



Evaluation of biospheric components in Earth system models using modern and palaeo-observations: the state-of-the-art

A. M. Foley^{1,*}, D. Dalmonech², A. D. Friend¹, F. Aires³, A. T. Archibald⁴, P. Bartlein⁵, L. Bopp⁶, J. Chappellaz⁷, P. Cox⁸, N. R. Edwards⁹, G. Feulner¹⁰, P. Friedlingstein⁸, S. P. Harrison¹¹, P. O. Hopcroft¹², C. D. Jones¹³, J. Kolassa³, J. G. Levine^{14,**}, I. C. Prentice¹⁵, J. Pyle⁴, N. Vázquez Riveiros¹⁶, E. W. Wolff^{14,***}, and S. Zaehle²

¹Department of Geography, University of Cambridge, Cambridge, UK

²Biogeochemical Integration Department, Max Planck Institute for Biogeochemistry, Jena, Germany

³Estellus, Paris, France

⁴Centre for Atmospheric Science, University of Cambridge, Cambridge, UK

⁵Department of Geography, University of Oregon, Eugene, Oregon, USA

⁶Laboratoire des Sciences du Climat et de l'Environnement, Gif sur Yvette, France

⁷UJF – Grenoble I and CNRS Laboratoire de Glaciologie et Géophysique de l'Environnement, Grenoble, France

⁸College of Engineering, Mathematics and Physical Sciences, University of Exeter, Exeter, UK

⁹Environment, Earth and Ecosystems, The Open University, Milton Keynes, UK

¹⁰Potsdam Institute for Climate Impact Research, Potsdam, Germany

¹¹Department of Biological Sciences, Macquarie University, Sydney, Australia and Geography and Environmental Sciences, School of Human and Environmental Sciences, Reading University, Reading, UK

¹²BRIDGE, School of Geographical Science, University of Bristol, Bristol, UK

¹³Met Office Hadley Centre, Exeter, UK

¹⁴British Antarctic Survey, Cambridge, UK

¹⁵AXA Chair of Biosphere and Climate Impacts, Department of Life Sciences and Grantham Institute for Climate Change, Imperial College, Silwood Park, UK and Department of Biological Sciences, Macquarie University, Sydney, Australia

¹⁶Godwin Laboratory for Palaeoclimate Research, Department of Earth Sciences, University of Cambridge, Cambridge, UK

* now at: Cambridge Centre for Climate Change Mitigation Research, Department of Land Economy, University of Cambridge, Cambridge, UK

** now at: School of Geography, Earth and Environmental Sciences, University of Birmingham, Birmingham, UK

*** now at: Department of Earth Sciences, University of Cambridge, Cambridge, UK

Correspondence to: A. M. Foley (amf62@cam.ac.uk)

Received: 3 June 2013 – Published in Biogeosciences Discuss.: 4 July 2013

Revised: 1 November 2013 – Accepted: 4 November 2013 – Published: 16 December 2013

Abstract. Earth system models (ESMs) are increasing in complexity by incorporating more processes than their predecessors, making them potentially important tools for studying the evolution of climate and associated biogeochemical cycles. However, their coupled behaviour has only recently been examined in any detail, and has yielded a very wide range of outcomes. For example, coupled climate–carbon cycle models that represent land-use change simulate total land carbon stores at 2100 that vary by as much as 600 Pg C, given the same emissions scenario. This large

uncertainty is associated with differences in how key processes are simulated in different models, and illustrates the necessity of determining which models are most realistic using rigorous methods of model evaluation. Here we assess the state-of-the-art in evaluation of ESMs, with a particular emphasis on the simulation of the carbon cycle and associated biospheric processes. We examine some of the new advances and remaining uncertainties relating to (i) modern and palaeodata and (ii) metrics for evaluation. We note that the practice of averaging results from many models is unreliable

and no substitute for proper evaluation of individual models. We discuss a range of strategies, such as the inclusion of pre-calibration, combined process- and system-level evaluation, and the use of emergent constraints, that can contribute to the development of more robust evaluation schemes. An increasingly data-rich environment offers more opportunities for model evaluation, but also presents a challenge. Improved knowledge of data uncertainties is still necessary to move the field of ESM evaluation away from a “beauty contest” towards the development of useful constraints on model outcomes.

1 Introduction

Earth system models (ESMs), which use sets of equations to represent atmospheric, oceanic, cryospheric, and biospheric processes and interactions (Claussen et al., 2002; Le Treut et al., 2007; Lohmann et al., 2008), are intended as tools for the study of the Earth system. The current generation of ESMs are substantially more complex than their predecessors in terms of land and ocean biogeochemistry, and can also account for land cover change, which is an important driver of the climate system through both biophysical and biogeochemical feedbacks. Yet their coupled behaviour has only recently begun to be explored.

The carbon cycle is a central feature of current ESMs, and the representation and quantification of climate-carbon cycle feedbacks involving the biosphere has been a primary goal of recent ESM development. ESM results submitted to the Coupled Model Intercomparison Project Phase 5 (CMIP5) simulate total land carbon stores in 2100 that vary by as much as 600 Pg C across models with the ability to represent land-use change, even when forced with the same anthropogenic emissions (Jones et al., 2013). This indicates that there are large uncertainties associated with how carbon cycle processes are represented in different models. In addition to these uncertainties in the biogeochemical climate-vegetation feedbacks, there are considerable uncertainties in the biogeophysical feedbacks (Willeit et al., 2013).

Robust evaluation of a model’s ability to simulate key carbon cycle processes is therefore a critical component of efforts to model future climate-carbon cycle dynamics. Robust evaluation establishes the confidence which can be placed on a given model’s projection of future behaviours and states of the system. However, evaluation is complicated by the fact that ESMs differ in their level of complexity. To take the example of land cover, while some models only account for biophysical effects (e.g. related to changes in surface albedo), some ESMs also account for biogeochemical effects (e.g. principally a change in carbon storage following land conversion). Another example is the representation of nutrient cycles. Not all ESMs include nutrient cycles. Current model projections that do include the coupling between ter-

restrial carbon and nitrogen (and in some cases phosphorus) cycles suggest that taking nutrient limitations into account attenuates possible future carbon cycle responses. This is because soil nitrogen tends to limit the ability of plants to respond positively to increases in atmospheric CO₂, reducing CO₂ fertilisation, and, conversely, tends to limit ecosystem carbon losses with temperature increases, as these also increase rates of nitrogen mineralisation. The reduction in CO₂ fertilisation is found to dominate, leading to a stronger accumulation of CO₂ in the atmosphere by the end of the 21st century than is projected by carbon cycle models that do not include nutrient feedbacks (Sokolov et al., 2008; Thornton et al., 2009; Zaehle et al., 2010).

Evaluation studies in climate modelling have highlighted how choice of methodology can significantly impact the conclusions reached concerning model skill (e.g. Radic and Clarke, 2011; Foley et al., 2013). Several studies have found that the mean of an ensemble of models outperforms all or most single models of that ensemble (e.g. Evans, 2008; Pincus et al., 2008). However, Schaller et al. (2011) demonstrated that although the multi-model mean outperforms individual models when the ability to reproduce global fields of climate variables is evaluated, it does not consistently outperform the individual models when the ability to simulate regional climatic features is evaluated. This highlights the need for robust assessments of model skill. Model evaluations which use inappropriate metrics or fail to consider key aspects of the system have the potential to lead to overconfidence in model projections. In particular, the averaging of results from different models is not an adequate substitute for proper evaluation of each model in turn.

Developing robust approaches to model evaluation, that is, approaches which reduce the data- and metric-dependency of statements about model skill, is challenging for reasons that are not exclusive to carbon cycle modelling but applicable across all aspects of Earth system modelling. Data sets may lack uncertainty estimates, significantly reducing their usefulness for model evaluation. Critical analysis may be required to reconcile differences between data sets intended to describe similar phenomena, such as temperature reconstructions based on different indicators (Mann et al., 2009). Furthermore, there are many metrics in use in model evaluation and often, the rationale for applying a specific metric is unclear. This paper considers these issues, along with strategies for improvement.

Overview of this paper

Knowledge of the system under observation is essential for the assessment of model performance (Oreskes et al., 1994). We therefore begin with a discussion of some challenges associated with the use of *modern and palaeodata* in model evaluation. Data validity (Sargent, 2010) is a crucial aspect. Key issues include uncertainties associated with our understanding of the changes captured in each type of record,

mismatches between available data and what is required for evaluation, and the challenges of using data collected at a specific spatial or temporal scale to develop larger-scale tests of model behaviour.

Next, we consider *metrics for model evaluation*. Metrics are simple formulae or mathematical procedures that measure the similarity or difference between two data sets. Whether using classical metrics (such as root mean square error, correlation, or model efficiency), or advanced analytical techniques (such as artificial neural networks), to compare models with data and quantify model skill, it is necessary to be aware of the statistical properties of metrics, as well as the properties of the model variables under consideration and the limitations of the evaluation data sets. Otherwise, there is a strong potential to draw false conclusions concerning model skill. Recent attempts to provide a benchmarking framework for land surface model evaluation indicate a move toward setting community-accepted standards (Randerson et al., 2009; Luo et al., 2012; Kelley et al., 2012). However, different levels of complexity in ESMs, different parameterisation procedures and modelling approaches, the validity of data, and an unavoidable level of subjectivity complicate the task of identifying universally applicable procedures.

Finally, *recommendations for more robust evaluation* are discussed. We note that evaluation can be process-based (“bottom-up”) or system-level (“top-down”) (Fig. 1). Evaluation can utilise pre-calibration, and/or emergent constraints across multiple models. A combination of approaches can increase our understanding of a model’s ability to simulate processes across multiple temporal and spatial scales.

Consideration will also be given to how key questions arising in the paper could potentially be resolved through *coordinated research activities*.

2 The role of data sets in ESM evaluation

ESMs aim to simulate a highly complex system. Non-linearities in the system imply that even a small change in one of the components might unexpectedly influence another component (Roe and Baker, 2007). As such, robust model evaluation is critical to assist in understanding the behaviour of ESMs and the limitations of what we can and cannot represent quantitatively. The development of such approaches to model evaluation requires consideration of many different data types.

Modern and palaeodata are both used for model evaluation, although each kind of data has advantages and limitations (Table 1). Experimental data provide benchmarks for a range of carbon cycle-relevant processes (e.g. physiologically-based responses of ecosystems to warming and CO₂ increase) that cannot be tested in other ways. However, for processes that are biome-specific, the limited geographical scope of the relatively few existing records is problematic. Data sets also exist with more global

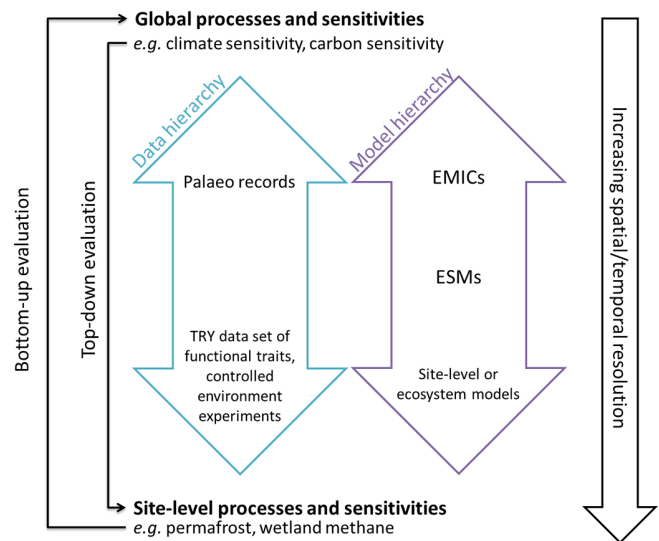


Fig. 1. Conceptual diagram of hierarchical approach to model evaluation on different spatial and temporal scales.

coverage, documenting changes in the recent past (last 30–50 yr), but an inherent limitation of these data sets is that they sample the carbon cycle response to a limited range of variation in atmospheric CO₂ concentration and climate.

Palaeoclimate evaluation is an important test of how well ESMs reproduce climate changes (e.g. Braconnot et al., 2012). The past does not provide direct analogues for the future, but does offer the opportunity to examine climate changes that are as large as those anticipated during the 21st century, and to evaluate climate system feedbacks with response times longer than the instrumental period (e.g. cryosphere, ocean circulation, some components of the carbon cycle).

2.1 Modern data sets

Evaluation analysis can benefit from modern data sets, to test and constrain components within ESMs in a hierarchical approach (Leffelaar, 1990; Wu and David, 2002). Recent initiatives in land and ocean model evaluation and benchmarking (land: Randerson et al., 2009; Luo et al., 2012; Kelley et al., 2012; Dalmonech and Zaehle, 2013; ocean: Najjaret al., 2007; Friedrichs et al., 2009; Stow et al., 2009; Bopp et al., 2013) give examples of suitable modern data sets for model evaluation and their use in diagnosing model inconsistencies with respect to behaviour of the carbon cycle. These include instrumental data, such as direct measurements of CO₂, and CH₄ spanning the last 30–50 yr, measurements from carbon flux monitoring networks, and satellite-based data of various kinds (Table 1).

Due to their detailed spatial coverage and high temporal resolution, satellite data sets offer the potential to explore the representation of processes in models in detail, and to reveal

Table 1. Summary of key data types for evaluation.

Type of data	Description	Examples	Advantages	Limitations
Modern last 30–50 yr	In situ instrumental data	Atmospheric CO ₂ , CH ₄ .	Direct observations of key variables, known uncertainties.	Local observations.
	Experimental data	Controlled environments (e.g. Phytotrons and glasshouses).	Provides new situations against which to test model behaviour, such as representing ecosystem-scale responses to combined environmental drivers.	Large-scale field experiments generally do not provide information across all biomes.
		Field experiments (e.g. Free Air Carbon dioxide Enrichment – FACE).		Interpretation of experiments may be ambiguous.
	Model-derived type I	Satellite-based data.	Excellent spatial and/or temporal resolution.	Lack full data-model independency, as data is model-derived (radiative transfer model converts radiation measurements into the parameter of interest). Inconsistent documentation of errors and uncertainties.
Model-derived type II	C-fluxes up-scaled data (e.g. MTE-MPI data set).	“Data-model” conceptual correspondence.	Lack fully data-model independency, as data is model-derived. Inconsistent documentation of errors and uncertainties.	
Palaeo	Reconstructions based on interpretation of biological or geochemical records.	Tree-ring data sets. Pollen and plant macrofossil data (e.g. BIOME 6000).	Tests ability to capture behaviour of the system outside modern range. Signal is large compared to noise.	Site-specific records (except for long-lived greenhouse gases), synthesis required to produce global estimates.
	Measurements of concentrations and isotopic ratios from ice cores.	Ice cores, e.g. Law Dome, EPICA.		Variable temporal resolution necessitates appropriate selection of data to address, e.g. rapid changes. Inconsistent documentation of errors and uncertainties.

compensating errors in ESMs. One of the main concerns is the lack of full consistency between what we can observe with different satellite sensors (e.g. top of the atmosphere reflectance) and what models actually simulate (e.g. net primary productivity). The lack of full independence between the data and the model is also an issue that often affects such comparisons. Satellite data are typically model-processed (type 1, Table 1), with some sort of model used to transform the direct measurements of the satellite into other parameters of interest. If, for example, a radiative transfer model is used to estimate the atmospheric or surface state from measured radiances, then there will likely be similarities between the functions used for the retrieval and those used in a climate model. This is not a major problem if the data are used in an

informed way, and indeed it presents opportunities (e.g. the estimated surface variable can be compared with a modelled variable without the model radiative transfer functions being involved). Statistical and change detection retrievals rely not on physical models but on statistical links between variables or on a modulation of the satellite signal. These two types of retrievals sometimes use model data for calibration but are otherwise independent of models. Statistical models in particular are not only useful to evaluate specific parameters in a model, but can also be used to perform process-based evaluations.

Uncertainty estimates are not always provided or propagated during the retrieval process. Nevertheless, modern data sets are a very rich data source with a number of useful

applications. For example, robust spatial and temporal information emerging from data can be used to rule out unreasonable simulations and diagnose model weaknesses. Satellite-based data sets of vegetation activity depict ecosystem response to climate variability at seasonal and interannual time scales and return patterns of forced variability that can be useful for model evaluation (Beck et al., 2011; Dahlke et al., 2012), even if bias within the data set is greater than data-model differences (e.g. Fig. 2).

Ecosystem observations, such as eddy covariance measurements of CO_2 and latent heat exchanges between the atmosphere and land, and ecosystem manipulation studies, such as drought treatments and free air CO_2 enrichment (FACE) experiments, provide a unique source of information to evaluate process formulations in the land component of ESMs (Friend et al., 2007; Bonan et al., 2012; de Kauwe et al., 2013). Manipulation experiments (e.g. FACE experiments: Nowak et al., 2004; Ainsworth and Long, 2005; Norby and Zak, 2011) are a particularly powerful test of key processes in ESMs and their constituent components, as shown by Zaehle et al. (2010) in relation to C-N cycle coupling, and de Kauwe et al. (2013) for carbon-water cycling. It should be expected that models would be able to reproduce experimental results involving manipulations of global change drivers such as CO_2 , temperature, rainfall, and N addition.

The application of such data for the evaluation of ESMs is challenging because of the limited spatial representativeness of the observations, resulting from the lack of any coherent global strategy for the placement of flux towers or experimental sites, and the high costs of running these facilities. Upscaling monitoring data using data-mining techniques and ancillary data, such as remote sensing and climate data, provides one possible means to bridge the gap between the spatial scale of observation and ESMs (Jung et al., 2011). However, this can be at the cost of introducing model assumptions and uncertainties that are difficult to quantify. Furthermore, such upscaling is near impossible for ecosystem manipulation experiments, as they are so scarce, and rarely performed following a comparable protocol. More and better-coordinated manipulation studies are needed to better constrain ESM prediction (Batterman and Larsen, 2011; Vicca et al., 2012). Hickler et al. (2008), for example, showed that the LPJ-GUESS model produced quantitatively realistic net primary production (NPP) enhancement due to CO_2 elevation in temperate forests, but also showed greatly different responses in boreal and tropical forests, for which no adequate manipulation studies exist. These predictions remain to be tested.

The interpretation of experiments is not unambiguous because it is seldom that just one external variable can be altered at a time. To give just one recent example, Bauerle et al. (2012) showed that the widely observed decline of Rubisco capacity ($V_{c,\max}$) in leaves during the transition from summer to autumn could be abated by supplementary light-

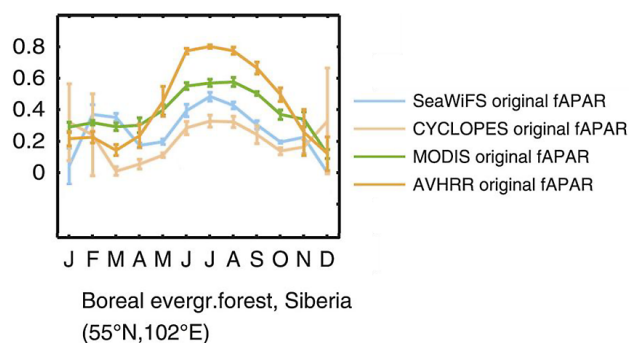


Fig. 2. Original fAPAR time series from a selected region (after Dahlke et al., 2012). © American Meteorological Society. Used with permission.

ing designed to maintain the summer photoperiod, and concluded that $V_{c,\max}$ is under photoperiodic control, asserting that models should include this effect. However, their treatment also inevitably increased total daily photosynthetically active radiation (PAR) in autumn. On the basis of the information given about the experimental protocol, these results could therefore also be interpreted as showing that seasonal variations in $V_{c,\max}$ are related to daily total PAR.

This example draws attention to a key principle for the use of experimental results in model evaluation, namely that such comparisons are only valid if the models explicitly follow the experimental protocol. It is not sufficient for models to attempt to reproduce the stated general conclusions of experimental studies. The possibility of invalid comparisons can most easily be avoided through the inclusion of experimentalists from the outset in model evaluation projects. This was the case, for example, in the FACE data-model comparison study of De Kauwe et al. (2013).

This example also illustrates a general challenge for the modelling community. One response to new experimental studies is to increase model complexity by adding new processes based on the ostensible advances in knowledge. However, we advise a more critical and cautious approach, employing case-by-case comparisons of model results and experiments, rather than general interpretation of experiments, to reduce the potential for ambiguities and avoid unnecessary complexity in models. Such an approach would prevent the occurrence of overparameterisation, the implications of which have been explored by Crout et al. (2009).

2.2 Palaeodata

The key purpose of palaeo-evaluation is to establish whether the model has the correct sensitivity for large-scale processes. Models are typically developed using modern observations (i.e. under a limited range of climate conditions and behaviours), but we need to determine how well they simulate a large climate change, to assess whether they can

capture the behaviour of the system outside the modern range. If our understanding of the physics and biology of the system is correct, models should be able to predict past changes as well as present behaviour.

Reconstructions of global temperature changes over the last 1500 yr (e.g. Mann et al., 2009) are primarily derived from tree-ring and isotopic records, while reconstructions of climates over the last deglaciation and the Holocene (e.g. Davis et al., 2003; Viau et al., 2008; Seppä et al., 2009) are primarily derived from pollen data, although other biotic assemblages and geochemical data have been used at individual sites (e.g. Larocque and Bigler, 2004; Hou et al., 2006; Millet et al., 2009). Marine sediment cores have been used extensively to generate sea-surface temperature reconstructions (e.g. Marcott et al., 2013), and to reconstruct different past climate variables (see review in Henderson, 2002) related to ocean conditions. For example, $\delta^{13}\text{C}$ is used in reconstructions of ocean circulation, marine productivity, and biosphere carbon storage (Oliver et al., 2010). However, the interpretation of these data is often not straightforward, since the measured indicators are frequently influenced by more than one climatic variable (e.g. the benthic $\delta^{18}\text{O}$ measured in foraminiferal shells contains information on both global sea level and deep water temperature). Errors associated with the data and their interpretation also need to be stated, as while analytical errors on the measurements are often small, errors in the calibrations used to obtain reconstructions tend to be much bigger. Therefore the incorporation of measured variables such as marine carbonate concentrations (e.g. Ridgwell et al., 2007), $\delta^{18}\text{O}$ (e.g. Roche et al., 2004) or $\delta^{13}\text{C}$ (e.g. Crucifix, 2005) as variables in models is an important advance, because it allows comparison of model outputs directly with data, rather than relying on a potentially flawed comparison between modelled variables and the same variables reconstructed from chemical or isotopic measurements.

Ice cores provide a polar contribution to climate response reconstruction, as well as crucial information on a range of climate-relevant factors. For example, responses to forcings by solar variability (through ^{10}Be), volcanism (through sulfate spikes), and changes in the atmospheric concentration of greenhouse gases (e.g. CO_2 , CH_4 , N_2O) and mineral dust can be assessed. CH_4 can be measured in both Greenland and Antarctic ice cores. CO_2 measurements require Antarctic cores, due to the high concentrations of impurities in Greenland samples, which lead to the in situ production of CO_2 (Tschumi and Stauffer, 2000; Stauffer et al., 2002). For the last few millennia, choosing sites with the highest snow accumulation rates yields decadal resolution. The highest resolution records to date are from Law Dome (MacFarling Meure et al., 2006), making these data more reliable, particularly for model evaluation (e.g. Frank et al., 2010). Further work at high accumulation sites would provide reassurance on this point. Over longer time periods, sites with progressively lower snow accumulation rates, and therefore lower intrinsic time resolution, have to be used. Through

the Holocene (last $\sim 11\,000$ yr) (Elsig et al., 2009), and the last deglaciation, i.e. the transition out of the Last Glacial Maximum (LGM) into the Holocene (Lourantou et al., 2010; Schmidt et al., 2012), there are now high-quality $^{13}\text{C}/^{12}\text{C}$ of CO_2 data available, as well as much improved information about the phasing between the change in Antarctic temperature and CO_2 (Pedro et al., 2012; Parrenin et al., 2013), and between CO_2 and the global mean temperature (Shakun et al., 2012).

Compared to the amount of effort spent on reconstructing past climates and atmospheric composition, comparatively few data sets provide information on different components of the terrestrial carbon cycle. Nevertheless, there are data sets – synthesised from many individual published studies – that provide information on changes in vegetation distribution (e.g. Prentice et al., 2000; Bigelow et al., 2003; Harrison and Sanchez Goñi, 2010; Prentice et al., 2011a), biomass burning (Power et al., 2008; Daniau et al., 2012), and peat accumulation (e.g. Yu et al., 2010; Charman et al., 2013). These data sets are important because they can be used to test the response of individual components of ESMs to changes in forcing.

The major advantage of evaluating models using the palaeorecord is that it is possible to focus on times when the signal is large compared to the noise. The change in forcing at the LGM relative to the pre-industrial control is of comparable magnitude, though opposite in direction, to the change in forcing from quadrupling CO_2 relative to that same control (Izumi et al., 2013). Thus, comparisons of palaeoclimatic simulations and observations since the LGM can provide a measure of individual model performance, discriminate between models, and allow diagnosis of the sources of model error for a range of climate states similar in scope to those expected in the future. For example, Harrison et al. (2013) evaluated mid-Holocene and LGM simulations from the CMIP5 archive, and from the second phase of the Palaeoclimate Modelling Intercomparison Project (PMIP2), against observational benchmarks, using goodness-of-fit and bias metrics. However, as is the case for many modern observational data sets (e.g. Kelley et al., 2012), not all published palaeoreconstructions provide adequate documentation of errors and uncertainties, and there is a lack of standardisation between data sets where such estimates are provided (e.g. Leduc et al., 2010; Bartlein et al., 2011). Reconstructions based on ice or sediment cores are intrinsically site-specific (except for the globally significant greenhouse gas records), therefore many records are required to synthesise regional or global distribution patterns and estimates (Fig. 3). Community efforts to provide high-quality compilations of already available data (e.g., Waelbroeck et al., 2009; Bartlein et al., 2010) make it possible to use palaeodata for model evaluation, but an increase in the coverage of palaeoreconstructions is still required to evaluate model behaviour at regional scales.

Unfortunately, most attempts to compare simulations and reconstructions using palaeodata have focused on purely

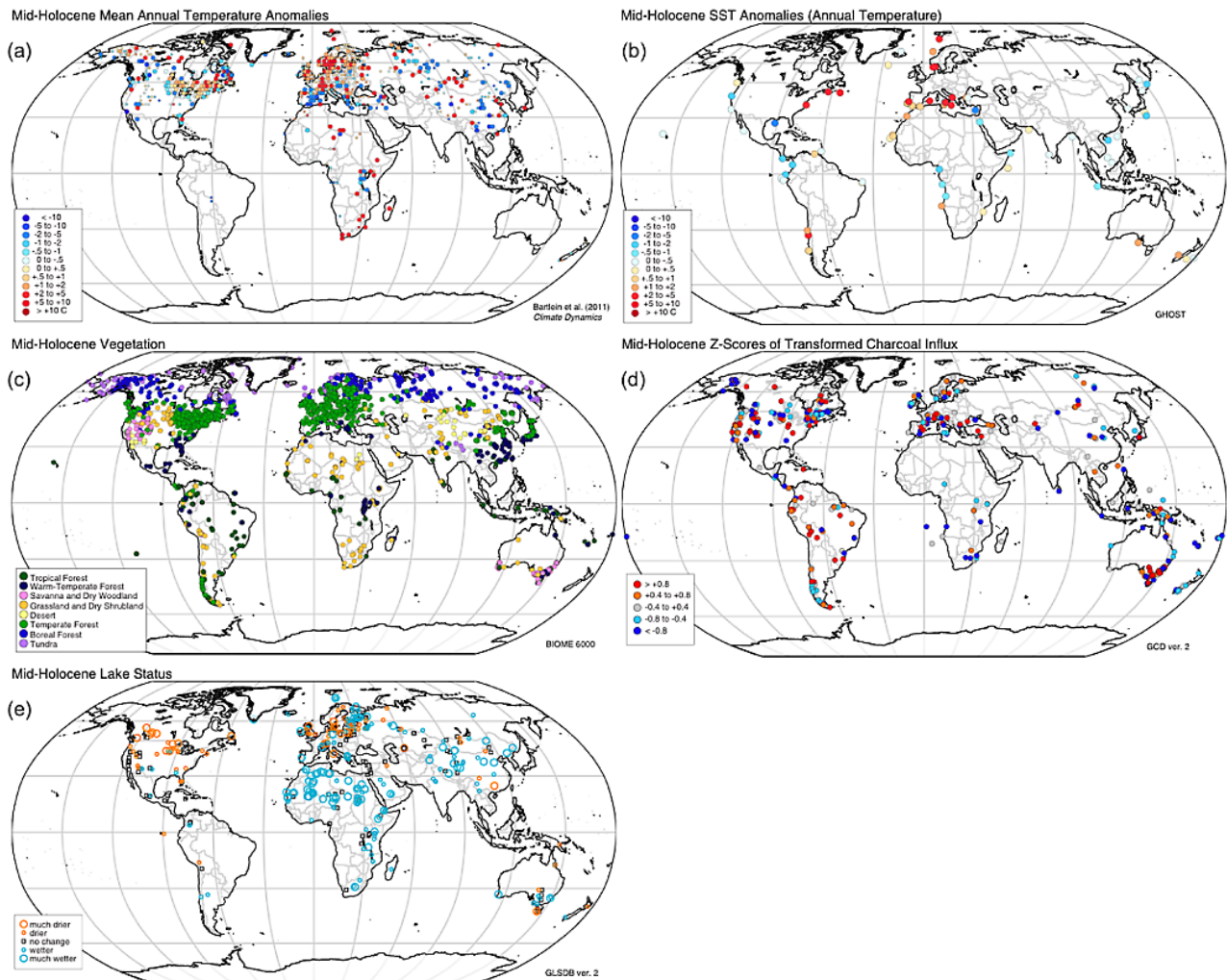


Fig. 3. Examples of global data sets documenting environmental conditions during the mid-Holocene (ca. 6000 yr ago) that can be used for benchmarking ESM simulations. In general, these are expressed as anomalies, i.e. the difference between the mid-Holocene and modern conditions: (a) pollen-based reconstructions of anomalies in mean annual temperature, (b) reconstructions of anomalies in sea-surface temperatures based on marine biological and chemical records, (c) pollen and plant macrofossil reconstructions of vegetation during the mid-Holocene, (d) charcoal records of the anomalies in biomass burning, and (e) anomalies of changes in the hydrological cycle based on lake-level records of the balance between precipitation and evaporation (after Harrison and Bartlein, 2012). (Reprinted from Harrison, S. P. and Bartlein, P.: Records from the Past, Lessons for the Future, in: The Future of the World's Climate, edited by: A. Henderson-Sellers and K. J. McGuffie, 403–436, Copyright© 2012, with permission from Elsevier.)

qualitative agreement of simulated and observed spatial patterns (e.g. Otto-Bleisner et al., 2007; Miller et al., 2010). There has been surprisingly little use of metrics for palaeodata-model comparisons (for exceptions see e.g. Guiot et al., 1999; Paul and Schäfer-Neth, 2004; Harrison et al., 2013). This situation probably reflects problems in developing meaningful ways of taking uncertainties into account in these comparisons. Quantitative assessments have generally focused on individual large-scale features of the climate system, for example the magnitude of insolation-induced increase in precipitation over northern Africa during the mid-Holocene (Joussaume et al., 1999; Jansen et al., 2007), zonal

cooling in the tropics at the LGM (Otto-Bleisner et al., 2009), or the amplification of cooling over Antarctica relative to the tropical oceans at the LGM (Masson-Delmotte et al., 2006; Braconnot et al., 2012). Comparisons of simulated vegetation changes have been based on assessments of the number of matches to site-based observations from a region (e.g. Harrison and Prentice, 2003; Wohlfahrt et al., 2004, 2008). Observational uncertainty is represented visually in such comparisons, and only used explicitly to identify extreme behaviour amongst the models. Nevertheless, the recent trend is towards explicit incorporation of uncertainties and systematic model benchmarking (Harrison et al., 2013; Izumi et al., 2013).

3 Key metrics for ESM evaluation

Many metrics have been proposed (Tables 2–4), and the choice of an appropriate metric in model evaluation is crucial because the use of inappropriate metrics can lead to overconfidence in model skill. The choice should be based on the properties of the data sets, the properties of the metric, and the specific objectives of the evaluation. Metric formalism – that is, the treatment of metrics as well-defined mathematical and statistical concepts – can help the interpretation of metrics, their analysis, or their combination into a “skill-score” (Taylor, 2001) in an objective way.

The use of metrics draws on the mathematical concept of “distance” ($d(x, y)$), expressed in terms of three characteristics: separation: $d(x, y) = 0 \iff x = y$, symmetry: $d(x, y) = d(y, x)$, and the triangle inequality $d(x, z) \leq d(x, y) + d(y, z)$. The two data sets could be two model outputs, where the metric is used to measure how similar the two models are, or one model output and one reference observation data set, where the metric is used to evaluate the model against real measurements. Three levels of metric complexity can be identified, relating to the state-space on which to apply the distance:

- Level 1 – “comparisons of raw biogeophysical variables”. Here the distance generally reflects errors and provides assessment of model performance where there is a reasonable degree of similarity between the model and reference data set (such as climate variables in weather models).
- Level 2 – “comparisons of statistics on biogeophysical variables”. Here the distance is measured on a statistical property of the data sets. This is particularly useful for models that are expected to characterise the statistical behaviour of a system (e.g. climate models). This level is appropriate for most of the biophysical variables simulated by ESMs.
- Level 3 – “comparisons of relationships among biogeophysical variables”. Here, the distance is diagnostic of relationships related to physical and/or biological processes and this level of comparison is therefore useful for understanding the behaviour of two data sets.

At all levels of metric complexity, the metric needs to be both synthetic enough to aid in understanding the similarities and differences between the two data sets, and to be understandable by non-specialists in order to facilitate its use by other communities. Next, the particular uses, advantages, and limitations of metrics in each level of metric complexity will be discussed.

3.1 Metrics on raw biogeophysical variables

Level 1 metrics are the most widely used. The distance measures the discrepancies between two data sets of a key bio-

geophysical variable. Discrepancies can be measured at site level or at pixel level for gridded data sets, and thus such comparisons can be used for model evaluation against sparse data, such as site-based NPP data (e.g. Zaehle and Friend, 2010), eddy-covariance data (e.g. Blyth et al., 2011), or atmospheric CO₂ concentration records at remote monitoring stations (e.g. Cadule et al., 2010; Dalmonech and Zaehle, 2013). Where there is sufficient data to make the calculation meaningful, comparisons can be made against spatial averages or global means of the biogeophysical variables. Comparisons can also be made in the time domain because climate change and climate variability act on Earth system components across a wide range of temporal scales. The distance can thus be measured on instantaneous variables or on time-averaged variables, such as annual means. Many distances, summarised in Table 2, can be considered to measure these discrepancies.

The Euclidean distance (Eq. 1) is the most commonly used distance. It is more sensitive to outliers than the Manhattan distance (Eq. 2). Both of these distances assume that direct comparisons of the data can be made. Some examples are reported in Jolliff et al. (2009), where the Euclidean distance is used to evaluate three ocean surface bio-optical fields.

In the case of the weighted Euclidean distance (Eq. 3), a weight is associated with each variable. This is useful for various reasons: (1) normalisation against a mean value provides a dimensionless metric and allows comparisons to be drawn between data sets with different orders of magnitude; (2) the weighting can take account of uncertainties in the reference data set (e.g. instrumental errors in an observational data set, or uncertainty in a model ensemble); and (3) this type of metric can be useful when the data have a different dynamical range. For example, in a time series of Northern Hemisphere monthly surface temperature, the variability is different for summer and winter, and it makes sense to normalise the differences by the variance.

The Chi-square “distance” (Eq. 4) is related to the Pearson Chi-square test or goodness-of-fit, and differs from previous distances discussed here as it measures the similarity between two probability density functions (PDFs), rather than between data points. It is particularly useful if the focus of the analysis is at the population level. Distances on PDFs are defined, in this paper, to be Level 2 metrics, but the Chi-square distance can be used when the geophysical variables are supposed to have a particular shape (e.g. an atmospheric profile of temperature). Equation (5) can also be used, in particular to facilitate the symmetry property of distances.

The Tchebychev distance (Eq. 6) can be used, for example, to identify the maximum annual discrepancy in a climatic run. It can be useful if the focus is on extreme events.

The Mahalanobis distance (Eq. 7) is particularly suitable if variables have very different units, as each one will be normalised by its variance, and/or if they are correlated with each other, since the distance takes these correlations into account. High correlation between two data sets has no impact

Table 2. Summary of Level 1 metrics (x, y represent points while D_1, D_2 are data sets).

Metric	Equation	Suitability
Euclidean distance	$d(x, y) = \sqrt{\sum_{i=1}^n (x_i - y_i)^2}$	(1) More sensitive to outliers compared to the Manhattan distance. Like the Manhattan distance, it also supposes that direct comparisons of the variables can be made.
Manhattan distance	$d(x, y) = \sum_{i=1}^n x_i - y_i $	(2) Implicitly supposes that x and y are comparable, so is not suited to mixed variables (e.g. variables with different units).
Weighted Euclidean distance	$d(x, y) = \sqrt{\sum_{i=1}^n w_i (x_i - y_i)^2}$	(3) Uncertainty in the reference data set, such as instrumental errors in an observation data set, or model uncertainty in a model ensemble, can be accounted for using a model efficiency metric: e.g. $w_i = \frac{1}{\sigma_i^2}$ where σ_i = uncertainty.
Chi-squared “distance”	$d(x, y) = \sqrt{\sum_{i=1}^n \frac{1}{2} \frac{(x_i - y_i)^2}{y_i}}$	(4) Measures how similar two PDFs (probability distribution functions) are. Particularly useful if the focus of the analysis is
	$d(x, y) = \sqrt{\sum_{i=1}^n \frac{1}{2} \frac{(x_i - y_i)^2}{x_i + y_i}}$	(5) at population level. Alternatively Eq. (5) can be used to facilitate the symmetry property of distances.
Tchebychev distance	$d(x, y) = \max_i (x_i - y_i)$	(6) Useful for extreme events, or maximum annual discrepancy in a climatic run.
Mahalanobis distance	$d(x, y) = \sqrt{(x - y)^T \mathbf{A}^{-1} (x - y)}$ where \mathbf{A}^{-1} = covariance matrix of x or y .	(7) Particularly useful if x or y include coordinates with very different units (each one will be normalised by its variance), if they are correlated one with each other (since the distance takes into account these correlations), and for combination of multiple sources of information.
Normalised mean error	$\frac{1}{E} \sum_e \frac{(D_{1,e} - D_{2,e})}{D_1 D_2}$ where E is the total number of samples in D_1 and D_2 .	(8) Applies the distance over the entirety of two data sets D_1 and D_2 .

on the distance computed, compared to two independent data sets. This distance is directly related to the quality criterion of the variational assimilation and Bayesian formalism that optimally combines weather forecast and real observations. This criterion needs to take into account the covariance matrices and the uncertainties of the state variables.

Interesting links can be established between metrics and the operational developments of the numerical weather prediction centres. The Mahalanobis distance is well suited for Gaussian distributions (meaning here that the data/model misfit distribution follows a Gaussian distribution with covariance matrix \mathbf{A} , e.g. Min and Hense, 2007). General Bayesian formalism can be used to generalise this distance to more complex distributions. The Mahalanobis distance and the more general Bayesian framework are particularly suitable to treat several evaluation issues at once, such as the quantification of multiple sources of error and uncertainty in models or the combination of multiple sources of information (including the acquisition of new information). For instance, Rowlands et al. (2012) use a goodness-of-fit statistic similar

to the Mahalanobis distance applied to surface temperature data sets.

We present here distances between two points, possibly multivariate. Some metrics use these distances and have been defined over the two whole data sets D_1 and D_2 . For example, the Normalized Mean Error (NME) is a normalisation of the bias between the two data sets (Eq. 8). Several other distances exist in the literature that have been applied in different scientific fields and that are not listed here (e.g. Deza and Deza, 2006). However most of these distances are particular cases or an extension of the preceding distances.

3.2 Metrics on statistical properties

Level 2 metrics, summarised in Table 3, use statistical quantities estimated for two data sets D_1 and D_2 . Some of the metrics presented in the previous section can then be applied to the selected statistics. For instance, the PDF can be estimated for both data sets and the Chi-Square distance can be used to measure their discrepancy. For example, Anav et al. (2013) compared the PDFs of gross primary production

Table 3. Summary of Level 2 metrics.

Metric	Equation	Suitability
Chi-squared distance	See Eq. (4).	For PDFs of two data sets (e.g. observed and modelled data).
Kullback–Leibler divergence	$d(p \parallel q) = \int p(x) \frac{p(x)}{q(x)} dx$ (9)	For PDFs of two data sets, p and q .
Variance	Depends on application.	Suitable if long observational record is available. Use the diagnostic that best suits the application.

(GPP) and leaf area index (LAI) from the CMIP5 model simulations with two selected data sets.

The Kullback–Leibler divergence (Eq. 9) is based on information theory and can also be used to measure the similarity of two PDFs. The Kolmogorov–Smirnov distance can be used when it is of interest to measure the maximum difference between the cumulative distributions. Tchebychev or other distances acting on estimated seasons are also considered here to be Level 2 metrics, since the seasons are statistical quantities estimated on D_1 and D_2 (although very close to level 1 raw geophysical variables). Similarly the distance can operate on derived variables from the original time series as decomposed signals in the frequency domain. Cadule et al. (2010), for example, analysed model performance in terms of representing the long-term trend and the seasonal signal of the atmospheric CO₂ record.

The variance of data and model is often used to formulate metrics for the quantification of the data-model similarity. In coupled systems, the use of a metric based on distance can become inadequate; the metric no longer facilitates definite conclusions on the model error, because it includes an unknown parameter in the form of the unforced variability. Furthermore, when applied to spatial fields, as variance is strongly location-dependent, a global spatial variance can be misleading. Gleckler et al. (2008) proposed a more suitable model variability index which has been applied to climatic variables, but is also highly applicable to several of the biogeophysical and biogeochemical variables simulated by land and ocean coupled models, and thus relevant to the carbon cycle. The metric can also focus on extreme events, with the distance acting on the percentile, assuming that the length of the records is sufficient to characterise these extremes.

3.3 Metrics on relationships

Level 3 metrics, summarised in Table 4, focus on relationships. The aim here is to diagnose a physical or a biophysical process that is particularly important, such as the link between two variables in the climate system. Various “relationship diagnostics” have been used, and are summarised in Table 4.

The correlation between two variables is a very simple and widely used metric; it satisfies the need to compare the data-

model phase correspondence of a particular biogeophysical variable. In this case parametric statistics such as the Pearson correlation coefficient (Eq. 10), or non-parametric statistics such as the Spearman correlation coefficient, are directly used as a metric. This is particularly used to evaluate the correspondence of the mean seasonal cycle of several variables, from precipitation (Taylor, 2001), to LAI, GPP (Anav et al., 2013), and atmospheric CO₂ (Dalmonech and Zaehle, 2013).

The sensitivity of one variable to another can be estimated using simple to very complex techniques (Aires and Rossow, 2003). It can be obtained by dividing concomitant perturbations of the two variables using spatial or temporal differences (Eq. 11), or by perturbing a model and measuring the impact when reaching equilibrium. The first approach can be used to evaluate, for example, site-level manipulative experiments to estimate carbon sensitivity to soil temperature or nitrogen deposition in terrestrial ecosystem models (e.g. Luo et al., 2012).

From the linear regression of two variables the slope or bias can be compared for D_1 or D_2 (Eq. 12). The slope is very close to the concept of sensitivity, but sensitivities are very dependent on the way they are measured. For example, sensitivity of the atmospheric CO₂ to climatic fluctuations may depend on the timescales they are calculated on (Cadule et al., 2010). An alternative, when more than two variables are involved in the physical or biophysical relationship under study, is a multiple linear regression (Eq. 13), or any other linear or nonlinear regression model such as neural networks. See, for example, the results obtained at site-level by Moffat et al. (2010).

Pattern-oriented approaches use graphs to identify particular patterns in the data set. These graphs aim at capturing relationships of more than two variables. For example, in Bony and Dufresne (2005), the tropical circulation is first decomposed into dynamical regimes using mid-tropospheric vertical velocity and then the sensitivity of the cloud forcing to a change in local sea surface temperature (SST) is examined for each dynamical regime. Moise and Delage (2011) proposed a metric that assesses the similarity of field structure of rainfall over the South Pacific Convergence Zone in terms of errors in replacement, rotation, volume, and pattern. The same metric could be applied to ocean Sea-viewing Wide Field-of-view Sensor (SeaWiFS) satellite-based fields

Table 4. Summary of Level 3 metrics.

Metric	Equation or method	Suitability
Pearson correlation coefficient	$\rho_{x,y} = \frac{\text{cov}(x,y)}{\sigma_x\sigma_y}$	(10) Very simple – measures only correlation, not the causality. Assumes data are normally distributed.
Sensitivity	$\frac{\Delta x}{\Delta y}$ From empirical increments or model experiments.	(11) Simple technique – univariate only. Does not consider possible interactions and non-linearities.
Linear regression	For two variables a and b , $a = c \cdot b + d$ where c = slope and d = bias.	Very simple. (12) Can become multivariate but has linear limitation. Provides “coincidence” information, not causality information.
Sensitivity (more advanced)	e.g. multiple linear regression $a = C \cdot B + d$ where B is now a vector of bio-geophysical variables. e.g. nonlinear model such as neural network.	(13) Nonlinear model that provides access to threshold, interactions and saturation behaviours. The metric can then be defined as the percentage of variance of a explained by B in the data and in the model. Still not causal.
Pattern-oriented approaches	Various methods.	Very process oriented, but requires a good understanding a priori of what needs to be examined.
Clustering algorithms	e.g. K-means, self-organising maps Uses a similarity distance, similar to level 1 metric.	Ideal for obtaining a limited set of prototypes, describing the variability of the data sets as much as possible.

in areas where particular spatial structures emerge. These powerful techniques could be more widely applied to evaluating ESM processes.

Clustering algorithms have been used to obtain weather regimes based only on the samples of a data set. For example, Jakob and Tselioudis (2003) and Ch eruy and Aires (2009) obtained cloud regimes based on cloud properties (optical thickness, cloud top pressure). The same methodology can be used in D_1 and D_2 and the two sets of regimes can be compared. The regimes can also be obtained on one data set and only the regime frequencies of the two data sets are compared. Abramowitz and Gupta (2008) applied a distance metric to compare several density functions of modelled net ecosystem exchange (NEE) clustered using the “self-organising map” technique.

It is often difficult to use a real mathematical distance to measure the discrepancy between the two “relationship diagnostics”. Although very useful for understanding differences in the physical behaviour, the simple comparison of two graphs (for D_1 and D_2) is not entirely satisfactory since it does not allow combination of multiple metrics or definition of scoring systems. In this paper, it is not possible to list all the ways to define a rigorous distance on each one of the relationship diagnostics that have been presented: Euclidean distance can be used on the regression parameters or the sensitivity coefficients, or two weather regime frequencies can be measured using confusion matrices (e.g. Aires et al., 2011). The distance needs to be adapted to the rela-

tionship diagnostic. The most limiting factor to this type of approach for ESM evaluation is that the relationship obtained might be not robust enough (i.e. statistically significant), or not easily framed within a process-based context.

4 A framework for robust model evaluation

Robust model evaluation relies on a combination of approaches, each informed using appropriate data and metrics (Fig. 4). Calibration and, ideally, pre-calibration (Sect. 4.2.2) must first be employed to rule out implausible outcomes, using data independent of that which may be subsequently used in model evaluation. Then, evaluation approaches must be a combination of process-focussed and system-wide, to ensure that both the representation of processes and the balance between them are realistic in the model. Optionally, the results of different model evaluation tests can be combined into a single model score, perhaps for the purposes of weighting future projections. When employed as part of a multi-model ensemble, the simulation can also contribute to the calculation of emergent constraints, which can then be used in subsequent model development (Sect. 4.3.3).

4.1 Recommendations for improved data availability and usage

The increasingly data-rich environment is both an opportunity and a challenge, in that it offers more opportunities for

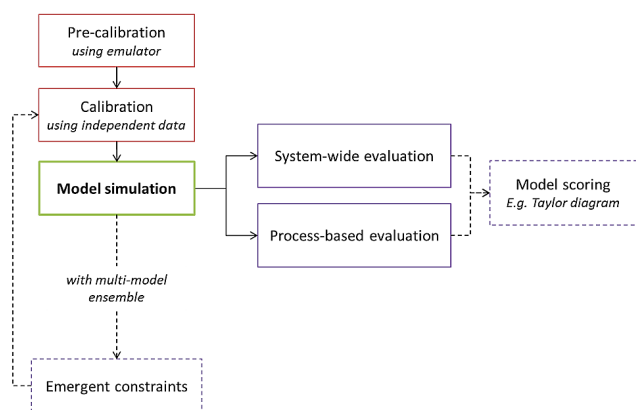


Fig. 4. Schematic diagram of model evaluation approaches, with optional approaches indicated by dashed lines.

model validation but requires more knowledge about the generation of data sets and their uncertainties in order to determine the best data set for evaluation of specific process representations. While improved documentation of data sets would go some way to alleviating the latter problem, there is scope for improved collaboration between the modelling and observational communities to develop an appropriate benchmarking system, that evolves to reflect new model developments (such as representing ecosystem-scale responses to combined environmental drivers) not addressed by existing benchmarks.

4.1.1 Coordinating data collection efforts

A key question for both the modelling and data communities to address together is how well model evaluation requirements and data availability are reconciled. There is an ongoing need for new and better data sets for model evaluation: data sets that are appropriately documented and for which useful information about errors and uncertainties are provided. The temporal and spatial coverage of data sets also needs to be sufficient to capture potential climatic perturbations, a point that is illustrated in the evaluation of marine productivity. Modelling studies offer conflicting evidence of the behaviour of this key variable in controlling marine carbon fluxes and exchanges of carbon with the atmosphere under a changing climate (e.g. Sarmiento et al., 2004; Steinacher et al., 2010; Taucher et al., 2012; Laufkötter et al., 2013), therefore model evaluation is essential. Recent compilations of observations of marine-productivity proxies give us a reasonably well-documented picture of qualitative changes in productivity over the last glacial-interglacial transition (e.g. Kohfeld et al., 2005), and in response to Heinrich events (e.g. Mariotti et al., 2012). These data sets are being used to evaluate the same ESMs used to predict changes in NPP in response to climate change (e.g. Bopp et al., 2003; Mariotti et al., 2012), and these studies show reasonable

agreement. On more recent timescales, remote sensing observations of ocean colour have been used to infer decadal changes in marine NPP. Studies show an increase in the extent of oligotrophic gyres over 1997–2008 with the SeaWiFS data (Polovina et al., 2008). However, on longer timescales, and using Coastal Zone Color Scanner (CZCS) and SeaWiFS data sets, analysis yields contrasting results of increase or decrease of NPP from 1979–1985 to 1998–2002 (Gregg et al., 2003; Antoine et al., 2005). Henson et al. (2010) have shown, based on a statistical analysis of biogeochemical model outputs, that an NPP time series of ~ 40 yr is needed to detect any global-warming-induced changes in NPP, highlighting the need for continued, focused data collection efforts.

4.1.2 Maximising the usefulness of current data in modelling studies

Modelling studies should be designed in a manner that makes the best use of the available data. For example, equilibrium model simulations of the distant past require time-slice reconstructions for evaluating processes relating to the carbon cycle. These reconstructions rely on synchronisation of records from ice cores, marine sediments, and terrestrial sequences, to take account of differences between forcings and responses in different archives, which is a significant effort even within a particular palaeo-archive, let alone across multiple archives. Yet the strength of palaeodata is precisely that it offers information about rates of change, and such information is discarded in a time-slice simulation. For that reason, the increasing use of transient model runs to simulate past climate and environmental changes is a particularly important development.

There is also an increasing need for forward modelling to simulate the quantities that are actually measured, such as isotopes in ice cores and pollen abundances. Ice core gas concentration measurements are unusual because what is measured is what we want to know, and is a variable that ESMs yield as a direct output. This is not generally the case, nor are all model setups easily able to simulate even the trace-gas isotopic data that are available from ice. A corollary is that we need to recognise the difficulty of trying to use palaeodata to reconstruct quantities that are essentially model constructs, for example inferring the strength of the meridional overturning circulation (MOC) from the $^{231}\text{Pa}/^{230}\text{Th}$ ratio in marine sediment cores (McManus et al., 2004). In the latter context, direct simulation of the $^{231}\text{Pa}/^{230}\text{Th}$ ratio is necessary to deconvolute the multiple competing processes (Siddall et al., 2005, 2007).

4.1.3 Using data availability to inform model development

Model development should also focus on incorporating processes that, at least collectively, are constrained by a wealth of data. Notable examples are processes such as those

governing methane (CH_4) emissions (e.g. from wetlands and permafrost) and the removal of methane from the atmosphere (e.g. via oxidation by the hydroxyl radical and atomic chlorine). There are four main observational constraints on the CH_4 budget with which we can evaluate the performance of ESMs: the concentration, $[\text{CH}_4]$; its isotopic composition with respect to carbon and deuterium, $\delta^{13}\text{CH}_4$ and δD (CH_4); and CH_4 fluxes at measurement sites. We have no natural record of CH_4 fluxes so their use in ESM evaluation is limited to the relatively recent period in which they have been measured, though measurements of CH_4 fluxes at specific sites can be used to verify spatial and seasonal distributions of CH_4 emissions inferred from tall tower and satellite measurements of $[\text{CH}_4]$, by inverse modelling. However, a range of $[\text{CH}_4]$, $\delta^{13}\text{CH}_4$, and δD (CH_4) records are available, spanning up to 800 000 yr in the case of polar ice cores, which can be used to evaluate the ability of ESMs to capture changes to the CH_4 budget in response to past changes in climate. The variety of climatic changes we can probe, from large glacial-interglacial changes spanning thousands of years to substantial changes over just a few tens of years at the beginning of Dansgaard–Oeschger events, and still more rapid, subtle changes following volcanic eruptions, enables us to evaluate the ability of ESMs to capture both the observed size and speed of changes known to have taken place. The complementary natures of the $[\text{CH}_4]$, $\delta^{13}\text{CH}_4$, and δD (CH_4) constraints is key to ESM evaluation. Each CH_4 source and sink affects these three constraints in different ways. As such, scenarios that explain only one set of observations can be eliminated. For instance, an increase in CH_4 emissions from tropical wetlands, biomass burning, or methane hydrates could explain an increase in $[\text{CH}_4]$, but of these only an increase in biomass burning emissions could explain an accompanying enrichment in $\delta^{13}\text{CH}_4$. Of course, more than one factor can change at a time, but the key point is that the most rigorous test of ESM performance utilises all three constraints and, therefore, ESMs should track the influence of each source and sink.

4.2 Recommendations for model calibration

4.2.1 Key principles of model calibration

Model evaluation is closely linked to model calibration. ESMs contain a large number of (sometimes poorly constrained) parameters, resulting from incomplete knowledge of certain processes or from the simplification of complex processes, which can be calibrated in order to improve model behaviour. In general, model calibration should follow a number of fundamental guiding principles. The principles detailed here are mostly based on the discussion in Petoukhov et al. (2000) for the CLIMBER-2 model.

First, parameters which are well constrained from observations or from theory must not be used for model calibration. Normally it would be physically inappropriate to mod-

ify the values of fundamental constants, for example, or use a value for a parameter which is different from the accepted empirical measurement just to improve the performance of the model.

Second, whenever possible, parameterisations and sub-modules should be tuned separately against observations rather than in the coupled system. In the case of parameterisations, this ensures that they represent the physical behaviour of the process described rather than their effect on the coupled system. The same principle should be applied as far as possible to the individual sub-modules of any ESM to make sure that their behaviour is self-consistent and to facilitate calibration of the much more complex fully coupled system.

Third, parameters must describe physical processes rather than unexplained differences between geographic regions. It is preferable for the model to represent the physical behaviour of the system rather than apply hidden flux corrections.

Fourth, the number of tuning parameters must be smaller than the predicted degrees of freedom. However, this is usually large for ESMs.

Finally, one of the key challenges relating to data used in ESM evaluation is to what extent ESM development and evaluation data are independent. In principle, the same observational data should not be used for calibration and evaluation. This is difficult to enforce in practice, however. Even if the observational data are divided into two parts, with one part used for calibration and the other for evaluation, for example, any mismatch in the evaluation will likely lead to a readjustment of model tuning parameters, making the evaluation not completely independent of the calibration procedure (Oreskes et al., 1994). Standard leave-one-out cross-validation techniques divide calibration data sets into multiple subsets, sequentially testing the calibration on each left-out subset (in the limit each data point) in turn but in Earth system modelling the subsets are unlikely to be fully independent.

4.2.2 Utilising pre-calibration to constrain implausible outcomes

The essence of pre-calibration is to apply weak constraints to model inputs in the initial ensemble design, and weak constraints on the model outputs to rule out implausible regions of input and output spaces (Edwards et al., 2011).

The pre-calibration approach is based on relatively simple statistical modelling tools and robust scientific judgements, but avoids the formidable challenges of applying full Bayesian calibration to a complex model (Rougier, 2007). A large set of model experiments sampling the variability in multiple input parameter values with the full simulator, here the ESM, is used to derive a statistical surrogate model or “emulator” of the dependence of key model outputs on uncertain model inputs. The choice of sampling points must be

highly efficient to span the input space and is usually based on Latin hypercube designs. The resulting emulator is computationally many orders of magnitude faster than the original model and can therefore be used for extensive, multi-dimensional sensitivity analyses to understand the behaviour of the model. Holden et al. (2010, 2013a, b) demonstrated the approach in constraining glacial and future terrestrial carbon storage.

The process is usually iterative, in that a large proportion of the initial parameter space may be deemed implausible, but one or more subsequent simulated ensembles can be designed by rejection sampling from the emulator to locate the not-implausible region of parameter space. The resulting simulated ensembles are then used to refine the emulator and the definition of the implausible space. The final output is an emulator of model behaviour and an ensemble of simulations, corresponding to a subset of parameter space that is deemed “plausible” in the sense that simulations from the identified parameter region do not disagree with a set of observational metrics by more than is deemed reasonable for the given simulator. The level of agreement is therefore dependent on the model and represents an assessment of the expected magnitude of its structural error (i.e. error due to choices for how processes are represented and relate to one another). The plausible ensemble, however, is a general result for the model that can be applied to any relevant prediction problems, and embodies an estimate of the structural and parametric error inherent in the model predictions.

Ideally, pre-calibration is a first step in a full Bayesian calibration analysis. The advantage of the logistic mapping or pure rejection sampling approach used is that, because no weighting is applied, a subsequent Bayesian calibration can be applied to refine the evaluation without any need to unravel convolution effects or the multiple use of constraints. In practice, however, the pre-calibration step can be sufficient to extract all the information that is readily available from top-down constraints given the magnitude of uncertainties in inputs and of structural errors in intermediate complexity ESMs.

4.3 Recommendations for model evaluation methodologies

4.3.1 Process-based (bottom-up) evaluation

Both bottom-up and top-down evaluation are required for evaluating ESMs: the first approach can give process-by-process information but not the balance between them; the second will give the balance but not the single terms. When bottom-up, process-based improvements can be shown to have top-down, system-level benefits, then we know our multi-pronged evaluation has worked.

Bottom-up, process-based evaluation will often require combinations of data to create the appropriate metrics as it is more likely to focus on the sensitivity of one output vari-

able to changes in a single input. For example, to assess if a model has the right sensitivity of NPP to precipitation a test could be to compute the partial derivative of NPP with respect to precipitation at constant values of temperature, radiation etc. for both the model and the observations (Randerson et al., 2009). This approach requires processing a data set of, in this case, NPP, to combine it with precipitation data to derive a relationship. The same NPP data could be combined with temperature data to derive a similar NPP(T) relationship. This is much more likely to isolate at least a small number of processes than simply comparing simulated NPP to an observational map or time series.

It is also common for model development to focus on specific features or aspects of the model in order to have faith in the model’s ability to make projections. For example, climate modelling centres may focus on the ability of their GCMs to represent coupled phenomena such as ENSO, or the timing and intensity of monsoon systems. In this way, bottom-up evaluation pinpoints important model processes, and helps to confirm that the model is a sufficiently accurate representation of the real system, giving the right results for the right reasons. However, a key limitation of this approach is that the relevant observations needed to assess a particular process may not exist.

Process-based evaluation requires metrics based on process-based sensitivities, as described in Sect. 3.2. Sensitivity analysis (e.g. Saltelli et al., 2000; Zaehle et al., 2005) may be useful to determine the parameters and processes to focus on in a bottom-up evaluation. In this approach, a simple statistical model is used to represent the physical relationships in the reference data set. A similar model is calibrated on the model simulations and the complex multivariate and non-linear relationships can then be compared.

Measuring these sensitivities allows prioritisation of the important parameters to validate in the model and isolate processes not well simulated in the model. For example, Aires et al. (2013) used neural networks to develop a reliable statistical model for the analysis of land–atmosphere interactions over the continental US in the North American Regional Reanalysis (NARR) data set. Such sensitivity analyses enable identification of key factors in the system and in this example, characterisation of rainfall frequency and intensity according to three factors: cloud triggering potential, low-level humidity deficit, and evaporative fraction.

4.3.2 System-level (top-down) evaluation

Top-down constraints tend to focus on whole-system behaviour and are more likely to involve evaluation of spatial or time-series data. Typical quantities used for top-down evaluations include surface temperature, pressure, precipitation, and wind speed maps. Observational data sets exist for many of these quantities throughout the atmosphere, so zonal-mean, or 3-dimensional comparisons are also possible (Randall et al., 2007). Anav et al. (2013) extend this approach

to assess new biogeochemical outputs of CMIP5 ESMs, such as distribution and time evolution of carbon stores and fluxes.

The appropriate choice of metrics is important, as discussed in Sect. 3. A correlation coefficient might seem an obvious choice to assess the seasonal cycle of a given variable, but a model with the right phase of seasonal cycle but a magnitude 5 times too big/small would score a high correlation coefficient, while a model with the correct magnitude but lagged by just one month would score poorly. To overcome these limitations of correlation-based metrics, additional metrics such as mean error should be included in the analysis to aid interpretation of the correlation, while lag errors could first be corrected so that the correlation gives a more meaningful result. There are also many studies that have attempted to overcome this issue by presenting summary statistical metrics for multiple components across multiple models.

Taylor (2001) is one of the examples where a metric based on correlation and a distance metric have been developed as a skill score. Gleckler et al. (2008) use Taylor diagrams to compare the performance of models in terms of both the magnitude and phase of the seasonal cycle. Reichler and Kim (2008) normalise model error variance on a grid-point basis to come up with a composite score and measure progress in model skill between generations of IPCC reports. Such scoring systems can be useful to synthesise the results of numerous metric comparisons, but should be used with caution as they can be hard to interpret – it is not always clear what model failing has led to a low score. The choice of which observations to use in the weighting is also subjective.

Model errors will inevitably evolve in time, affecting the reliability of simulations of future Earth system states. Measuring this type of uncertainty is an extremely difficult challenge. Presently, the best approach is to use expert judgement to estimate the growth of errors beyond the known forcing space, and this logic underpins the large, subjective choice of input ranges in the precalibration technique. Palaeoclimate analysis expands the space of forcings applied to the Earth system, such that possible future states might be more likely to occur inside the envelope of testable simulations. With sufficient high-quality data, it would be possible to cross-validate predictions against extreme past states that stretched the envelope in the most appropriate way. The Paleocene–Eocene Thermal Maximum offers perhaps the best opportunity for this, due to the large difference from current climate and atmospheric CO₂ conditions. An advanced theoretical approach is the “Reification” technique of Goldstein and Rougier (2009), which allows the error in a given model to be successively related to more and more accurate models, but its implementation is very much under development (see Williamson et al., 2012).

4.3.3 The role of emergent constraints in model evaluation

Emergent constraints (Table 5) can also provide valuable information for model evaluation, as they convert the extensive short-timescale information available for the contemporary period into longer-timescale constraints on the Earth system sensitivities that are most important for the 21st and 22nd centuries (e.g. climate sensitivity to CO₂, or carbon cycle sensitivity to climate). Observational data on short timescales do not relate directly to these sensitivities, and analogue approaches, which evaluate ESM sensitivity against known changes in the past, are also limited by observational data, as the analogue events in Earth’s past are not as well characterised as those in the contemporary period.

Emergent constraints relate some observable aspect of the contemporary Earth system to a key system sensitivity, using an ensemble of Earth system simulations (Collins et al., 2012). The archetypal example of this relates the magnitude of the snow-albedo feedback to the size of the seasonal cycle in snow cover in the Northern Hemisphere, across more than twenty GCMs (Hall and Qu, 2006). Since the seasonal cycle of snow-cover can be estimated from observations, this model-derived relationship provides a means to convert observations to a constraint on the size of the snow-albedo feedback in the real climate system, for which there is no direct reliable measurement. A similar emergent constraint has been used to relate the sensitivity of the interannual variability in atmospheric CO₂ to the loss of carbon from tropical land under climate change (Cox et al., 2013).

In general terms, such emergent constraint methods build on the realisation that analysis of short-time fluctuations in a system can assist in determining the sensitivity of that system to external forcing (Leith, 1975). Conversely, valuable information is unnecessarily lost when taking long-term trends and ignoring the shorter-timescale variations about these trends. Such emergent constraints utilise the large differences amongst ESM projections to reduce uncertainties in the sensitivities of the real Earth system to anthropogenic forcing.

5 Outlook

Although the current generation of ESMs encompass a wide range of processes, they are likely to become increasingly complex as processes that are currently being explored in, for example, dynamic global vegetation models such as better representation of nutrient cycles (e.g. Gotangco Castillo et al., 2012), fire (e.g. Thonicke et al., 2010; Prentice et al., 2011b; Pfeiffer et al., 2013), permafrost (e.g. Lawrence et al., 2012; Schaphoff et al., 2013) and wetland dynamics (e.g. Collins et al., 2011), or dust- (e.g. Shannon and Lunt, 2010), vegetation-climate interactions (Quillet et al., 2010), and aerosol-climate interactions (Woodage et al.,

Table 5. Summary of evaluation methodologies.

Type of evaluation	Description	Examples	Advantages	Limitations
Process-based, “bottom-up”	Looks at relationships between variables in a way that isolates a single process, or small number of processes.	NPP vs. precip. (Randerson et al., 2009). Magnitude of seasonal cycle of T_{air} vs. T_{surf} to evaluate insulation by snow pack.	Pinpoints important model processes. “Right answer for right reason.” Easy to interpret, e.g. can see if response is too big or too small for a given input.	Only targets a small part of the model. Relevant observations may not exist. Even when process representation is close to perfect, this does not ensure overall balance between them is right.
System-level, “top-down”	Compares large-scale model outputs that emerge from interactions between many processes within the model with relevant observations.	Global patterns of temperature, precipitation, etc. Seasonal cycle of carbon fluxes.	Evaluates end-result, i.e. quantities that we actually want the model to predict. Assesses overall balance between many (possibly finely balanced) processes.	Compensating errors: “Right answer for wrong reason”. Hard to interpret as offers no indication of what is causing an error and how to fix it.
Multi-model emergent constraints	Robust (across models) relationship between a quantity we can observe and a future change we want to predict.	Hall and Qu (2006): seasonal cycle of snow albedo. Cox et al. (2013): IAV of tropical carbon fluxes.	No requirement for models to be right – models might be wrong on individual basis regarding magnitude of response but the relationship may be robust. Guides where we want observational effort.	Relies on “bad” models more than “good” ones to derive regression. May get false confidence if models systematically wrong (e.g. all lack long-term carbon release from permafrost).

2010; Bellouin et al., 2011) are incorporated. This growing complexity has the potential to mask model errors, making robust evaluation of the model and its components increasingly necessary.

Common to any dynamical system under evaluation, key challenges include choosing the most important variables in the system, identifying the fundamental relationships, estimating non-linear and multivariate sensitivities, and analysing the interactions between processes. We have outlined how approaches such as pre-calibration and robust calibration, along with a combination of process- and system-level evaluation with relevant data, can be used to characterise model skill. We have also illustrated the usefulness of emergent constraints to further refine model outcomes.

A combination of approaches can greatly increase our understanding of a model’s ability to realistically simulate processes across multiple temporal and spatial scales. For example, both locally and globally, the net terrestrial carbon budget fluxes are a small difference between large uptake (photosynthesis) and release (respiration) terms. Even if each pro-

cess could be modelled with high precision, the net balance could still be poorly constrained. Hence, single, process-based tests are necessary but not sufficient. Conversely, observations of the seasonal cycle or interannual variability of carbon balance constrain overall terrestrial carbon balance, but do not provide detail about the processes contributing to it. It is theoretically possible to simulate the carbon balance with a number of different combinations of the components; therefore there is the potential to get the right answer for the wrong reasons. Different parameter combinations are potentially able to recreate the historical record of atmospheric CO_2 concentration (Sitch et al., 2008; Booth et al., 2012). Furthermore, some of the most accurate features of climate simulations (such as the pattern of near-surface temperatures) are poor predictors of the sensitivity of the terrestrial carbon balance to increasing CO_2 . It is thus eminently possible to get a skilful simulation of the present through the cancellation of multiple errors. A combination of “bottom-up” constraints on the processes and “top-down” constraints on the balance between them is essential, to give confidence that the model gives the right behaviour for the right reason.

A key limitation of current model evaluation approaches is that the widely used statistical measures of sensitivities are based on “coincident increments”, such as correlations, not on causality. A very interesting extension of sensitivities would investigate causal links among the important parameters in the system. Some tentative studies have investigated measures such as Granger causality; see Notaro et al. (2006) for an application to vegetation patterns. However, a more complete framework needs to be used (Pearl, 2009). Due to the complexity of this type of work, a close collaboration of climate-carbon cycle scientists and statisticians would be required.

Model complexity and structure has to be kept in mind when making comparisons of skill with respect to any given metric across a range of ESMs. Comparing models of different complexity could create an artificially large model spread, that does not reflect current process knowledge. However, comparing only models of similar complexity could lead to underestimation of the true uncertainty in model projections due to structural similarities between models and restricted sample size.

Benchmarking models against a set of well-chosen observations (Sect. 2), and using appropriate metrics (Sect. 3), should be considered a vital step in any model evaluation. While individual metrics might be each easily interpreted, a combination of many different metrics could be a challenge to interpret, particularly when very different scores in metrics that measure different aspects of model performance need to be reconciled. Therefore, while it may be tempting to simply evaluate the performance of the model against every data set that can be found (and indeed a “perfect” model should be able to withstand such a test), if this comes at the expense of being able to interpret the results then it may be more beneficial to focus on a smaller set of tests which target key model outputs. This level of discrimination is inevitably an expert judgement, but is necessary if the field of ESM evaluation is to move from “beauty contest” to constraint.

Acknowledgements. This paper emerged from the GREENCYCLESII mini-conference “Evaluation of Earth system models using modern and palaeo-observations” held at Clare College, Cambridge, UK, in September 2012. We would like to thank the Marie Curie FP7 Research and Training Network GREENCYCLESII for providing funding which made this meeting possible. Research leading to these results has received funding from the European Community’s Seventh Framework Programme (FP7 2007–2013) under grant agreement no. 238366. The work of C. D. Jones was supported by the Joint DECC/Defra Met Office Hadley Centre Climate Programme (GA01101). N. R. Edwards acknowledges support from FP7 grant no. 265170 (ERMITAGE). N. Vázquez Riveiros acknowledges support from the AXA Research Fund and the Newton Trust.

Edited by: V. Brovkin

References

- Abramowitz, G. and Gupta, H.: Toward a model space and model independence metric, *Geophys. Res. Lett.*, 35, L05705, doi:10.1029/2007GL032834, 2008.
- Ainsworth, E. A. and Long, S. P.: What have we learned from 15 years of free-air CO₂ enrichment (FACE)? A meta-analytic review of the responses of photosynthesis, canopy properties and plant production to rising CO₂, *New Phytol.*, 165, 351–371, doi:10.1111/j.1469-8137.2004.01224.x, 2005.
- Aires, F. and Rossow, W. B.: Inferring instantaneous, multivariate and nonlinear sensitivities for the analysis of feedback processes in a dynamical system: Lorenz model case-study, *Q. J. Roy. Meteor. Soc.*, 129, 239–275, doi:10.1256/qj.01.174, 2003.
- Aires, F., Bernardo, F., Brogniez, H., and Prigent, C.: An Innovative Calibration Method for the Inversion of Satellite Observations, *J. Appl. Meteorol.*, 49, 2458–2473, doi:10.1175/2010JAMC2435.1, 2010.
- Aires, F., Marquisseau, F., Prigent, C., and Sèze, G.: A land and ocean microwave cloud classification algorithm derived from AMSU-A and -B, trained using MSG-SEVIRI infrared and visible observations, *Mon. Weather Rev.*, 139, 2347–2366, doi:10.1175/MWR-D-10-05012.1, 2011.
- Aires, F., Gentine, P., Findell, K., Lintner, B. R., and Kerr, C.: Neural network-based sensitivity analysis of summertime convection over continental US, *J. Climate*, doi:10.1175/JCLI-D-13-00161.1, 2013.
- Anav, A., Friedlingstein, P., Kidston, M., Bopp, L., Ciais, P., Cox, P., Jones, C., Jung, M., Myneni, R., and Zhu, Z.: Evaluating the land and ocean components of the global carbon cycle in the CMIP5 Earth System Models, *J. Climate*, 26, 6801–6843, doi:10.1175/JCLI-D-12-00417.1, 2013.
- Antoine, D., Morel, A., Gordon, H. R., Banzon, V. F., and Evans, R. H.: Bridging ocean color observations of the 1980s and 2000s in search of long-term trends, *J. Geophys. Res.-Oceans*, 110, C06009, doi:10.1029/2004JC002620, 2005.
- Bartlein, P. J., Harrison, S. P., Brewer, S., Connor, S., Davis, B. A. S., Gajewski, K., Guiot, J., Harrison-Prentice, T. I., Henderson, A., Peyron, O., Prentice, I. C., Scholze, M., Seppä, H., Shuman, B., Sugita, S., Thompson, R. S., Vial, A. E., Williams, J., and Wu, H.: Pollen-based continental climate reconstructions at 6 and 21 ka: a global synthesis, *Clim. Dynam.*, 37, 775–802, doi:10.1007/s00382-010-0904-1, 2011.
- Batterman, S. A. and Larsen, K. S.: Integrating empirical studies and global models to improve climate change predictions, *Eos. T. Am. Geophys. Un.*, 92, 353–353, doi:10.1029/2011EO410011, 2011.
- Bauerle, W. L., Oren, R., Way, D. A., Qian, S. S., Stoy, P. C., Thornton, P. E., Bowden, J. D., Hoffman, F. M., and Reynolds, R. F.: Photoperiodic regulation of the seasonal pattern of photosynthetic capacity and the implications for carbon cycling, *P. Natl. Acad. Sci. USA*, 109, 8612–8617, doi:10.1073/pnas.1119131109, 2012.
- Beck, H. E., McVicar, T. R., Van Dijk, A. I. J. M., Schellekens, J., De Jeu, R. A. M., and Bruijnzeel, L. A.: Global evaluation of four AVHRR–NDVI data sets: Intercomparison and assessment against Landsat imagery, *Remote Sens. Environ.*, 115, 2547–2563, doi:10.1016/j.rse.2011.05.012, 2011.
- Bellouin, N., Rae, J., Jones, A., Johnson, C., Haywood, J., and Boucher, O.: Aerosol forcing in the Climate Model Intercom-

- parison Project (CMIP5) simulations by HadGEM2-ES and the role of ammonium nitrate, *J. Geophys. Res.*, 116, D20206, doi:10.1029/2011JD016074, 2011.
- Bigelow, N. H., Brubaker, L. B., Edwards, M. E., Harrison, S. P., Prentice, I. C., Anderson, P. M., Andreev, A. A., Bartlein, P. J., Christensen, T. R., Cramer, W., Kaplan, J. O., Lozhkin, A. V., Matveyeva, N. V., Murray, D. F., McGuire, A. D., Razzhivin, V. Y., Ritchie, J. C., Smith, B., Walker, D. A., Gajewski, K., Wolf, V., Holmqvist, B. H., Igarashi, Y., Kremenetskii, K., Paus, A., Pisaric, M. F. J., and Volkova, V. S.: Climate change and Arctic ecosystems: 1. vegetation changes north of 55° N between the last glacial maximum, mid-Holocene, and present, *J. Geophys. Res.*, 108, 8170, doi:10.1029/2002JD002558, 2003.
- Blyth, E., Clark, D. B., Ellis, R., Huntingford, C., Los, S., Pryor, M., Best, M., and Sitch, S.: A comprehensive set of benchmark tests for a land surface model of simultaneous fluxes of water and carbon at both the global and seasonal scale, *Geosci. Model Dev.*, 4, 255–269, doi:10.5194/gmd-4-255-2011, 2011.
- Bonan, G. B., Oleson, K. W., Fisher, R. A., Lasslop, G., and Reichstein, M.: Reconciling leaf physiological traits and canopy flux data: Use of the TRY and FLUXNET databases in the Community Land Model version 4, *J. Geophys. Res.-Biogeo.*, 117, G02026, doi:10.1029/2011JG001913, 2012.
- Bony, S. and Dufresne, J.-L.: Marine boundary layer clouds at the heart of tropical cloud feedback uncertainties in climate models, *Geophys. Res. Lett.*, 32, L20806, doi:10.1029/2005GL023851, 2005.
- Booth, B. B. B., Jones, C. D., Collins, M., Totterdell, I. J., Cox, P. M., Sitch, S., Huntingford, C., Betts, R. A., Harris, G. R., and Lloyd, J.: High sensitivity of future global warming to land carbon cycle processes, *Environ. Res. Lett.*, 7, 024002, doi:10.1088/1748-9326/7/2/024002, 2012.
- Bopp, L., Kohfeld, K. E., Le Quééré, C., and Aumont, O.: Dust impact on marine biota and atmospheric CO₂ during glacial periods, *Paleoceanography*, 18, 1046, doi:10.1029/2002PA000810, 2003.
- Bopp, L., Resplandy, L., Orr, J. C., Doney, S. C., Dunne, J. P., Gehlen, M., Halloran, P., Heinze, C., Ilyina, T., Séférian, R., Tjiputra, J., and Vichi, M.: Multiple stressors of ocean ecosystems in the 21st century: projections with CMIP5 models, *Biogeosciences*, 10, 6225–6245, doi:10.5194/bg-10-6225-2013, 2013.
- Braconnot, P., Harrison, S. P., Kageyama, M., Bartlein, P. J., Masson-Delmotte, V., Abe-Ouchi, A., Otto-Bliesner, B., and Zhao, Y.: Evaluation of climate models using palaeoclimatic data, *Nat. Clim. Change.*, 2, 417–424, doi:10.1038/nclimate1456, 2012.
- Cadule, P., Friedlingstein, P., Bopp, L., Sitch, S., Jones, C. D., Ciais, P., Piao, S. L., and Peylin, P.: Benchmarking coupled climate-carbon models against long-term atmospheric CO₂ measurements, *Global Biogeochem. Cy.*, 24, GB2016, doi:10.1029/2009GB003556, 2010.
- Charman, D. J., Beilman, D. W., Blaauw, M., Booth, R. K., Brewer, S., Chambers, F. M., Christen, J. A., Gallego-Sala, A., Harrison, S. P., Hughes, P. D. M., Jackson, S. T., Korhola, A., Mauquoy, D., Mitchell, F. J. G., Prentice, I. C., van der Linden, M., De Vleeschouwer, F., Yu, Z. C., Alm, J., Bauer, I. E., Corish, Y. M. C., Garneau, M., Hohl, V., Huang, Y., Karofeld, E., Le Roux, G., Loisel, J., Moschen, R., Nichols, J. E., Nieminen, T. M., MacDonald, G. M., Phadtare, N. R., Rausch, N., Sillasoo, Ü., Swindles, G. T., Tuittila, E.-S., Ukonmaanaho, L., Väliranta, M., van Bellen, S., van Geel, B., Vitt, D. H., and Zhao, Y.: Climate-related changes in peatland carbon accumulation during the last millennium, *Biogeosciences*, 10, 929–944, doi:10.5194/bg-10-929-2013, 2013.
- Chéruy, F. and Aires, F.: Cluster analysis of cloud properties over the southern European Mediterranean area in observations and a model, *Mon. Weather Rev.*, 137, 3161–3176, doi:10.1175/2009MWR2882.1, 2009.
- Claussen, M., Mysak, L., Weaver, A., Crucifix, M., Fichet, T., Loutre, M.-F., Weber, S., Alcamo, J., Alexeev, V., Berger, A., Calov, R., Ganopolski, A., Goosse, H., Lohmann, G., Lunkeit, F., Mokhov, I., Petoukhov, V., Stone, P., and Wang, Z.: Earth system models of intermediate complexity: closing the gap in the spectrum of climate system models, *Clim. Dynam.*, 18, 579–586, doi:10.1007/s00382-001-0200-1, 2002.
- Collins, M., Chandler, R. E., Cox, P. M., Huthnance, J. M., Rougier, J., and Stephenson, D. B.: Quantifying future climate change, *Nat. Clim. Change*, 2, 403–409, doi:10.1038/nclimate1414, 2012.
- Collins, W. J., Bellouin, N., Doutriaux-Boucher, M., Gedney, N., Halloran, P., Hinton, T., Hughes, J., Jones, C. D., Joshi, M., Liddicoat, S., Martin, G., O'Connor, F., Rae, J., Senior, C., Sitch, S., Totterdell, I., Wiltshire, A., and Woodward, S.: Development and evaluation of an Earth-System model – HadGEM2, *Geosci. Model Dev.*, 4, 1051–1075, doi:10.5194/gmd-4-1051-2011, 2011.
- Cox, P. M., Pearson, D., Booth, B. B., Friedlingstein, P., Huntingford, C., Jones, C. D., and Luke, C. M.: Sensitivity of tropical carbon to climate change constrained by carbon dioxide variability, *Nature*, 494, 341–344, doi:10.1038/nature11882, 2013.
- Crout, N. M. J., Tarsitano, D., and Wood, A. T.: Is my model too complex? Evaluating model formulation using model reduction, *Environ. Model. Softw.*, 24, 1–7, doi:10.1016/j.envsoft.2008.06.004, 2009.
- Crucifix, M.: Distribution of carbon isotopes in the glacial ocean: a model study, *Paleoceanography*, 20, PA4020, doi:10.1029/2005PA001131, 2005.
- Dahlke, C., Loew, A., and Reick, C.: Robust identification of global greening phase patterns from remote sensing vegetation products, *J. Climate*, 25, 8289–8307, doi:10.1175/JCLI-D-11-00319.1, 2012.
- Dalmonech, D. and Zaehle, S.: Towards a more objective evaluation of modelled land-carbon trends using atmospheric CO₂ and satellite-based vegetation activity observations, *Biogeosciences*, 10, 4189–4210, doi:10.5194/bg-10-4189-2013, 2013.
- Daniau, A.-L., Bartlein, P. J., Harrison, S. P., Prentice, I. C., Brewer, S., Friedlingstein, P., Harrison-Prentice, T. I., Inoue, J., Marlon, J. R., Mooney, S., Power, M. J., Stevenson, J., Tinner, W., Andrić, M., Atanassova, J., Behling, H., Black, M., Blarquez, O., Brown, K. J., Carcaillet, C., Colhoun, E., Colombaroli, D., Davis, B. A. S., D'Costa, D., Dodson, J., Dupont, L., Eshetu, Z., Gavin, D. G., Genries, A., Gebru, T., Haberle, S., Hallett, D. J., Horn, S., Hope, G., Katamura, F., Kennedy, L., Kershaw, P., Krivonogov, S., Long, C., Magri, D., Marinova, E., McKenzie, G. M., Moreno, P. I., Moss, P., Neumann, F. H., Norström, E., Paitre, C., Rius, D., Roberts, N., Robinson, G., Sasaki, N., Scott, L., Takahara, H., Terwilliger, V., Thevenon, F., Turner, R. B., Valsecchi, V. G., Vannièrè, B., Walsh, M., Williams, N., Zhang, Y.: Predictability

- of biomass burning in response to climate changes, *Global Biogeochem. Cy.*, 26, GB4007, doi:10.1029/2011GB004249, 2012.
- Davis, B. A. S., Brewer, S., Stevenson, A. C., and Guiot, J.: The temperature of Europe during the Holocene reconstructed from pollen data, *Quaternary Sci. Rev.*, 22, 1701–1716, doi:10.1016/S0277-3791(03)00173-2, 2003.
- De Kauwe, M. G., Medlyn, B. E., Zaehle, S., Walker, A. P., Dietze, M. C., Hickler, T., Jain, A. K., Luo, Y., Parton, W. J., Prentice, I. C., Smith, B., Thornton, P. E., Wang, S., Wang, Y.-P., Wårlind, D., Weng, E., Crous, K. Y., Ellsworth, D. S., Hanson, P. J., Kim, H.-S., Warren, J. M., Oren, R., and Norby, R. J.: Forest water use and water use efficiency at elevated CO₂: a model-data inter-comparison at two contrasting temperate forest FACE sites, *Glob. Change Biol.*, 19, 1759–1779, doi:10.1111/gcb.12164, 2013.
- Deza, E. and Deza, M.-M.: Chapter 17 – Distances and similarities in data analysis, in: *Dictionary of Distances*, Elsevier, Amsterdam, 217–229, 2006.
- Elsig, J., Schmitt, J., Leuenberger, D., Schneider, R., Eyer, M., Leuenberger, M., Joos, F., Fischer, H., and Stocker, T. F.: Stable isotope constraints on Holocene carbon cycle changes from an Antarctic ice core, *Nature*, 461, 507–510, doi:10.1038/nature08393, 2009.
- Edwards, N. R., Cameron, D., and Rougier, J.: Precalibrating an intermediate complexity climate model, *Clim. Dynam.*, 37, 1469–1482, doi:10.1007/s00382-010-0921-0, 2011.
- Evans, J. P.: 21st century climate change in the Middle East, *Climatic Change*, 92, 417–432, doi:10.1007/s10584-008-9438-5, 2008.
- Foley, A., Fealy, R., and Sweeney, J.: Model skill measures in probabilistic regional climate projections for Ireland, *Clim. Res.*, 56, 33–49, doi:10.3354/cr01140, 2013.
- Frank, D. C., Esper, J., Raible, C. C., Büntgen, U., Trouet, V., Stocker, B., and Joos, F.: Ensemble reconstruction constraints on the global carbon cycle sensitivity to climate, *Nature*, 463, 527–530, doi:10.1038/nature08769, 2010.
- Friedrichs, M. a. M., Carr, M.-E., Barber, R. T., Scardi, M., Antoine, D., Armstrong, R. a., Asanuma, I., Behrenfeld, M. J., Buitenhuis, E. T., Chai, F., Christian, J. R., Ciotti, A. M., Doney, S. C., Dowell, M., Dunne, J., Gentili, B., Gregg, W., Hoepffner, N., Ishizaka, J., Kameda, T., Lima, I., Marra, J., Mélin, F., Moore, J. K., Morel, A., O'Malley, R. T., O'Reilly, J., Saba, V. S., Schmeltz, M., Smyth, T. J., Tjiputra, J., Waters, K., Westberry, T. K. and Winguth, A.: Assessing the uncertainties of model estimates of primary productivity in the tropical Pacific Ocean, *J. Mar. Syst.*, 76, 113–133, doi:10.1016/j.jmarsys.2008.05.010, 2009.
- Friend, A. D., Arneeth, A., Kiang, N. Y., Lomas, M., Ogée, J., Rödenbeck, C., Running, S. W., Santaren, J.-D., Sitch, S., Viovy, N., Ian Woodward, F., and Zaehle, S.: FLUXNET and modelling the global carbon cycle, *Glob. Change Biol.*, 13, 610–633, doi:10.1111/j.1365-2486.2006.01223.x, 2007.
- Gleckler, P. J., Taylor, K. E., and Doutriaux, C.: Performance metrics for climate models, *J. Geophys. Res.*, 113, D06104, doi:10.1029/2007JD008972, 2008.
- Goldstein, M. and Rougier, J.: Reified Bayesian modelling and inference for physical systems, *J. Stat. Plan. Infer.*, 139, 1221–1239, doi:10.1016/j.jspi.2008.07.019, 2009.
- Gotangco Castillo, C. K., Levis, S., and Thornton, P.: Evaluation of the new CNDV option of the Community Land Model: effects of dynamic vegetation and interactive nitrogen on CLM4 means and variability*, *J. Climate*, 25, 3702–3714, doi:10.1175/JCLI-D-11-00372.1, 2012.
- Gregg, W. W., Conkright, M. E., Ginoux, P., O'Reilly, J. E., and Casey, N. W.: Ocean primary production and climate: global decadal changes, *Geophys. Res. Lett.*, 30, 1809, doi:10.1029/2003GL016889, 2003.
- Guiot, J., Boreux, J. J., Braconnot, P., and Torre, F.: Data-model comparison using fuzzy logic in paleoclimatology, *Clim. Dynam.*, 15, 569–581, doi:10.1007/s003820050301, 1999.
- Hall, A. and Qu, X.: Using the current seasonal cycle to constrain snow albedo feedback in future climate change, *Geophys. Res. Lett.*, 33, L03502, doi:10.1029/2005GL025127, 2006.
- Harrison, S. P. and Bartlein, P.: Records from the past, lessons for the future, in: *The Future of the World's Climate*, edited by: Henderson-Sellers, A. and McGuffie, K. J., Elsevier, 403–436, 2012.
- Harrison, S. P. and Prentice, C. I.: Climate and CO₂ controls on global vegetation distribution at the last glacial maximum: analysis based on palaeovegetation data, biome modelling and palaeoclimate simulations, *Glob. Change Biol.*, 9, 983–1004, doi:10.1046/j.1365-2486.2003.00640.x, 2003.
- Harrison, S. P. and Sanchez Goñi, M. F.: Global patterns of vegetation response to millennial-scale variability and rapid climate change during the last glacial period, *Quaternary Sci. Rev.*, 29, 2957–2980, doi:10.1016/j.quascirev.2010.07.016, 2010.
- Harrison, S. P., Bartlein, P. J., Brewer, S., Prentice, I. C., Boyd, M., Hessler, I., Holmgren, K., Izumi, K., and Willis, K.: Model benchmarking with glacial and mid-Holocene climates, *Clim. Dynam.*, doi:10.1007/s00382-013-1922-6, 2013.
- Henderson, G. M.: New oceanic proxies for paleoclimate, *Earth Planet. Sc. Lett.*, 203, 1–13, doi:10.1016/S0012-821X(02)00809-9, 2002.
- Henson, S. A., Sarmiento, J. L., Dunne, J. P., Bopp, L., Lima, I., Doney, S. C., John, J., and Beaulieu, C.: Detection of anthropogenic climate change in satellite records of ocean chlorophyll and productivity, *Biogeosciences*, 7, 621–640, doi:10.5194/bg-7-621-2010, 2010.
- Hickler, T., Smith, B., Prentice, I. C., MjöFors, K., Miller, P., Arneeth, A., and Sykes, M. T.: CO₂ fertilization in temperate FACE experiments not representative of boreal and tropical forests, *Glob. Change Biol.*, 14, 1531–1542, doi:10.1111/j.1365-2486.2008.01598.x, 2008.
- Holden, P. B., Edwards, N. R., Oliver, K. I. C., Lenton, T. M., and Wilkinson, R. D.: A probabilistic calibration of climate sensitivity and terrestrial carbon change in GENIE-1, *Clim. Dynam.*, 35, 785–806, doi:10.1007/s00382-009-0630-8, 2010.
- Holden, P. B., Edwards, N. R., Gerten, D., and Schaphoff, S.: A model-based constraint on CO₂ fertilisation, *Biogeosciences*, 10, 339–355, doi:10.5194/bg-10-339-2013, 2013a.
- Holden, P. B., Edwards, N. R., Müller, S. A., Oliver, K. I. C., Death, R. M., and Ridgwell, A.: Controls on the spatial distribution of oceanic $\delta^{13}\text{C}_{\text{DIC}}$, *Biogeosciences*, 10, 1815–1833, doi:10.5194/bg-10-1815-2013, 2013b.
- Hou, J., Huang, Y., Wang, Y., Shuman, B., Oswald, W. W., Faison, E., and Foster, D. R.: Postglacial climate reconstruction based on compound-specific D/H ratios of fatty acids from Blood Pond, New England, *Geochem. Geophys. Geos.*, 7, Q03008, doi:10.1029/2005GC001076, 2006.

- Izumi, K., Bartlein, P. J., and Harrison, S. P.: Consistent large-scale temperature responses in warm and cold climates, *Geophys. Res. Lett.*, 40, 1817–1823, doi:10.1002/grl.50350, 2013.
- Jansen, E., Overpeck, J., Briffa, K. R., Duplessy, J.-C., Joos, F., Masson-Delmotte, V., Olago, D., Otto-Bliesner, B., Peltier, W. R., Rahmstorf, S., Ramesh, R., Raynaud, D., Rind, D., Solomina, O., Villalba, R., and Zhang, D.: Palaeoclimate, 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: Solomon, S., Qin, D., Manning, M., Chen, Z., Marquis, M., Averyt, K. B., Tignor, M., and Miller, H. L., Cambridge University Press, Cambridge, UK and New York, NY, USA, 2007.
- Jakob, C. and Tselioudis, G.: Objective identification of cloud regimes in the Tropical Western Pacific, *Geophys. Res. Lett.*, 30, 2082, doi:10.1029/2003GL018367, 2003.
- Jolliff, J. K., Kindle, J. C., Shulman, I., Penta, B., Friedrichs, M. A. M., Helber, R., and Arnone, R. A.: Summary diagrams for coupled hydrodynamic-ecosystem model skill assessment, *J. Mar. Syst.*, 76, 64–82, doi:10.1016/j.jmarsys.2008.05.014, 2009.
- Jones, C., Gregory, J., Thorpe, R., Cox, P., Murphy, J., Sexton, D., and Valdes, P.: Systematic optimisation and climate simulation of FAMOUS, a fast version of HadCM3, *Clim. Dynam.*, 25, 189–204, doi:10.1007/s00382-005-0027-2, 2005.
- Jones, C., Robertson, E., Arora, V., Friedlingstein, P., Shevliakova, E., Bopp, L., Brovkin, V., Hajima, T., Kato, E., Kawamiya, M., Liddicoat, S., Lindsay, K., Reick, C. H., Roelandt, C., Segsneider, J., and Tjiputra, J.: 21st Century compatible CO₂ emissions and airborne fraction simulated by CMIP5 Earth System models under 4 representative concentration pathways, *J. Climate*, 26, 4398–4413, doi:10.1175/JCLI-D-12-00554.1, 2013.
- Joussaume, S., Taylor, K. E., Braconnot, P., Mitchell, J. F. B., Kutzbach, J. E., Harrison, S. P., Prentice, I. C., Broccoli, A. J., Abe-Ouchi, A., Bartlein, P. J., Bonfils, C., Dong, B., Guiot, J., Herterich, K., Hewitt, C. D., Jolly, D., Kim, J. W., Kislov, A., Kitoh, A., Loutre, M. F., Masson, V., McAvaney, B., McFarlane, N., de Noblet, N., Peltier, W. R., Peterschmitt, J. Y., Pollard, D., Rind, D., Royer, J. F., Schlesinger, M. E., Syktus, J., Thompson, S., Valdes, P., Vettoretti, G., Webb, R. S., and Wypytta, U.: Monsoon changes for 6000 years ago: results of 18 simulations from the Paleoclimate Modeling Intercomparison Project (PMIP), *Geophys. Res. Lett.*, 26, 859–862, doi:10.1029/1999GL900126, 1999.
- Jung, M., Reichstein, M., Margolis, H. A., Cescatti, A., Richardson, A. D., Arain, M. A., Arneeth, A., Bernhofer, C., Bonal, D., Chen, J., Gianelle, D., Gobron, N., Kiely, G., Kutsch, W., Lasslop, G., Law, B. E., Lindroth, A., Merbold, L., Montagnani, L., Moors, E. J., Papale, D., Sottocornola, M., Vaccari, F., and Williams, C.: Global patterns of land-atmosphere fluxes of carbon dioxide, latent heat, and sensible heat derived from eddy covariance, satellite, and meteorological observations, *J. Geophys. Res. Biogeo.*, 116, G00J07, doi:10.1029/2010JG001566, 2011.
- Kelley, D. I., Prentice, I. C., Harrison, S. P., Wang, H., Simard, M., Fisher, J. B., and Willis, K. O.: A comprehensive benchmarking system for evaluating global vegetation models, *Biogeosciences*, 10, 3313–3340, doi:10.5194/bg-10-3313-2013, 2013.
- Kohfeld, K. E., Quéré, C. L., Harrison, S. P., and Anderson, R. F.: Role of marine biology in glacial–interglacial CO₂ cycles, *Science*, 308, 74–78, doi:10.1126/science.1105375, 2005.
- Larocque, I. and Bigler, C.: Similarities and discrepancies between chironomid- and diatom-inferred temperature reconstructions through the Holocene at Lake 850, northern Sweden, *Quatern. Int.*, 122, 109–121, doi:10.1016/j.quaint.2004.01.033, 2004.
- Laufkötter, C., Vogt, M., and Gruber, N.: Long-term trends in ocean plankton production and particle export between 1960–2006, *Biogeosciences Discuss.*, 10, 5923–5975, doi:10.5194/bgd-10-5923-2013, 2013.
- Lawrence, D. M., Slater, A. G., and Swenson, S. C.: Simulation of present-day and future permafrost and seasonally frozen ground conditions in CCSM4, *J. Climate*, 25, 2207–2225, doi:10.1175/JCLI-D-11-00334.1, 2012.
- Leduc, G., Schneider, R., Kim, J.-H., and Lohmann, G.: Holocene and Eemian sea surface temperature trends as revealed by alkenone and Mg/Ca paleothermometry, *Quaternary Sci. Rev.*, 29, 989–1004, doi:10.1016/j.quascirev.2010.01.004, 2010.
- Leffelaar, P. A.: On scale problems in modelling: an example from soil ecology, in: *Theoretical Production Ecology?: Reflections and Prospects*, edited by: Rabbinge, J. G. R., Pudoc, Wagenin- gen, 57–73, available at: <http://edepot.wur.nl/171931> (last access: 24 June 2013), 1990.
- Leith, C. E.: Climate response and fluctuation dissipation, *J. Atmos. Sci.*, 32, 2022–2026, doi:10.1175/1520-0469(1975)032<2022:CRAFD>2.0.CO;2, 1975.
- Le Treut, H., Somerville, R., Cubasch, U., Ding, Y., Mauritzen, C., Mokssit, A., Peterson, T., and Prather, M.: Historical overview of climate change, 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: Solomon, S., Qin, D., Manning, M., Chen, Z., Marquis, M., Averyt, K. B., Tignor, M., and Miller, H. L., Cambridge University Press, Cambridge, UK and New York, NY, USA, 2007.
- Lohmann, G., Schulz, M., and Otto-Bliesner, B.: Earth System Models: Their Use and Reliability in Past Climate Reconstructions and Future Predictions, available at: <http://pages-142.unibe.ch/cgi-bin/WebObjects/products.woa/wa/product?id=300> (last access: 3 April 2013), 2008.
- Lourantou, A., Lavric, J. V., Köhler, P., Barnola, J.-M., Paillard, D., Michel, E., Raynaud, D., and Chappellaz, J.: Constraint of the CO₂ rise by new atmospheric carbon isotopic measurements during the last deglaciation, *Global Biogeochem. Cy.*, 24, GB2015, doi:10.1029/2009GB003545, 2010.
- Luo, Y. Q., Randerson, J. T., Abramowitz, G., Bacour, C., Blyth, E., Carvalhais, N., Ciais, P., Dalmonech, D., Fisher, J. B., Fisher, R., Friedlingstein, P., Hibbard, K., Hoffman, F., Huntzinger, D., Jones, C. D., Koven, C., Lawrence, D., Li, D. J., Mahecha, M., Niu, S. L., Norby, R., Piao, S. L., Qi, X., Peylin, P., Prentice, I. C., Riley, W., Reichstein, M., Schwalm, C., Wang, Y. P., Xia, J. Y., Zaehle, S., and Zhou, X. H.: A framework for benchmarking land models, *Biogeosciences*, 9, 3857–3874, doi:10.5194/bg-9-3857-2012, 2012.
- MacFarling Meure, C., Etheridge, D., Trudinger, C., Steele, P., Langenfelds, R., Van Ommen, T., Smith, A., and Elkins, J.: Law dome CO₂, CH₄ and N₂O ice core records extended to 2000 years BP, *Geophys. Res. Lett.*, 33, L14810, doi:10.1029/2006GL026152, 2006.

- McManus, J. F., Francois, R., Gherardi, J.-M., Keigwin, L. D., and Brown-Leger, S.: Collapse and rapid resumption of Atlantic meridional circulation linked to deglacial climate changes, *Nature*, 428, 834–837, doi:10.1038/nature02494, 2004.
- Mann, M. E., Zhang, Z., Rutherford, S., Bradley, R. S., Hughes, M. K., Shindell, D., Ammann, C., Faluvegi, G., and Ni, F.: Global signatures and dynamical origins of the Little Ice Age and medieval climate anomaly, *Science*, 326, 1256–1260, doi:10.1126/science.1177303, 2009.
- Marcott, S. A., Shakun, J. D., Clark, P. U., and Mix, A. C.: A reconstruction of regional and global temperature for the past 11 300 years, *Science*, 339, 1198–1201, doi:10.1126/science.1228026, 2013.
- Mariotti, V., Bopp, L., Tagliabue, A., Kageyama, M., and Swingedouw, D.: Marine productivity response to Heinrich events: a model-data comparison, *Clim. Past*, 8, 1581–1598, doi:10.5194/cp-8-1581-2012, 2012.
- Masson-Delmotte, V., Kageyama, M., Braconnot, P., Charbit, S., Krinner, G., Ritz, C., Guilyardi, E., Jouzel, J., Abe-Ouchi, A., Crucifix, M., Gladstone, R. M., Hewitt, C. D., Kitoh, A., LeGrande, A. N., Marti, O., Merkel, U., Motoi, T., Ohgaito, R., Otto-Bliesner, B., Peltier, W. R., Ross, I., Valdes, P. J., Vettoretti, G., Weber, S. L., Wolk, F., and Yu, Y.: Past and future polar amplification of climate change: climate model intercomparisons and ice-core constraints, *Clim. Dynam.*, 27, 437–440, doi:10.1007/s00382-006-0149-1, 2006.
- Miller, G. H., Alley, R. B., Brigham-Grette, J., Fitzpatrick, J. J., Polyak, L., Serreze, M. C., and White, J. W. C.: Arctic amplification: can the past constrain the future?, *Quaternary Sci. Rev.*, 29, 1779–1790, doi:10.1016/j.quascirev.2010.02.008, 2010.
- Millet, L., Arnaud, F., Heiri, O., Magny, M., Verneaux, V., and Desmet, M.: Late-Holocene summer temperature reconstruction from chironomid assemblages of Lake Anterne, northern French Alps, *Holocene*, 19, 317–328, doi:10.1177/0959683608100576, 2009.
- Min, S.-K. and Hense, A.: Hierarchical evaluation of IPCC AR4 coupled climate models with systematic consideration of model uncertainties, *Clim. Dynam.*, 29, 853–868, doi:10.1007/s00382-007-0269-2, 2007.
- Moffat, A. M., Beckstein, C., Churkina, G., Mund, M., and Heimann, M.: Characterization of ecosystem responses to climatic controls using artificial neural networks, *Glob. Change Biol.*, 16, 2737–2749, doi:10.1111/j.1365-2486.2010.02171.x, 2010.
- Moise, A. F. and Delage, F. P.: New climate model metrics based on object-orientated pattern matching of rainfall, *J. Geophys. Res.*, 116, D12108, doi:10.1029/2010JD015318, 2011.
- Najjar, R. G., Jin, X., Louanchi, F., Aumont, O., Caldeira, K., Doney, S. C., Dutay, J.-C., Follows, M., Gruber, N., Joos, F., Lindsay, K., Maier-Reimer, E., Matear, R. J., Matsumoto, K., Monfray, P., Mouchet, A., Orr, J. C., Plattner, G.-K., Sarmiento, J. L., Schlitzer, R., Slater, R. D., Weirig, M.-F., Yamanaka, Y. and Yool, A.: Impact of circulation on export production, dissolved organic matter, and dissolved oxygen in the ocean: Results from Phase II of the Ocean Carbon-cycle Model Intercomparison Project (OCMIP-2), *Global Biogeochem. Cy.*, 21, GB3007, doi:10.1029/2006GB002857, 2007.
- Norby, R. J. and Zak, D. R.: Ecological lessons from Free-Air CO₂ Enrichment (FACE) experiments, *Annu. Rev. Ecol. Evol. S.*, 42, 181–203, doi:10.1146/annurev-ecolsys-102209-144647, 2011.
- Norby, R. J., DeLucia, E. H., Gielen, B., Calfapietra, C., Giardina, C. P., King, J. S., Ledford, J., McCarthy, H. R., Moore, D. J. P., Ceulemans, R., De Angelis, P., Finzi, A. C., Karnosky, D. F., Kubiske, M. E., Lukac, M., Pregitzer, K. S., Scarascia-Mugnozza, G. E., Schlesinger, W. H., and Oren, R.: Forest response to elevated CO₂ is conserved across a broad range of productivity, *P. Natl. Acad. Sci. USA*, 102, 18052–18056, doi:10.1073/pnas.0509478102, 2005.
- Notaro, M., Liu, Z., and Williams, J. W.: Observed vegetation–climate feedbacks in the United States, *J. Climate*, 19, 763–786, doi:10.1175/JCLI3657.1, 2006.
- Nowak, R. S., Ellsworth, D. S., and Smith, S. D.: Functional responses of plants to elevated atmospheric CO₂ – do photosynthetic and productivity data from FACE experiments support early predictions?, *New Phytol.*, 162, 253–280, doi:10.1111/j.1469-8137.2004.01033.x, 2004.
- Oliver, K. I. C., Hoogakker, B. A. A., Crowhurst, S., Henderson, G. M., Rickaby, R. E. M., Edwards, N. R., and Elderfield, H.: A synthesis of marine sediment core $\delta^{13}\text{C}$ data over the last 150 000 years, *Clim. Past*, 6, 645–673, doi:10.5194/cp-6-645-2010, 2010.
- Oreskes, N., Shrader-Frechette, K., and Belitz, K.: Verification, validation, and confirmation of numerical models in the Earth sciences, *Science*, 263, 641–646, doi:10.1126/science.263.5147.641, 1994.
- Otto-Bliesner, B. L., Hewitt, C. D., Marchitto, T. M., Brady, E., Abe-Ouchi, A., Crucifix, M., Murakami, S., and Weber, S. L.: Last Glacial Maximum ocean thermohaline circulation: PMIP2 model intercomparisons and data constraints, *Geophys. Res. Lett.*, 34, L12706, doi:10.1029/2007GL029475, 2007.
- Otto-Bliesner, B. L., Schneider, R., Brady, E. C., Kucera, M., Abe-Ouchi, A., Bard, E., Braconnot, P., Crucifix, M., Hewitt, C. D., Kageyama, M., Marti, O., Paul, A., Rosell-Melé, A., Waelbroeck, C., Weber, S. L., Weinelt, M., and Yu, Y.: A comparison of PMIP2 model simulations and the MARGO proxy reconstruction for tropical sea surface temperatures at last glacial maximum, *Clim. Dynam.*, 32, 799–815, doi:10.1007/s00382-008-0509-0, 2009.
- Parrenin, F., Masson-Delmotte, V., Köhler, P., Raynaud, D., Pailard, D., Schwander, J., Barbante, C., Landais, A., Wegner, A., and Jouzel, J.: Synchronous change of atmospheric CO₂ and Antarctic temperature during the last deglacial warming, *Science*, 339, 1060–1063, 2013.
- Paul, A. and Schäfer-Neth, C.: How to combine sparse proxy data and coupled climate models, *Quaternary Sci. Rev.*, 24, 1095–1107, doi:10.1016/j.quascirev.2004.05.010, 2005.
- Pearl, J.: *Causality: Models, Reasoning, and Inference*, Cambridge University Press, 2000.
- Pedro, J. B., Rasmussen, S. O., and van Ommen, T. D.: Tightened constraints on the time-lag between Antarctic temperature and CO₂ during the last deglaciation, *Clim. Past*, 8, 1213–1221, doi:10.5194/cp-8-1213-2012, 2012.
- Petoukhov, V., Ganopolski, A., Brovkin, V., Claussen, M., Eliseev, A., Kubatzki, C., and Rahmstorf, S.: CLIMBER-2: a climate system model of intermediate complexity, Part I: model description

- and performance for present climate, *Clim. Dynam.*, 16, 1–17, doi:10.1007/PL00007919, 2000.
- Pfeiffer, M., Spessa, A., and Kaplan, J. O.: A model for global biomass burning in preindustrial time: LPJ-LMfire (v1.0), *Geosci. Model Dev.*, 6, 643–685, doi:10.5194/gmd-6-643-2013, 2013.
- Pincus, R., Batstone, C. P., Hofmann, R. J. P., Taylor, K. E., and Glecker, P. J.: Evaluating the present-day simulation of clouds, precipitation, and radiation in climate models, *J. Geophys. Res.*, 113, D14209, doi:10.1029/2007JD009334, 2008.
- Polovina, J. J., Chai, F., Howell, E. A., Kobayashi, D. R., Shi, L., and Chao, Y.: Ecosystem dynamics at a productivity gradient: a study of the lower trophic dynamics around the northern atolls in the Hawaiian Archipelago, *Prog. Oceanogr.*, 77, 217–224, doi:10.1016/j.pocean.2008.03.011, 2008.
- Power, M., Marlon, J., Ortiz, N., Bartlein, P., Harrison, S., Mayle, F., Ballouche, A., Bradshaw, R., Carcaillet, C., Cordova, C., Mooney, S., Moreno, P., Prentice, I., Thonicke, K., Tinner, W., Whitlock, C., Zhang, Y., Zhao, Y., Ali, A., Anderson, R., Beer, R., Behling, H., Briles, C., Brown, K., Brunelle, A., Bush, M., Camill, P., Chu, G., Clark, J., Colombaroli, D., Connor, S., Daniiau, A.-L., Daniels, M., Dodson, J., Doughty, E., Edwards, M., Finsinger, W., Foster, D., Frechette, J., Gaillard, M.-J., Gavin, D., Gobet, E., Haberle, S., Hallett, D., Higuera, P., Hope, G., Horn, S., Inoue, J., Kaltenrieder, P., Kennedy, L., Kong, Z., Larsen, C., Long, C., Lynch, J., Lynch, E., McGlone, M., Meeks, S., Mensing, S., Meyer, G., Minckley, T., Mohr, J., Nelson, D., New, J., Newnham, R., Noti, R., Oswald, W., Pierce, J., Richard, P., Rowe, C., Sanchez Goñi, M., Shuman, B., Takahara, H., Toney, J., Turney, C., Urrego-Sanchez, D., Umbanhowar, C., Vandergoes, M., Vanniore, B., Vescovi, E., Walsh, M., Wang, X., Williams, N., Wilmshurst, J., and Zhang, J.: Changes in fire regimes since the Last Glacial Maximum: an assessment based on a global synthesis and analysis of charcoal data, *Clim. Dynam.*, 30, 887–907, doi:10.1007/s00382-007-0334-x, 2008.
- Prentice, I. C., Jolly, D., and BIOME 6000 participants: Mid-Holocene and Glacial-Maximum vegetation geography of the northern continents and Africa, *J. Biogeogr.*, 27, 507–519, 2000.
- Prentice, I. C., Harrison, S. P., and Bartlein, P. J.: Global vegetation and terrestrial carbon cycle changes after the last ice age, *New Phytol.*, 189, 988–998, doi:10.1111/j.1469-8137.2010.03620.x, 2011a.
- Prentice, I. C., Kelley, D. I., Foster, P. N., Friedlingstein, P., Harrison, S. P., and Bartlein, P. J.: Modeling fire and the terrestrial carbon balance, *Global Biogeochem. Cy.*, 25, GB3005, doi:10.1029/2010GB003906, 2011b.
- Quillet, A., Peng, C., and Garneau, M.: Toward dynamic global vegetation models for simulating vegetation–climate interactions and feedbacks: recent developments, limitations, and future challenges, *Environ. Rev.*, 18, 333–353, doi:10.1139/A10-016, 2010.
- Radic, V. and Clarke, G. K. C.: Evaluation of IPCC models' performance in simulating late-twentieth-century climatologies and weather patterns over North America, *J. Climate*, 24, 5257–5274, doi:10.1175/JCLI-D-11-00011.1, 2011.
- Randall, D. A., Wood, R. A., Bony, S., Colman, R., Fichet, T., Fyfe, J., Kattsov, V., Pitman, A., Shukla, J., Srinivasan, J., Stouffer, R. J., Sumi, A., and Taylor, K. E.: 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: Solomon, S., Qin, D., Manning, M., Chen, Z., Marquis, M., Averyt, K. B., Tignor, M., and Miller, H. L., Cambridge University Press, Cambridge, UK and New York, NY, USA, 2007.
- Randerson, J. T., Hoffman, F. M., Thornton, P. E., Mahowald, N. M., Lindsay, K., Lee, Y.-H., Nevison, C. D., Doney, S. C., Bonan, G., Stöckli, R., Covey, C., Running, S. W., and Fung, I. Y.: Systematic assessment of terrestrial biogeochemistry in coupled climate–carbon models, *Glob. Change Biol.*, 15, 2462–2484, doi:10.1111/j.1365-2486.2009.01912.x, 2009.
- Reichler, T. and Kim, J.: Uncertainties in the climate mean state of global observations, reanalyses, and the GFDL climate model, *J. Geophys. Res.-Atmos.*, 113, D05106, doi:10.1029/2007JD009278, 2008.
- Ridgwell, A.: Interpreting transient carbonate compensation depth changes by marine sediment core modeling, *Paleoceanography*, 22, PA4102, doi:10.1029/2006PA001372, 2007.
- Roche, D., Paillard, D., and Cortijo, E.: Constraints on the duration and freshwater release of Heinrich event 4 through isotope modelling, *Nature*, 432, 379–382, doi:10.1038/nature03059, 2004.
- Roe, G. H. and Baker, M. B.: Why is climate sensitivity so unpredictable?, *Science*, 318, 629–632, doi:10.1126/science.1144735, 2007.
- Rougier, J.: Probabilistic inference for future climate using an ensemble of climate model evaluations, *Climatic Change*, 81, 247–264, doi:10.1007/s10584-006-9156-9, 2007.
- Rowlands, D. J., Frame, D. J., Ackerley, D., Aina, T., Booth, B. B., Christensen, C., Collins, M., Faull, N., Forest, C. E., Grandey, B. S., Gryspeerdt, E., Highwood, E. J., Ingram, W. J., Knight, S., Lopez, A., Massey, N., McNamara, F., Meinshausen, N., Piani, C., Rosier, S. M., Sanderson, B. M., Smith, L. A., Stone, D. A., Thurston, M., Yamazaki, K., Yamazaki, Y. H., and Allen, M. R.: Broad range of 2050 warming from an observationally constrained large climate model ensemble, *Nat. Geosci.*, 5, 256–260, doi:10.1038/ngeo1430, 2012.
- Saltelli, A., Chan, K., and Scott, E. M.: *Sensitivity Analysis: Gauging the Worth of Scientific Models*, Wiley, 2000.
- Sargent, R. G.: Verification and validation of simulation models, in: *Proceedings of the 2010 Winter Simulation Conference (WSC)*, 166–183, 2010.
- Sarmiento, J. L., Slater, R., Barber, R., Bopp, L., Doney, S. C., Hirst, A. C., Kleypas, J., Matear, R., Mikolajewicz, U., Monfray, P., Soldatov, V., Spall, S. A., and Stouffer, R.: Response of ocean ecosystems to climate warming, *Global Biogeochem. Cy.*, 18, GB3003, doi:10.1029/2003GB002134, 2004.
- Schaller, N., Mahlstein, I., Cermak, J., and Knutti, R.: Analyzing precipitation projections: a comparison of different approaches to climate model evaluation, *J. Geophys. Res.*, 116, D10118, doi:10.1029/2010JD014963, 2011.
- Schaphoff, S., Heyder, U., Ostberg, S., Gerten, D., Heinke, J., Lucht, W.: Contribution of permafrost soils to the global carbon budget, *Environ. Res. Lett.*, 8, 014026, doi:10.1088/1748-9326/8/1/014026, 2013.
- Schmidt, G. A., Jungclaus, J. H., Ammann, C. M., Bard, E., Braconnot, P., Crowley, T. J., Delaygue, G., Joos, F., Krivova, N. A., Muscheler, R., Otto-Bliessner, B. L., Pongratz, J., Shindell, D. T., Solanki, S. K., Steinhilber, F., and Vieira, L. E. A.: Climate forcing reconstructions for use in PMIP simulations of

- the Last Millennium (v1.1), *Geosci. Model Dev.*, 5, 185–191, doi:10.5194/gmd-5-185-2012, 2012.
- Seppä, H., Bjune, A. E., Telford, R. J., Birks, H. J. B., and Veski, S.: Last nine-thousand years of temperature variability in Northern Europe, *Clim. Past*, 5, 523–535, doi:10.5194/cp-5-523-2009, 2009.
- Shakun, J. D., Clark, P. U., He, F., Marcott, S. A., Mix, A. C., Liu, Z., Otto-Bliessner, B., Schmittner, A., and Bard, E.: Global warming preceded by increasing carbon dioxide concentrations during the last deglaciation, *Nature*, 484, 49–54, doi:10.1038/nature10915, 2012.
- Shannon, S. and Lunt, D. J.: A new dust cycle model with dynamic vegetation: LPJ-dust version 1.0, *Geosci. Model Dev.*, 4, 85–105, doi:10.5194/gmd-4-85-2011, 2011.
- Siddall, M., Henderson, G. M., Edwards, N. R., Frank, M., Müller, S. A., Stocker, T. F., and Joos, F.: $^{231}\text{Pa}/^{230}\text{Th}$ fractionation by ocean transport, biogenic particle flux and particle type, *Earth Planet. Sc. Lett.*, 237, 135–155, doi:10.1016/j.epsl.2005.05.031, 2005.
- Siddall, M., Stocker, T. F., Henderson, G. M., Joos, F., Frank, M., Edwards, N. R., Ritz, S. P., and Müller, S. A.: Modeling the relationship between $^{231}\text{Pa}/^{230}\text{Th}$ distribution in North Atlantic sediment and Atlantic meridional overturning circulation, *Paleoceanography*, 22, 1–14, doi:10.1029/2006PA001358, 2007.
- Sitch, S., Huntingford, C., Gedney, N., Levy, P. E., Lomas, M., Piao, S. L., Betts, R., Ciais, P., Cox, P., Friedlingstein, P., Jones, C. D., Prentice, I. C., and Woodward, F. I.: Evaluation of the terrestrial carbon cycle, future plant geography and climate-carbon cycle feedbacks using five Dynamic Global Vegetation Models (DGVMs), *Glob. Change Biol.*, 14, 2015–2039, doi:10.1111/j.1365-2486.2008.01626.x, 2008.
- Sokolov, A. P., Kicklighter, D. W., Melillo, J. M., Felzer, B. S., Schlosser, C. A., and Cronin, T. W.: Consequences of considering carbon–nitrogen interactions on the feedbacks between climate and the terrestrial carbon cycle, *J. Climate*, 21, 3776–3796, doi:10.1175/2008JCLI2038.1, 2008.
- Stauffer, B., Fluckiger, J., Monnin, E., Schwander, J., Barnola, J.-M., and Chappellaz, J.: Atmospheric CO_2 , CH_4 and N_2O records over the past 60 000 years based on the comparison of different polar ice cores, *Ann. Glaciol.*, 35, 202–208, doi:10.3189/172756402781816861, 2002.
- Steinacher, M., Joos, F., Frölicher, T. L., Bopp, L., Cadule, P., Cocco, V., Doney, S. C., Gehlen, M., Lindsay, K., Moore, J. K., Schneider, B., and Segsneider, J.: Projected 21st century decrease in marine productivity: a multi-model analysis, *Biogeosciences*, 7, 979–1005, doi:10.5194/bg-7-979-2010, 2010.
- Stow, C. A., Jolliff, J., McGillicuddy, D. J., Doney, S. C., Allen, J. I., Friedrichs, M. A. M., Rose, K. A. and Wallhead, P.: Skill assessment for coupled biological/physical models of marine systems, *J. Mar. Syst.*, 76, 4–15, doi:10.1016/j.jmarsys.2008.03.011, 2009
- Taucher, J., Schulz, K. G., Dittmar, T., Sommer, U., Oschlies, A., and Riebesell, U.: Enhanced carbon overconsumption in response to increasing temperatures during a mesocosm experiment, *Biogeosciences*, 9, 3531–3545, doi:10.5194/bg-9-3531-2012, 2012.
- Taylor, K. E.: Summarizing multiple aspects of model performance in a single diagram, *J. Geophys. Res.*, 106, 7183–7192, doi:10.1029/2000JD900719, 2001.
- Thonicke, K., Spessa, A., Prentice, I. C., Harrison, S. P., Dong, L., and Carmona-Moreno, C.: The influence of vegetation, fire spread and fire behaviour on biomass burning and trace gas emissions: results from a process-based model, *Biogeosciences*, 7, 1991–2011, doi:10.5194/bg-7-1991-2010, 2010.
- Thornton, P. E., Doney, S. C., Lindsay, K., Moore, J. K., Mahowald, N., Randerson, J. T., Fung, I., Lamarque, J.-F., Fedema, J. J., and Lee, Y.-H.: Carbon-nitrogen interactions regulate climate-carbon cycle feedbacks: results from an atmosphere-ocean general circulation model, *Biogeosciences*, 6, 2099–2120, doi:10.5194/bg-6-2099-2009, 2009.
- Tschumi, J. and Stauffer, B.: Reconstructing past atmospheric CO_2 concentration based on ice-core analyses: open questions due to in situ production of CO_2 in the ice, *J. Glaciol.*, 46, 45–53, doi:10.3189/172756500781833359, 2000.
- Van Hoof, T. B., Kaspers, K. A., Wagner, F., Van De Wal, R. S. W., Kürschner, W. M., and Visscher, H.: Atmospheric CO_2 during the 13th century AD: reconciliation of data from ice core measurements and stomatal frequency analysis, *Tellus B*, 57, 351–355, doi:10.1111/j.1600-0889.2005.00154.x, 2005.
- Viau, A. E., Gajewski, K., Sawada, M. C., and Bunbury, J.: Low- and high-frequency climate variability in eastern Beringia during the past 25 000 years, *Can. J. Earth Sci.*, 45, 1435–1453, 2008.
- Vicca, S., Gilgen, A. K., Camino Serrano, M., Dreesen, F. E., Dukes, J. S., Estiarte, M., Gray, S. B., Guidolotti, G., Hoepfner, S. S., 5 Leakey, A. D. B., Ogaya, R., Ort, D. R., Ostrogovic, M. Z., Rambal, S., Sardans, J., Schmitt, M., Siebers, M., van der Linden, L., van Straaten, O., and Granier, A.: Urgent need for a common metric to make precipitation manipulation experiments comparable, *New Phytol.*, 195, 518–522, doi:10.1111/j.1469-8137.2012.04224.x, 2012.
- Waelbroeck, C., Paul, A., Kucera, M., Rosell-Melé, A., Weinelt, M., Schneider, R., Mix, A. C., Abelmann, A., Armand, L., Bard, E., Barker, S., Barrows, T. T., Benway, H., Cacho, I., Chen, M.-T., Cortijo, E., Crosta, X., de Vernal, A., Dokken, T., Duprat, J., Elderfield, H., Eynaud, F., Gersonde, R., Hayes, A., Henry, M., Hillaire-Marcel, C., Huang, C.-C., Jansen, E., Juggins, S., Kallel, N., Kiefer, T., Kienast, M., Labeyrie, L., Leclaire, H., Londeix, L., Mangin, S., Matthiessen, J., Marret, F., Meland, M., Morey, A. E., Mulitza, S., Pflaumann, U., Pisias, N. G., Radi, T., Rochon, A., Rohling, E. J., Saffi, L., Schäfer-Neth, C., Solignac, S., Spero, H., Tachikawa, K., and Turon, J.-L.: Constraints on the magnitude and patterns of ocean cooling at the Last Glacial Maximum, *Nat. Geosci.*, 2, 127–132, doi:10.1038/geo411, 2009.
- Willeit, M., Ganopolski, A., and Feulner, G.: Asymmetry and uncertainties in biogeophysical climate–vegetation feedback over a range of CO_2 forcings, *Biogeosciences Discuss.*, 10, 12967–13013, doi:10.5194/bgd-10-12967-2013, 2013.
- Williamson, D., Goldstein, M., and Blaker, A.: Fast linked analyses for scenario-based hierarchies, *J. Roy. Stat. Soc. C-App.*, 61, 665–691, doi:10.1111/j.1467-9876.2012.01042.x, 2012.
- Wohlfahrt, J., Harrison, S. P., and Braconnot, P.: Synergistic feedbacks between ocean and vegetation on mid- and high-latitude climates during the mid-Holocene, *Clim. Dynam.*, 22, 223–238, doi:10.1007/s00382-003-0379-4, 2004.
- Wohlfahrt, J., Harrison, S. P., Braconnot, P., Hewitt, C. D., Kitoh, A., Mikolajewicz, U., Otto-Bliessner, B. L., and Weber, S. L.: Evaluation of coupled ocean–atmosphere simulations of the mid-Holocene using palaeovegetation data from the

- Northern Hemisphere extratropics, *Clim. Dynam.*, 31, 871–890, doi:10.1007/s00382-008-0415-5, 2008.
- Woodage, M. J., Slingo, A., Woodward, S., and Comer, R. E.: U. K. HiGEM: simulations of desert dust and biomass burning aerosols with a high-resolution atmospheric GCM, *J. Climate*, 23, 1636–1659, doi:10.1175/2009JCLI2994.1, 2010.
- Wu, J. and David, J. L.: A spatially explicit hierarchical approach to modeling complex ecological systems: theory and applications, *Ecol. Model.*, 153, 7–26, doi:10.1016/S0304-3800(01)00499-9, 2002.
- Yu, Z., Loisel, J., Brosseau, D. P., Beilman, D. W., and Hunt, S. J.: Global peatland dynamics since the Last Glacial Maximum, *Geophys. Res. Lett.*, 37, L13402, doi:10.1029/2010GL043584, 2010.
- Zaehle, S. and Friend, A. D.: Carbon and nitrogen cycle dynamics in the O-CN land surface model: 1. Model description, site-scale evaluation, and sensitivity to parameter estimates, *Global Biogeochem. Cy.*, 24, GB1005, doi:10.1029/2009GB003521, 2010.
- Zaehle, S., Sitch, S., Smith, B., and Hatterman, F.: Effects of parameter uncertainties on the modeling of terrestrial biosphere dynamics, *Global Biogeochem. Cy.*, 19, GB3020, doi:10.1029/2004GB002395, 2005.
- Zaehle, S., Friedlingstein, P., and Friend, A. D.: Terrestrial nitrogen feedbacks may accelerate future climate change, *Geophys. Res. Lett.*, 37, L01401, doi:10.1029/2009GL041345, 2010.