1751          *Environmental Science & Policy, 13*(5), 351-361, 2010.

1752    Pitman, A. J. and R. J. Stouffer,  Abrupt change in climate and climate models,
1753          *Hydrology and Earth System Sciences, 10*(6), 903-912, 2006.

1754    Popper, K. R., Realism and the aim of science, Routledge, London and New York,
1755          1983.

1756    Popper, K. R., Conjectures and refutations: the growth of scientific knowledge,
1757          Routledge, London and New York, 2002.

1758    Popper, K. R., The logic of scientific discovery, Routledge, London and New York,
1759          2005.

1760    Psillos, S., Scientific realism: how science tracks truth, Routledge, London and New
1761          York, 1999.

1762    Raisanen, J. and T. N. Palmer,  A probability and decision-model analysis of a
1763          multimodel ensemble of climate change simulations, *J. Climate, 14*(15), 3212-
1764          3226, 2001.

Randall, D. A., R. A. Wood, S. Bony, R. Coleman, T. Fichefet, J. Fyfe, V. Kattsof, A. Pitman, J. Shukla, J. Srinivasan, R. J. Stouffer, A. Sumi and K. E. Taylor, Climate Models and Their Evaluation, in *Climate Change* 2007: *The Physical Science Basis. Contribution of Working Group I to the Fourth Assessment Report of the Intergovernmental Panel on Climate Change*, edited by S. Solomon, D. Qin, M. Manning, Z. Chen, M. Marquis, K. M. Averyt, M. Tignor and H. L. Miller, Cambridge University Press, Cambridge, United Kingdom and New York, NY, USA, 2007.

Randall, D., M. Khairoutdinov, A. Arakawa and W. Grabowski, Breaking the Cloud Parameterization Deadlock, *Bull. Amer. Meteor. Soc., 84*(11), 1547-1564, 2003.

Rial, J. A., R. A. Pielke, M. Beniston, M. Claussen, J. Canadell, P. Cox, H. Held, N. De Noblet-Ducoudre, R. Prinn, J. F. Reynolds and J. D. Salas, Nonlinearities, feedbacks and critical thresholds within the Earth's climate system, *Climatic Change, 65*(1-2), 11-38, 2004.

Rougier, J., Probabilistic inference for future climate using an ensemble of climate model evaluations, *Climatic Change, 81*(3), 247-264, 2007.

Rowbottom, D., Popper's Critical Rationalism: A Philosophical Investigation, Routledge, New York, 2011.

Saatsi, J., On the Pessimistic Induction and Two Fallacies, *Philosophy of Science, 72*(5), 1088-1098, 2005.

Salmon, W. C., Rational Prediction, *British Journal for the Philosophy of Science, 32*(2), 115-125, 1981.

Sanderson, B. M., R. Knutti, T. Aina, C. Christensen, N. Faull, D. J. Frame, W. J. Ingram, C. Piani, D. A. Stainforth, D. A. Stone and M. R. Allen, Constraints on Model Response to Greenhouse Gas Forcing and the Role of Subgrid-Scale Processes, *J. Climate, 21*(11), 2384-2400, 2008.

Schmidt, G. A., Enhancing the relevance of palaeoclimate model/data comparisons for assessments of future climate change, *Journal of Quaternary Science, 25*(1), 79-87, 2010.

Schmittner, A., M. Latif and B. Schneider, Model projections of the North Atlantic thermohaline circulation for the 21st century assessed by observations, *Geophysical Research Letters, 32*(23), L23710, 2005.

Schwartz, S., R. J. Charlson and H. Rodhe, Quantifying climate change - too rosy a picture?, *Nature Reports Climate Change, 2*, 23-24, 2007.

Semenov, V. A., M. Latif, D. Dommenget, N. S. Keenlyside, A. Strehz, T. Martin and W. Park, The Impact of North AtlanticΓÇôArctic Multidecadal Variability on Northern Hemisphere Surface Air Temperature, *J. Climate, 23*(21), 5668-5677, 2010.

1804    Shackley, S., Epistemic lifestyles in climate change modeling, in Changing the
1805            atmosphere, edited by C. A. Miller and Edwards P.N., pp. 109-133, The MIT
1806            Press, Cambridge Mass. and London, 1997.

1807    Shukla, J., T. DelSole, M. Fennessy, J. Kinter and D. Paolino,  Climate model fidelity
1808            and projections of climate change, *Geophys. Res. Lett., 33*(7), L07702, 2006.

1809    Siddall, M., A. be-Ouchi, M. Andersen, F. Antonioli, J. Bamber, E. Bard, J. Clark, P.
1810            Clark, P. Deschamps, A. Dutton, M. Elliot, C. Gallup, N. Gomez, J. Gregory,
1811            P. Huybers, K. Kawarnura, M. Kelly, K. Lambeck, T. Lowell, J. Milrovica, B.
1812            Otto-Bliesner, D. Richards, J. Stanford, C. Stirling, T. Stocker, A. Thomas, B.
1813            Thompson, T. Tornqvist, N. V. Riveiros, C. Waelbroeck, Y. Yokoyama and S.
1814            Y. Yu,  The sea-level conundrum: case studies from palaeo-archives, *Journal
1815            of Quaternary Science, 25*(1), 19-25, 2010.

1816    Smith, L. A., Predictability past, predictability present, in Predictability of Weather
1817            and Climate, edited by T. Palmer and R. Hagedorn, pp. 217-250, Cambridge
1818            University Press, Cambridge, 2006.

1819    Snyder, C. W.,  The value of paleoclimate research in our changing climate, *Climatic
1820            Change, 100*(3-4), 407-418, 2010.

1821    Sober, E., Bayesianism — its Scope and Limits, in Bayes' Theorem, edited by R.
1822            Swinburne, pp. 21-38, Oxford University Press, Oxford, 2002.

1823    Spicer, R. A., A. Ahlberg, A. B. Herman, C. C. Hofmann, M. Raikevich, P. J. Valdes
1824            and P. J. Markwick,  The Late Cretaceous continental interior of Siberia: A
1825            challenge for climate models, *Earth and Planetary Science Letters, 267*(1-2),
1826            228-235, 2008.

1827    Stainforth, D. A., T. Aina, C. Christensen, M. Collins, N. Faull, D. J. Frame, J. A.
1828            Kettleborough, S. Knight, A. Martin, J. M. Murphy, C. Piani, D. Sexton, L. A.
1829            Smith, R. A. Spicer, A. J. Thorpe and M. R. Allen,  Uncertainty in predictions
1830            of the climate response to rising levels of greenhouse gases, *Nature,
1831            433*(7024), 403-406, 2005.

1832    Stainforth, D. A., M. R. Allen, E. R. Tredger and L. A. Smith,  Confidence,
1833            uncertainty and decision-support relevance in climate predictions,
1834            *Philosophical Transactions of the Royal Society A-Mathematical Physical and
1835            Engineering Sciences, 365*(1857), 2145-2161, 2007a.

1836    Stainforth, D. A., T. E. Downing, R. Washington, A. Lopez and M. New,  Issues in
1837            the interpretation of climate model ensembles to inform decisions,
1838            *Philosophical Transactions of the Royal Society A-Mathematical Physical and
1839            Engineering Sciences, 365*(1857), 2163-2177, 2007b.

1840    Stanford, P. K., Exceeding Our Grasp: Science, History, and the Problem of
1841            Unconceived Alternatives, Oxford University Press, New York, 2006.

1842    Stott, P. A., J. F. B. Mitchell, M. R. Allen, T. L. Delworth, J. M. Gregory, G. A.
1843            Meehl and B. D. Santer,  Observational Constraints on Past Attributable

1844    Warming and Predictions of Future Global Warming, *J. Climate, 19*(13),
1845         3055-3069, 2006.

1846    Strevens, M., Probability and chance, in Macmillan Encyclopedia of Philosophy, vol.
1847         8, edited by D. M. Borchert, pp. 24-40, MacMillan Reference USA. Thomson
1848         Gale., New York, 2006a.

1849    Strevens, M., The Bayesian Approach to the Philosophy of Science, in Encyclopedia
1850         of Philosophy, edited by D. M. Borchert, Macmillan Reference, Detroit,
1851         2006b.

1852    Swanson, K. L., G. Sugihara and A. A. Tsonis,  Long-term natural variability and
1853         20th century climate change, *Proceedings of the National Academy of
1854         Sciences, 106*(38), 16120-16123, 2009.

1855    Tebaldi, C. and R. Knutti,  The use of the multi-model ensemble in probabilistic
1856         climate projections, *Philosophical Transactions of the Royal Society A-
1857         Mathematical Physical and Engineering Sciences, 365*(1857), 2053-2075,
1858         2007.

1859    Tebaldi, C., R. L. Smith, D. Nychka and L. O. Mearns,  Quantifying uncertainty in
1860         projections of regional climate change: A Bayesian approach to the analysis of
1861         multimodel ensembles, *J. Climate, 18*(10), 1524-1540, 2005.

1862    Thagard, P. R.,  Best Explanation - Criteria for Theory Choice, *Journal of Philosophy,
1863         75*(2), 76-92, 1978.

1864    Valdes, P.,  Built for stability, *Nature Geosci, 4*(7), 414-416, 2011.

1865    Vineberg, S.,  Dutch Book Arguments (on line), *the Stanford Encyclopedia of
1866         Philosophy, http://plato. stanford. edu/archives/sum2011/entries/dutch-book/,*
1867         2011.

1868    Weisberg, M.,  Robustness analysis, *Philosophy of Science, 73*(5), 730-742, 2006.

1869    Wigley, T. M. L.,  Climate-Change - Outlook Becoming Hazier, *Nature, 369*(6483),
1870         709-710, 1994.

1871    Winsberg, E.,  Models of success versus the success of models: Reliability without
1872         truth, *Synthese, 152*(1), 1-19, 2006.

1873    Worrall, J., Error, tests and theory confirmation, in Error and Inference: Recent
1874         Exchanges on Experimental Reasoning, Reliability, and the Objectivity and
1875         Rationality of Science, edited by D. G. Mayo and A. Spanos, pp. 125-154,
1876         Cambridge University Press, New York, 2010.

1877    Wu, Z., N. Huang, J. Wallace, B. Smoliak and X. Chen,  On the time-varying trend in
1878         global-mean surface temperature, *Climate Dynamics,* 1-15, 2011.

1879    Wunsch, C.,  Towards understanding the Paleocean, *Quaternary Science Reviews,
1880         29*(17-18), 1960-1967, 2010.

1881    Wyatt, M., S. Kravtsov and A. Tsonis,  Atlantic Multidecadal Oscillation and
1882            Northern Hemisphere's climate variability, *Climate Dynamics,* 1-21, 2011.
1883

**Captions**

*Fig. 1 Temperature changes relative to the corresponding average for 1901-1950 (°C) from decade to decade from 1906 to 2005 over the entire globe, global land area and the global ocean. The black line indicates observed temperature change, while the colored bands show the combined range covered by 90% of recent model simulations. Red indicates simulations that include natural and human factors, while blue indicates simulations that include only natural factors. Dashed black lines indicate decades and continental regions for which there are substantially fewer observations. Adapted from Hegerl et al., FAQ9.2, Fig. 1 (2007, p. 703).*

*Fig. 2 Projected 21st century global mean temperatures changes for various greenhouse gas emission scenarios.  Solid lines are multi-model global averages of surface warming for scenarios A2, A1B and B1, shown as continuations of the 20th-century simulations. These projections also take into account emissions of short-lived GHGs and aerosols. The pink line is not a scenario, but is for Atmosphere-Ocean General Circulation Model (AOGCM) simulations where atmospheric concentrations are held constant at year 2000 values. The bars at the right of the figure indicate the best estimate (solid line within each bar) and the likely range assessed for the six SRES marker scenarios at 2090-2099. All temperatures are relative to the period 1980-1999. Adapted from the  Synthesis Report for IPCC AR4, Fig. 3.2 (2007, p. 7).*

*Fig. 3 Evolution of the MOC at 30°N in simulations with the suite of comprehensive coupled climate models from 1850 to 2100 using 20th Century Climate in Coupled Models (20C3M) simulations for 1850 to 1999 and the SRES A1B emissions scenario for 1999 to 2100. Some of the models continue the integration to year 2200 with the forcing held constant at the values of year 2100. Observationally based estimates of late-20th century MOC are shown as vertical bars on the left. Adapted from Meehl et al., Fig. 10.15 (2007b, p. 773),  who build on Schmittner et al. (2005).*

*Fig. 4. Root-mean-square (RMS) error of 1980–99 surface temperature (averaged over space, relative to the 40-year reanalysis of the European Centre of Medium range Weather Forecast) shown as a function of the number of models included in the model average. Panel (a) shows the December-January-February period (DJF), panel (b) the June-July-August (JJA) period. Red dashed lines indicate the range covered by randomly sampling the models for the subset; the red solid line indicates the average. The RMS error converges to a constant value that is more than half of the initial value for one model. The black dashed line is a theoretical RMS. If the model biases were independent, then the RMS error for a large sample of models should decrease with the square root of the number of models (dotted). The blue line results if the models are sorted by how well they agree with DJF and JJA observations combined, and it indicates that the average of a few good models outperforms an average of more models with poorer performance. Adapted from Knutti et al., Figs 3(c) and 3(d) (2010, p. 2744).*

*Fig. 5 Constraints on the radiative forcing from the observed atmospheric and oceanic warming. Probability density functions (PDF) for the total (anthropogenic and natural) radiative forcing (a–c) and the indirect aerosol forcing (d–f) in the year*
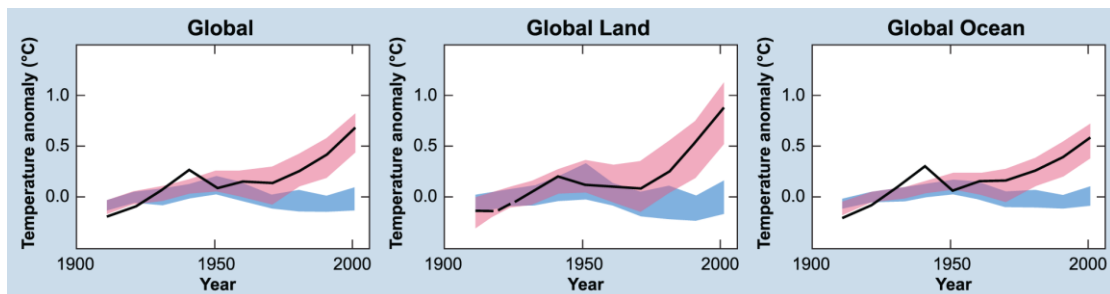
1931 *2000 are based on 25,000 Monte Carlo simulations. The initially assumed PDFs are*
1932 *given in a and d. The requirement that the model matches the temperature*
1933 *observations strongly narrows the PDFs (b and e). If in addition the climate*
1934 *sensitivity is restricted to the range adopted by the IPCC (1.5–4.5 K), the PDFs in c*
1935 *and f are obtained. Adapted from Knutti et al., Fig. 2 (2002, p. 720).*
1936
1937 *Fig. 6. The response to parameter perturbations: the frequency distribution of*
1938 *simulated climate sensitivity using all model versions (black), all model versions*
1939 *except those with perturbations to the cloud-to-rain conversion threshold (red), and*
1940 *all model versions except those with perturbations to the entrainment coefficient*
1941 *(blue). Adapted from Stainforth et al, Fig. 2(a) (2005, p. 404).*
1942
1943 **Figures**



1944
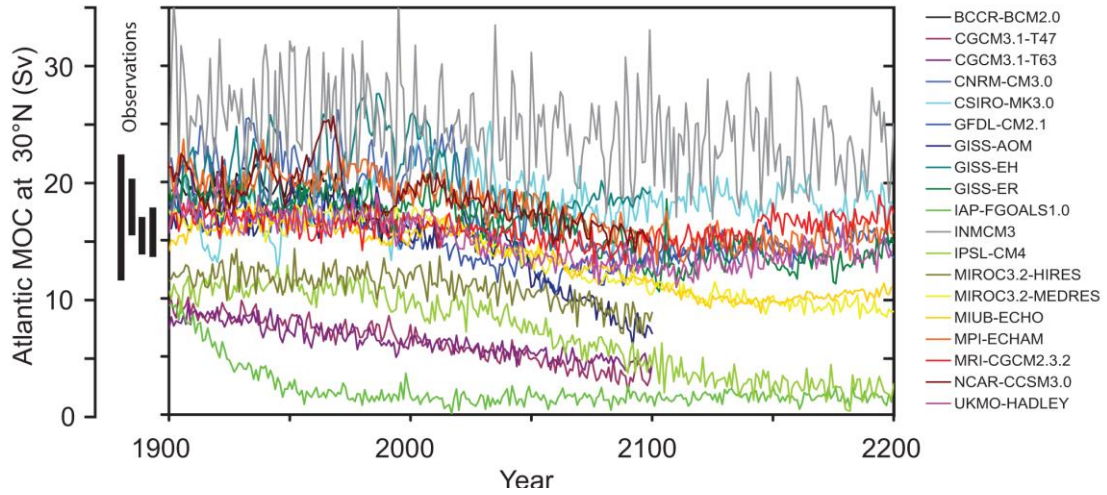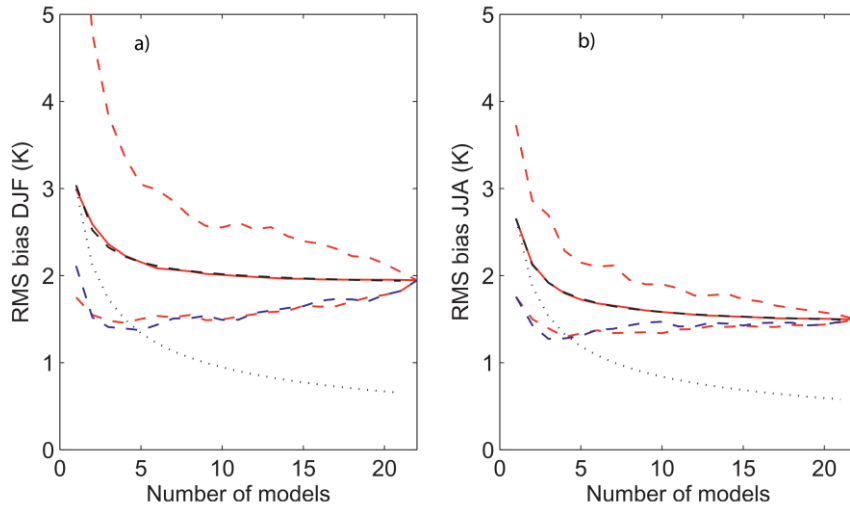1945 *Fig. 1*
1946



1947

1948 *Fig. 2*
1949

1950
1951    *Fig. 3*
1952



1953
1954    *Fig. 4*
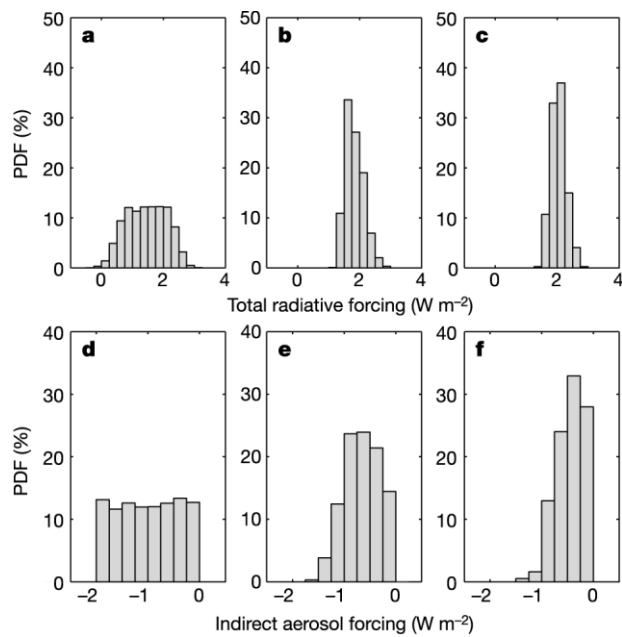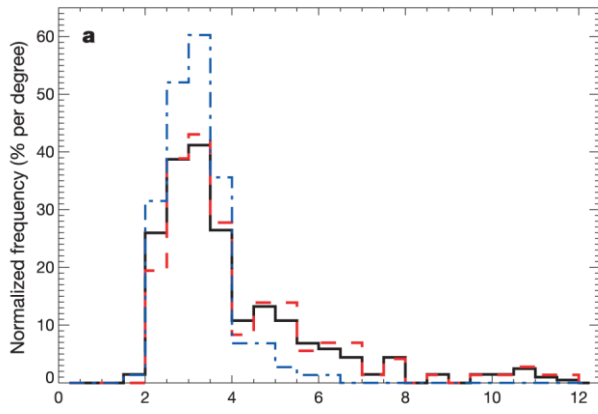1955



1956
1957    *Fig. 5*

1958



1959
1960    *Fig. 6*