



# Towards a more objective evaluation of modelled land-carbon trends using atmospheric CO<sub>2</sub> and satellite-based vegetation activity observations

D. Dalmonech and S. Zaehle

Max Planck Institute for Biogeochemistry, Biogeochemical Systems Department, Hans-Knöll-Str. 10, 07745 Jena, Germany

Correspondence to: D. Dalmonech (ddalmo@bgc-jena.mpg.de)

Received: 12 September 2012 – Published in Biogeosciences Discuss.: 15 November 2012

Revised: 23 April 2013 – Accepted: 16 May 2013 – Published: 25 June 2013

**Abstract.** Terrestrial ecosystem models used for Earth system modelling show a significant divergence in future patterns of ecosystem processes, in particular the net land–atmosphere carbon exchanges, despite a seemingly common behaviour for the contemporary period. An in-depth evaluation of these models is hence of high importance to better understand the reasons for this disagreement.

Here, we develop an extension for existing benchmarking systems by making use of the complementary information contained in the observational records of atmospheric CO<sub>2</sub> and remotely sensed vegetation activity to provide a novel set of diagnostics of ecosystem responses to climate variability in the last 30 yr at different temporal and spatial scales. The selection of observational characteristics (traits) specifically considers the robustness of information given that the uncertainty of both data and evaluation methodology is largely unknown or difficult to quantify.

Based on these considerations, we introduce a baseline benchmark – a minimum test that any model has to pass – to provide a more objective, quantitative evaluation framework. The benchmarking strategy can be used for any land surface model, either driven by observed meteorology or coupled to a climate model.

We apply this framework to evaluate the offline version of the MPI Earth System Model’s land surface scheme JS-BACH. We demonstrate that the complementary use of atmospheric CO<sub>2</sub> and satellite-based vegetation activity data allows pinpointing of specific model deficiencies that would not be possible by the sole use of atmospheric CO<sub>2</sub> observations.

## 1 Introduction

The terrestrial and oceanic biospheres currently absorb almost half of the fossil-fuel emissions, and thereby buffer the atmospheric CO<sub>2</sub> increase and reduce the rate of climate change (Cox et al., 2000; Raupach et al., 2008; Le Quéré et al., 2009). Because of the strong interactions between the biosphere net carbon (C) uptake and climate in particular on land (Cox et al., 2000; Friedlingstein et al., 2006; Arora et al. 2013), projections of future climate changes from Earth system models (ESMs) need to accurately simulate the processes that control the evolution of the terrestrial net C balance. However, despite a seemingly common behaviour of C cycle models for the contemporary period, estimates of the future C land balance by different terrestrial biosphere models (TBMs) diverge significantly. This divergence contributes strongly to the overall uncertainty in the future evolution of the global carbon cycle (Friedlingstein et al., 2006; Sitch et al., 2008; Arora et al., 2013). The apparently contradictory behaviour underlines the difficulty of constraining future projections of terrestrial models with current observations. This calls for an in-depth model evaluation that focusses on the model’s capacity to simulate key features of C-cycle-related processes rather than simply ensuring that the easily diagnosed simulated net land–atmosphere C exchange agrees with estimates inferred from observations.

Several global model evaluation analyses have been published in the last decades with respect to land model performances of the carbon cycle (Anav, et al., 2013; Cadule et al., 2009; Blyth et al., 2009; Randerson et al., 2009; Heimann et al., 1998). However, they differ with respect to reference dataset used, selection of the observational traits as well as

their computation, and mathematical formulations used to quantify the data–model mismatch. These differences cause uncertainty when it comes to ranking several land surface models or to analyse the outcome from different evaluation works. Recent model benchmarking initiatives (Randerson et al., 2009; Luo et al., 2012) have therefore underlined the need for the development of a standard set of tests and metrics applicable to any land surface model at different spatial and temporal scale.

In addition to a lack of standards, a key challenge in evaluating global biosphere models comes from the uncertainties in observations. From a perspective of data–model mismatch quantification, given uncertainties in data and observation, operators to link model and data exist. However, data error and structural errors are often not known or provided quantitatively (e.g. Raupach et al., 2005).

This study is an attempt to move toward a more robust and a more objective evaluation framework by defining novel tests/diagnostics and quantitative model performance measures that are robust against these mentioned unquantifiable uncertainties. We first selected a parsimonious number of reference datasets that are as much as possible direct observations. In first instance, upscaled products such as that from Beer et al. (2009) were not used as the fraction of gap-filled information is not quantified. Atmospheric CO<sub>2</sub> and remote sensing data of vegetation activity were selected to take advantage of their spatial and temporal coverage and the complementarity of their information content.

Atmospheric CO<sub>2</sub> measurements and transport modelling that links surface fluxes to these measurements are a valuable approach to evaluate TBMs since the atmospheric CO<sub>2</sub> retains the signature of terrestrial ecosystem response to climate variability (Heimann et al., 1998; Randerson et al., 2009; Cadule et al., 2010). However, atmospheric CO<sub>2</sub> observations alone do not allow inference of the contribution of vegetation and soil components to the observed signal, such that a good fit might hide compensating model errors. Remote-sensing observations of vegetation activity may provide complementary information as they reflect the climate- and disturbance-related seasonal and interannual trends of vegetation greenness (Peñuelas et al., 2009; Richardson et al., 2009).

Rather than comparing average quantities, the analyses presented here analyse how much relevant and robust information, which helps constraining model projections, can be extracted from observations. Hence, we select traits, in particular with respect to vegetation activity, that are based on the information of changes with time, correlations with covariates, and the sign of the changes, as well as based on metrics that are sensitive to difference in sign and phase. Phasing and extent of the climate variability simulated by Earth system models (ESMs) often differs from observed climate because of unforced variability (Deser et al., 2010). To circumvent the resulting mismatch from a direct comparison of ESM simulations and modern observations, and to make key

characteristics of the observations useful for the evaluation of ESMs, priority was given to traits and metrics that describe the relationship between climate variables and carbon cycle processes rather than direct comparison of observed and modelled time series.

The second innovation of our studies is that we impose a lower acceptable model performance measure (baseline benchmark) based on the assumption of a null model, i.e. a model that does not show any trend in the quantity under investigation. This lower boundary for each metric helps to avoid misleading interpretation of the number returned by the scores, and to provide a more informative and intuitively interpretable analysis of the model performance. With respect to the atmospheric CO<sub>2</sub> traits, the aim is to quantify how much information the land surface model adds to the signal of ocean and anthropogenic fossil-fuel emissions and thus to quantify how good the model is relative to the null hypothesis (null model). The working line is thus as follows: the analyses were performed on a seasonal and a de-seasonalized signal to better identify C-cycle patterns and the relationship between C-cycle-related processes and climate variability. As detailed in Sect. 2, we selected several characteristics (traits) of the observational data that are relevant to the biosphere's response to climate variability in terms of terrestrial C cycling patterns. We focussed on the the last three decades (1980–2010) since this is the period with the best data availability (Tables 2 and 3). For the selected tests, a list of comprehensive metrics was selected to quantify model performances according to the information content of identified traits. We then compared this metric to the reference value of the metric obtained according to the baseline benchmark to arrive at a final score for the model.

In Sect. 3 we discuss the potential strengths and limitations of the evaluation framework at the example of the the JSBACH land surface model of the MPI-ESM (Raddatz et al., 2007; Giorgetta et al., 2013) driven by reconstructed meteorology.

## 2 Materials and methods

### 2.1 Observational datasets

#### 2.1.1 Atmospheric CO<sub>2</sub>

Atmospheric CO<sub>2</sub> concentration recorded at remote measuring stations were obtained from the flask data/continuous measurements provided by different institutions (e.g. flask data of NOAA/CMDL's sampling network, update of Conway et al., 1994, Japan Meteorological Agency (JMA), Meteorological Service of Canada (MSC), and many others; see Rödenbeck, 2005). Simulated net land–atmosphere CO<sub>2</sub> fluxes for the period 1980 to 2009 were transported together with estimated net ocean CO<sub>2</sub> fluxes (Jacobson et al., 2007; Mikaloff Fletcher et al., 2006, 2007 – one of the best

available products based on Takahashi ocean dataset and involving several biogeophysical ocean models) and fossil-fuel fluxes (EDGARv.4.0, Olivier et al., 2001, <http://edgar.jrc.ec.europa.eu/faq.php>) by means of an atmospheric transport model (TM) to estimate atmospheric CO<sub>2</sub> record at the measuring stations. For our analysis, we used the TM3 model, version 3.7.22 (Rödenbeck et al., 2003), with a spatial resolution of 4° × 5° and driven by interannually varying wind fields of the NCEP reanalysis (Kalnay et al., 1996).

The model-based time series of CO<sub>2</sub> at the measuring stations were based on sampling simulated CO<sub>2</sub> abundance at the same time in which measurements were available in order to reduce the representation bias. The temporal resolution of CO<sub>2</sub> data is the original resolution as recorded at the monitoring stations (hourly to daily/weekly) and dependent of the specific station.

Stations were selected in order to cover representatively a latitudinal gradient (Table 1). Latitudinal and vertical transport of CO<sub>2</sub> differs among TMs (Yang et al., 2007), but these differences are difficult to quantify and attribute to particular model features (Gurney et al., 2003; Peylin et al., 2005). In remote stations with simple topography, different TMs tend to agree better and are expected to have less error. The selection of monitoring stations takes account of this by including mainly oceanic/island stations as these remote stations have a lower uncertainty and are only marginally influenced by local C sources or sinks (MPI Biogeochemistry, technical reports 5–6: <http://www.bgc-jena.mpg.de/bgc-systems/pmwiki2/pmwiki.php/Publications/TechnicalReports>).

Two estimates of the net land–atmosphere CO<sub>2</sub> flux obtained from inverting the observed atmospheric concentrations using atmospheric transport modelling (hereafter referred to as standard fluxes) were also transported using the same protocol as for the simulated TBM fluxes. These fluxes were taken from the Jena inversion system, which relies on the same TM3 transport model (Jena inversion version 3.7.22, available at <http://www.bgc-jena.mpg.de/~christian.roedenbeck/download-CO2/>, update of Rödenbeck et al., 2003; Rödenbeck, 2005, covering the periods 1996–2008 and 1981–2008, respectively). The standard fluxes were not used to derive an absolute benchmark *sensu strictu* but as reference to compute additional traits as reported in Sects. 2.4.1 and 2.4.5.

### 2.1.2 Vegetation activity datasets

To characterize seasonal and interannual changes in vegetation activity, we rely on two satellite-based products: the SeaWiFS-FAPAR (Gobron et al., 2006a, b), the fraction of photosynthetically active radiation absorbed by vegetation, and the longer GIMMS-NDVI collection *g* ([http://glcf.umd.edu/library/guide/GIMMSdocumentation\\_NDVIg\\_GLCF.pdf](http://glcf.umd.edu/library/guide/GIMMSdocumentation_NDVIg_GLCF.pdf)), which is the normalized difference vegetation index, retrieved from the AVHRR sensor records (Tucker et al., 2005; Beck et al., 2011). Both FAPAR and

NDVI provide a measure of greenness integrating canopy functioning. It has been previously shown that these quantities are nearly linearly related (Myneni and Williams, 1994). The selected FAPAR data were provided as 10-day-aggregated time series from September 1997 until June 2006 at a nominal spatial resolution of 2 km and were used to analyse the seasonal cycle of vegetation activity (Table 2). The GIMMS dataset contains biweekly data at a spatial resolution of 8 km from 1981 until 2006 and was used to estimate long-term changes in vegetation activity (Table 3).

Satellite data were aggregated at the spatial resolution of the TBM, including grid cells that are partially covered by bare soils. With this approach, the aggregated signal indirectly accounts for changes in vegetation activity and density. A simple gap-filling procedure based on 2nd degree polynomial interpolation in time was applied to replace bad-quality flag data. All data were aggregated at the monthly temporal resolution. In the case of GIMMS-NDVI, the maximum value composite (MVC) method was used (Holben, 1986). It is assumed that the process of temporal and spatial aggregation of satellite-based vegetation activity smoothes out noise in the data, and the uncertainty induced by the aggregation might be considered negligible for our purpose. Tropical areas were excluded from the analysis due to the high uncertainty in the interpretation of the satellite signal (Asner and Alencar, 2010) and high uncertainties in NDVI datasets in these regions (Huete et al., 2002; Brown et al., 2006).

### 2.2 The JSBACH model

JSBACH is the land surface model of the Max Planck Institute's Earth System Model (MPI-ESM) (Raddatz et al., 2006; Giorgetta et al., 2013) In this study we use the version that was used for the CMIP5 activity (JSBACH version 2.0). JSBACH considers 11 plant functional types, which occupy annually varying fractions (tiles) of a model grid cell, prescribed from land-use data (see Sect. 2.2). Phenology and C cycling is simulated explicitly for each tile, while the half-hourly fluxes of energy and water are calculated for each grid cell, based on the relevant average properties of vegetation and soils across the tiles. The land-use emissions are computed according to the method reported in Reick et al. in review. JSBACH is applied here in offline mode, i.e. driven by reconstructed daily meteorology (see Sect. 2.3), at the same spatial resolution of the CMIP5 simulations of the MPI-ESM (T63, corresponding to a 1.875° × 1.875° resolution at the equator).

### 2.3 Climate and land-use forcing

Meteorological forcing data (air temperature and humidity, shortwave and longwave incident radiation, precipitation, and surface wind speed) for 1860 to 2010 were derived from CRU-NCEP (CRU-NCEPv4Viovy, N. 2011, available from <http://dods.extra.cea.fr/data/p529viov/cruncep/>),

**Table 1.** List of selected atmospheric CO<sub>2</sub> monitoring stations and satellite-based vegetation activity datasets used in the analyses, as well as the time period used for elaborations.

Label	Name	Lat. (degree)	Lon. (degree)	Years of elaboration
ALT	Alert, Canada	82.45	−62.52	1982–2008
BRW	Point Barrow	71.32	156.6	1982–2008
STM	Station ‘M’, Atlantic	66	2	1982–2008
CBA	Cold Bay, Alaska	55.2	162.72	1982–2008
SHM	Shemya Island, Alaska	52.72	174.1	1985–2008
MHD	Mace Head, Ireland	53.33	9.9	1991–2008
AZR	Azores	38.75	27.08	1995–2008
KEY	Key Biscayne, Florida	25.67	−80.2	1982–2008
MLO	Mauna Loa, Hawaii	19.53	−155.58	1982–2008
KUM	Kumakahi	19.52	−154.82	1982–2008
GMI	Guam, Mariana Island, Pacific	13.43	144.78	1996–2008
RPB	Ragged Point, Barbados	13.17	−59.43	1987–2008
CHR	Christmas Island	1.7	−157.17	1982–2008
SEY	Mahe Island, Seychelles	−4.47	55.17	1996–2008
ASC	Ascension Island	−7.92	14.42	1982–2008
SMO	Tutuila, American Samoa, Pacific	−14.25	−170.57	1982–2008
PSA	Palmer station, Antarctica	−64.92	−64	1982–2008
HBA	Halley Bay, Antarctica	−75.67	−25.5	1996–2008
SPO	South Pole	−89.98	−24.8	1982–2008
GIMMS	GIMMS-NDVI	–	–	1982–2006
SW	SeaWiFS-FAPAR	–	–	1998–2005

and were aggregated via conservative regriding to the T63 resolution of the MPI-ESM grid at daily resolution. These data were used as model forcing as well as for the climate correspondence analysis. The standardized precipitation index SPI was computed from the precipitation record of the CRU observational dataset (Mckee et al., 1993; Lloyd-Hughes and Saunders, 2002). SPI is suitable as indicator of both dry and wet soil conditions. Irrespective of biomes or region, the 6-month cumulated precipitation data was used to compute the SPI for each grid cell (see Appendix A for more details). Land-cover and land-use change transition maps were derived from Hurtt et al. (2006).

## 2.4 Evaluation methodology

The analyses in this study focus on seasonal and interannual/decadal time scales. To identify these components from the observed and simulated atmospheric CO<sub>2</sub>, as well as vegetation activity and climatic drivers, a seasonal component (up to annual time scale) and an interannual time scale component were isolated using a filter implemented in the Fourier space. We followed the method and the cut-off values presented in Thoning et al. (1989), using Gaussian spectral weights (Rödenbeck et al., 2003). The outcome of the filtering is (i) a seasonal component with a mean of zero, which retains information up to the annual frequency with the very high frequency (daily to biweekly) removed, and (ii) a de-

seasonalized signal, which includes all the frequencies lower than the annual cycle – i.e. the interannual to decadal time scales. In terms of interannual variability, this approach of filtering is more advantageous than consideration of monthly anomalies since a de-seasonalized signal provides a better measure of the strength and persistence of interannual variability related to climatic and natural events as El Niño events and volcanic eruptions.

The analysis of seasonal patterns aims not only at the relative phasing of vegetation growth and ecosystem respiration and modelled phenology that affects the seasonal phasing of the net land–atmosphere C exchange (Prentice et al., 2000) but also at biogeophysical effects such as the water and energy exchanges (Notaro et al., 2007; Peñuelas et al., 2009). Interannual variability and long-term trends of net land–C exchanges and vegetation activity are an important and crucial aspect of the terrestrial ecosystem in a climate change context. Changes of vegetation activity might have implications to long-term potential for retaining more C in the system, contributing hence to the biosphere–atmosphere feedbacks and internal plant–soil feedbacks (Bonan, 2008).

In the following sections, we describe key features of the atmospheric CO<sub>2</sub> and vegetation activity obtained from the decomposed signals (Table 2: seasonal time scales; Table 3: interannual time scale). These traits are used to assess the capacity of the model to reproduce climate-variability-induced

effects on terrestrial ecosystems. In addition, traits characterizing the co-variability of vegetation features/atmospheric CO<sub>2</sub> and land climatic patterns are defined. Some of the selected traits were analysed separately in three time intervals (1982–1991, 1992–1997, and 1998–2006) according to two breakpoint events: the Mount Pinatubo eruption in 1991, and the El Niño event in 1997 – two of the most relevant natural events occurred in the last three decades.

The systematic quantitative assessment of the correspondence of anomalies and trends in simulated vegetation activity and net C exchange is performed using normalized metrics (see Appendix B for the mathematical description). The proposed selected traits and metrics are suitable to be applied to land surface models run in either offline or fully coupled mode because they are based on reproducing variability and/or statistical relationships with the driving climate rather than focusing on the absolute correspondence of the variables. This strategy reduces potential biases in the assessment due to uncertainty in the predicted climatic variability (Deser et al., 2010).

Geographical regions at the continental scale consistent with the regions used for the Transcom3 project (Gurney et al., 2002; Fig. 1) were used to determine the influence of net land–atmosphere CO<sub>2</sub> fluxes from a particular region to the signal at the monitoring stations following the procedure reported in Cadule et al. (2010). The characterization of vegetation activity was performed at grid-cell level and at regional level according to the same Transcom3 regions. The Transcom3 region maps were further intersected with the dominant vegetation map obtained from the Synmap vegetation classification of Jung et al. (2006) (see Appendix D). Grid cells with dominance of bare soil or ice, as well as grid cells with no valid observations, were excluded from the analyses.

#### 2.4.1 Seasonality of Atmospheric CO<sub>2</sub>

The model's capacity to simulate phase and amplitude of the mean seasonal cycle of atmospheric CO<sub>2</sub> (MSC) was evaluated using the Taylor score (Taylor, 2001). The selected metric gives more weight to the correspondence in phase instead of amplitude (Taylor, 2001), which is the more reliable feature of transport models (Stephens et al., 2007). Additional information on the land net C exchange is contained in the latitudinal gradient of the amplitude of the mean seasonal cycle (MSClg), which increases from the South Pole northwards because of the relatively higher land masses fraction in the Northern Hemisphere (NH). A metric based on the variance of the amplitude data was used to assess the model performance (Table 2 and Appendix B).

The relative contribution of the C fluxes from land (and ocean) Transcom3 regions to the seasonal cycle amplitude (MSCc) was computed using the atmospheric CO<sub>2</sub> record obtained by transporting the standard fluxes constrained on the period 1996–2008 as reference. This choice was made so

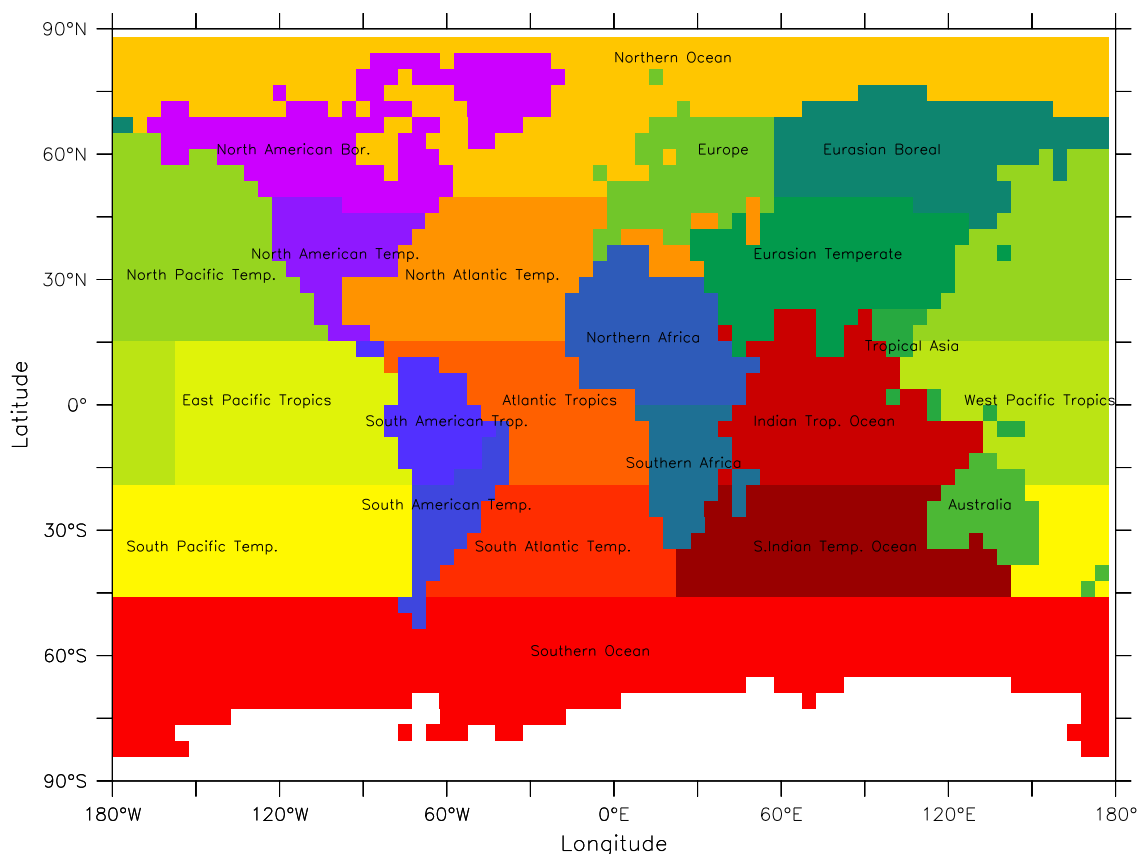
as to overlap with the time period for which the SeaWiFS-FAPAR data are available (see Sect. 2.4.2). The relative contribution of each region to each single monitoring station in both standard fluxes and modelled fluxes was compared using the Pearson correlation coefficient. This trait checks thus also for the existence of potential inconsistencies between the regional and seasonal distribution of net land-C fluxes from the model and estimated by the inversion of atmospheric observations.

Changes in the seasonal cycle over time, referred to as the monthly CO<sub>2</sub> trend (MT), are quantified as the year-to-year change in CO<sub>2</sub> concentration for each month. Previous works analysed solely the change in amplitude of the seasonal cycle in Mauna Loa as response to land surface warming (Myneni et al., 1997; Angert et al., 2005; Buermann et al., 2007), while we focus on decadal trends in long-term northern stations, which exhibit a clearer signal. This trait summarizes the seasonal change in the trend of land-C sink/sources in response to climatic drivers and natural disturbances in the extratropical latitudinal band. The model–data correspondence is analysed using the Pearson correlation coefficient.

The trend in the seasonal onset of net land-C uptake (C-dd) was computed as follows: for each year, the algorithm looks for the downward zero-crossing point of the seasonal time series of atmospheric CO<sub>2</sub>. The trend is thereafter computed on the extracted dates. This feature characterizes in particular the observed high-latitude ecosystem responses to recent land surface warming and it is indirectly linked to the beginning of the growing season (Keeling et al., 1996; Myneni et al., 1997). Because the years 1991–1993 – i.e. the years following the Mount Pinatubo eruption – are an anomaly in this trend (Lucht et al., 2002), these three years were excluded from the analysis. The analyses for the MT and C-dd traits focus on the stations in the extratropical latitudinal band with a clear signal from land and low contamination of the trends due to uncertainties in the fossil-fuel emissions (Table 2).

#### 2.4.2 Seasonality of vegetation activity

A direct comparison of absolute values of remote sensing data such as NDVI or FAPAR and their corresponding modelled variable might be not a viable strategy, first and foremost because of different retrieval and post-processing algorithms used to compute the final estimated FAPAR/NDVI in different satellites products, and to remove, for example, cloud contamination and atmospheric corruption, etc. (e.g. the intercomparison study of Dahlke et al., 2013). This implies that the outcome of a direct model–data comparison is dependent on the reference dataset used. In addition, the radiances recorded by satellites differ in the way that radiation extinction is computed at the land surface in land surface models. This difference does not allow a priori for a perfect match between data and model.



**Fig. 1.** Map of the land regions used for the regional benchmark of phenology and the analysis of the biosphere fluxes, as defined in the TransCom intercomparison studies (Gurney et al., 2002). The map shows the regions at TM3 resolution. Code: North American Boreal (NAB), North American Temperate (NATe), South American Tropical (SATr), South American Temperate (SATE), Northern Africa (NA), Southern Africa (SA), Eurasian Boreal (EAB), Eurasian Temperate (EATe), Tropical Asia (TrA), Australia (AUS), Europe (EUR). The ocean was considered as a single region.

However, as shown in Dahlke et al. (2013) for the seasonal information, the temporal evolution of the recorded signal is likely to be a robust feature among datasets and the temporal evolution of the modelled signal should resemble the reference dataset such that they can be evaluated by a metric that is independent from the absolute values of the time series. Because of the aforementioned reasons, we focused on metrics based on information on time and sign of changes as indicated by the satellite data.

With respect to the seasonal signal, as a first step, grid cells with only one detected growing season per year were selected by analysing the autocorrelation of the seasonal record and its significance. The shape of the seasonality of vegetation activity was then characterized by two robustly identifiable and meaningful phases of the phenological cycle: the time of the beginning of the vegetative growing season, hereafter referred to as time of onset ( $t$ -onset), and the time of the maximum FAPAR signal ( $t$ -max) (Randerson et al., 2009). Data and model signals characterized by mean amplitude of the seasonal record within 1 % of total FAPAR range were ex-

cluded from the analyses. The definition of the beginning of the growing season is a subjective matter and a direct and precise link to ground-level observation is difficult to identify (Lucht et al., 2002; Maignan et al., 2008; Verstraete et al., 2008). Analogously to the method of estimating the beginning of the net CO<sub>2</sub> uptake reported in Sect. 2.4.1, the proxy of the time of onset of vegetation activity is calculated on the seasonal signal, and corresponds to the point in time of the upward zero-crossing point of the seasonal curve (see Fig. A1).

Linear differences of the most frequent month of time of onset or maximum of FAPAR were computed between model and data. Consequently this metric ranges between one (no difference) to zero (6-month difference). The length of the growing season was not used as additional trait because it is poorly defined from satellite data as autumnal leaf colouring and the simultaneous presence of living and dead leaves confounds the satellite signal, in particular in temperate regions (Estrella and Menzel, 2006; Menzel et al., 2006).

### 2.4.3 Interhemispheric gradient and trend of atmospheric CO<sub>2</sub>

The long-term trend in atmospheric CO<sub>2</sub> (C-LTT), given known fossil fuel, land-use change emissions and net ocean carbon fluxes, is an indication of the long-term net C balance of the terrestrial biosphere (Prentice et al., 2000; le Quere et al., 2009). The trend was computed from the mean annual values of the de-seasonalized signals and compared directly to the observations for stations covering the period 1982–2008 (Table 1). The interhemispheric gradient in atmospheric CO<sub>2</sub> abundance (IHG) measures the north–south differences in atmospheric CO<sub>2</sub> caused by changing balance of the increasing fossil-fuel emissions in industrialized regions and the net ocean and land-C uptake. For each year, this trait was computed by subtracting the observed and modelled annual CO<sub>2</sub> concentration at the South Pole station (SPO) from the respective station concentrations, as in Cadule et al. (2010). The metric was based on the comparison of the standard deviation of modelled and observed data.

### 2.4.4 Trend of vegetation activity

Similar to atmospheric CO<sub>2</sub>, vegetation activity trends were computed from modelled and reference data. Beck et al. 2011 indicated that the GIMMS-NDVI dataset is suitable for assessing temporal changes of vegetation activity. However, due to the unknown uncertainty of the absolute NDVI values, the selected trait does not compare numerical trends. Instead, the selected metric focusses on the robust trends in the data and determines the spatial patterns of positive, negative, or no significant trend in vegetation signal from the GIMMS-NDVI dataset and compares this to the pattern in modelled FAPAR (Table 3). For each grid cell, the metric calculation was performed on annual values of the de-seasonalized vegetation time series. The non-parametric Mann–Kendall test was used to determine whether a positive (greening), negative (browning) trend or no significant trend was detected (two-tailed statistic). The advantage of this approach is that it is robust against satellite drift and high-model internal variability that is, for instance, induced by high variability in the climate simulated by an Earth system model. At the grid-cell level, the metric is a binary score which measures whether the model and data show a significant trend of the same sign. The global-scale metric is then a ranking of a percentage agreement for cells of a particular trend class.

### 2.4.5 Quantification of interannual variability: atmospheric CO<sub>2</sub> and vegetation activity relationship with land climate pattern

The relationship between the seasonality of phenology and local climatic drivers at grid-cell level was explored using the annual variations of the time of beginning of the growing season (*t*-onset; Table 2). The time series for the SeaWiFS-

FAPAR data is too short to allow for a trend analysis. Therefore the correlation of the *t*-onset with the annual temperature, given the annual SPI as conditional variable, was taken as a proxy. A ranking metric, analogous to the vegetation activity trend metric, was computed according to cell-by-cell agreement in terms of sign of the statistics, hence according to significantly positive, negative, or non-existent correlation.

Interannual variability in vegetation activity was assessed using de-seasonalized signals obtained from the GIMMS-NDVI/modelled FAPAR aggregated to the Transcom3 land region. Cross correlations between monthly records of vegetation activity and regional climatic variables, temperature and SPI, were computed with lags up to 24 months (Table 3). The South American Tropical region, Tropical Asia regions, and grid cells with dominance of tropical forests in Africa are excluded by the analysis (see Sect. 2.1.2).

The same approach was used to measure the relationship between atmospheric CO<sub>2</sub> growth rate and land surface climate (Table 3). The atmospheric CO<sub>2</sub> growth rate is well known to provide information on the interannual variability of the biospheric response to climate variability and in particular land response at the ENSO time scale (Keeling et al., 1995; Le Quere et al., 2003; Peylin et al., 2005). However, most of the land surface climate shows some coherence with this large-scale climatic feature (Buermann et al., 2003), such that the CO<sub>2</sub> signal in the atmosphere could be perfectly correlated, instantaneously or lagged, with climate over most of the land regions. To reduce this problem, an empirical orthogonal function (EOF) decomposition of the atmospheric CO<sub>2</sub> records, obtained by transporting the “inverted fluxes” from each land region, was computed. The three most contributing land regions (to at least 80 % of the variability in the observed total signal) for selected monitoring stations were determined and only these were used in the analysis (see Appendix C).

The obtained statistically significant cross correlations from data and model (vegetation and atmospheric CO<sub>2</sub> growth versus regional climate) were compared with a correlation metric in order to test if the model is able to return the coupled patterns with time lags (see Appendix C).

The use of inverted fluxes to determine the most contributing regions at interannual time scale and for the EOF decomposition does not affect significantly the results in terms of model behaviour evaluation. However, it changes the degree to which the observations can effectively constrain the model if in the model domain a region contributes less than inferred from the inverted fluxes.

The last selected feature of the carbon cycle uses the CO<sub>2</sub> growth rate to compute an apparent land-C cycle sensitivity to global temperature anomalies, defined as the slope of the annual CO<sub>2</sub> growth rate versus the aggregated annual land surface temperature. The record at the station of Mauna Loa (MLO) was used as proxy of evolution of globally averaged atmospheric CO<sub>2</sub> concentration (Zeng et al., 2005).

## 2.5 The baseline benchmark and the final scores

The reference minimum (baseline benchmark) concept applied in this study compares the skill of the model under investigation with the score of the metric obtained assuming a land biosphere that does not systematically contribute to any signal. For the C-cycle analyses, the baseline benchmark is set to be a biosphere without a terrestrial C-cycle ecosystem, implying that the signal or trend in the observations is driven by fluxes of fossil fuel and net ocean only (no-land case). Since this lower benchmark is applied based on the same TM for all the simulations, this further reduces the potential errors introduced by transport modelling uncertainties. Scaling the metric to the lower benchmark highlights the contribution of modelled land fluxes to match the observed trait of the data under consideration. In other words, the final score number is the metric for an individual trait, cleaned by the contribution of other CO<sub>2</sub> source/sink other than the modelled land fluxes. Only for the CO<sub>2</sub> drawdown test (C-dd; Table 2) is the baseline benchmark set as zero trend (i.e. there is no trend on land). A similar concept is applied for the vegetation activity traits: the lower benchmark is provided by the case with constant vegetation (no-change case). Only in the case of timing of the vegetation onset and maximum (*t*-onset and *t*-max traits; Table 2) is the baseline benchmark set as the maximal possible difference (6 months).

The final global model metrics *M* for each trait are computed as follows: first, the metrics are computed for the CO<sub>2</sub> signal in each monitoring station and in correspondence of each Transcom3 region for the vegetation-related traits (*M*<sub>or</sub> in Eq. 1). The same statistic is also applied for the null-model case to return the numerical metric value of the trait for the baseline benchmark case (*M*<sub>base</sub> in Eq. 1). The original metric is then scaled to a new, normalized metric (score) between 0 and 1 according to Eq. (1), where 1 indicates perfect data–model match and 0 indicates that the model is not able to perform better than a system without the representation of the land biosphere.

$$M = \frac{M_{or} - M_{base}}{1 - M_{base}} \quad (1)$$

Secondly, the model performances are summarized in a polar plot that goes radially from 0 (less skillful model), in the centre, to 1 (skillful). The global scores are derived as follows: for the satellite-based scores, the global score is the average of the scores computed for each Transcom3 region, with the exception of the ranking-based scores, which are already computed at global scales. For the CO<sub>2</sub>-station-based scores, the scores for each station were first averaged by latitudinal band, and the global score was then derived as the average of the scores computed by latitudinal band.

## 3 Results and discussion

In the following, we discuss the results of the above framework at the example of the JSBACH model. The results for the individual traits are summarized in Fig. 2. Table A1 reports the results of the baseline benchmarking for comparison. Table 4 reports results per latitudinal band with regards to CO<sub>2</sub> traits, and global scores for the vegetation traits.

In this section an in-depth analysis of the mechanisms behind data–model mismatch is not performed, but what we can learn from observations and how can we use them to quantify data–model differences is shown, as well as what the benchmarking framework can tell about potential areas of model deficiencies.

### 3.1 Seasonality of atmospheric CO<sub>2</sub> and vegetation activity (Table 2)

#### 3.1.1 Seasonality of atmospheric CO<sub>2</sub>

The Taylor diagram (Fig. 3a) reports the data–model correspondence in terms of phase and amplitude of the mean seasonal cycle (MSC). JSBACH is in general capable of simulating the phase of the seasonal cycle of CO<sub>2</sub>, with the exception of the stations south of the equator, which tend to be out of phase. At those stations, ocean fluxes dominate the signal, which can be seen in the large difference between the low original and the higher scaled metric (Table 4). The anti-correlation of the model's seasonality might further indicate either (or both) a high contribution of the signal from the Northern Hemisphere or reveal effective out-of-phase seasonal land-C fluxes. The in-depth analysis of the regional contribution to the mean seasonal cycle (the MSCc trait) indicates that the Eurasian Boreal and Eurasian Temperate regions slightly but systematically contribute more to the signal in the stations above the 50° N than inferred from observations (Fig. A2). At the southern stations, the model signal from the South American Temperate region clearly dominates the ocean signal (Fig. A2), suggesting that this region has a seasonal cycle of net land–atmosphere C fluxes inconsistent with the atmospheric record. This inconsistency leads to the low scores in the southern latitudinal band (Fig. 2, Table 4).

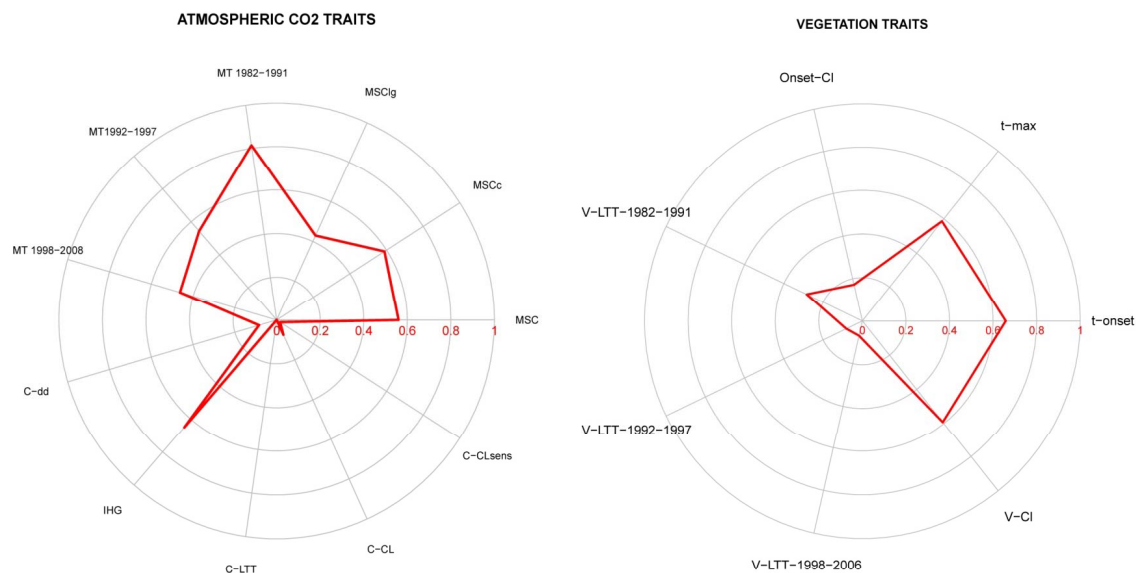
The model clearly overestimates the amplitude of the MSC across the global network of stations, as can be seen in the latitudinal gradient of the amplitude of the mean seasonal cycle (Fig. 3b). Although uncertainties in the transport model could partially contribute to this, the steep drop of the CO<sub>2</sub> concentration during the summer months (data not shown) are an indication that an overestimation of spring C uptake (i.e. too large global gross primary productivity) is responsible for the overestimation of the amplitude.



**Table 2.** List of atmospheric CO<sub>2</sub> and vegetation activity traits used for the analyses at the seasonal time scale. A detailed explanation of metrics can be found in Sect. 2 and Appendices B and C).

Seasonal time scales					
CO <sub>2</sub> Trait	Label	Test	Metric	Section	
				Methods	Results
Mean seasonal cycle	MSC	centred pattern variability	Taylor (2001) Eq. (B1)	2.4.1	3.1.1
Regional contribution to mean seas. cycle	MSCc	relative contribution	Pearson correlation $r$	2.4.1	3.1.1
Latitudinal gradient of MSC amplitude	MSClg	latitudinal pattern of amplitude	standard-deviation-based metric Eq. (B2)	2.4.1	3.1.1
Monthly CO <sub>2</sub> trend* (1982–91/1992–97/1998–2008)	MT	phase of the monthly pattern	Pearson correlation $r$	2.4.1	3.2
CO <sub>2</sub> drawdown points*	C-dd	direct comparison of numerical trend	single value comparison metric Eq. (B4)	2.4.1	3.5.1
Vegetation Trait	Label	Test	Metric	Methods	Results
Time of onset of phenology	$t$ -onset	most frequent month	absolute difference Eq. (B3)	2.4.2	3.1.2
Time of maximum activity of phenology	$t$ -max	most frequent month	absolute difference Eq. (B3)	2.4.2	3.1.2
$t$ -onset $\sim$ drivers	Onset-CL	occurrence of positive/negative/no correlations	spatial ranking	2.4.5	3.5.1

\* Trait applied to the stations ALT, BRW, STM.

**Fig. 2.** Global atmospheric CO<sub>2</sub> and vegetation activity scores for the JSBACH model according to the list of traits in Tables 2 and 3. The polar plot goes radially from 0 (less skillful model), in the centre, to 1 (skillful). Since we only consider one model here, we refer to the threshold value of 0.5 to indicate the model good/high performances and less good/low performances.

### 3.1.2 Seasonality of vegetation activity

Figure 4 shows that JSBACH simulates the time of onset with a systematic lag of 1 to 2 months over large areas of the Northern Hemisphere (NH). A major exception is the east

and south of the North American Temperate region, where the model tends to lead the observed growing season. Given the monthly temporal resolution of the analyses, these results in the NH are still in line with the good performance in terms of phase correspondence of the MSC of CO<sub>2</sub> at the northern

stations (Fig. 3a). However, results indicate that there is space to improve the modelled phenology.

In large parts of the tropical latitudinal band, most of the modelled signal is flat, in contrast to the detected seasonal cycle recorded in the SeaWiFS data in seasonally dry tropical areas (Fig. 4). In these areas, which are dominated by rain-deciduous vegetation, the occurrence of one growing season is driven by seasonality in rainfall. Similarly, the model signal in the Australian scrubland does not show any clear seasonality in contrast of the observations. The flat tropical signal and the detected differences up to 3–4 months in some southern regions are responsible for the low aggregated global model performance (Fig. 2, Table 4). The vegetation classes contributing most to the lower performances are deciduous broadleaved forests and grassland, probably mostly due to their geographical distribution and presence in drought-prone areas (see Fig. A3).

The timing of the maximum analysis ( $t$ -max; data not shown) returns similar geographical pattern to the “ $t$ -onset” trait, but the differences are generally slightly higher. This aspect partly relates to the less well defined nature of the timing of the maximum in regions with several months of full foliar coverage. At the global scale, there are no discernible differences between the two scores (Fig. 2, Table 4). These results show that the seasonality in the model is slightly lagged in time, but without strong distortions in the signal in the first period of the growing season in the Northern Hemisphere. An improvement in phenology parameterization in areas dominated by raingreen vegetation in the seasonally dry tropical latitudinal band and drought-prone shrublands is necessary. It is unclear from the analysis, however, whether a too low sensitivity of raingreen vegetation to soil moisture stress or a too low seasonal cycle of simulated soil moisture as a result of problems with the modelled soil hydrology is the cause of this phenomenon.

### 3.2 Monthly CO<sub>2</sub> trend

As example for the trend in the monthly CO<sub>2</sub> signal (MT), Fig. 5a displays the trend computed for observed signal at the Alert station (ALT) together with the contributions of the net land and ocean C fluxes and fossil-fuel emissions. For the selected northern stations, the observational analysis shows that, in particular in the summer months (June–July), the land is the most dominant contributor to the tendency towards a more pronounced seasonal cycle. That is to say, increased monthly land-C uptake rather than changes in ocean fluxes and fossil-fuel emissions are responsible for this trend. This feature is particularly strong in the period 1982–1991 and consistent across the selected stations, although this trend is not always statistically significant for all the months (Fig. 5a). The monthly CO<sub>2</sub> trend in the period 1992–1997 is less clear (data not shown), while the negative trend of the summer uptake occurs in that of 1998–2006, albeit weaker than for 1982–1991. The latter pattern likely

reflects the weakening of the positive land warming effect on phenology during the growing season, which was particularly apparent in the 1980s (Myneni et al., 1997).

Using an additional TM simulation, we verified that the observed weakening of the negative trend in summer is indeed mainly land induced and not induced by the interannual wind fields used in the transport model. The experimental results with constant wind (data not shown) confirmed that interannually varying transport can contribute but does not overwhelm the land-based trends in monthly CO<sub>2</sub> concentrations. Potential trends in the seasonality of fossil-fuel emissions (Blasing et al., 2005) are unlikely to strongly affect this trend (data not shown).

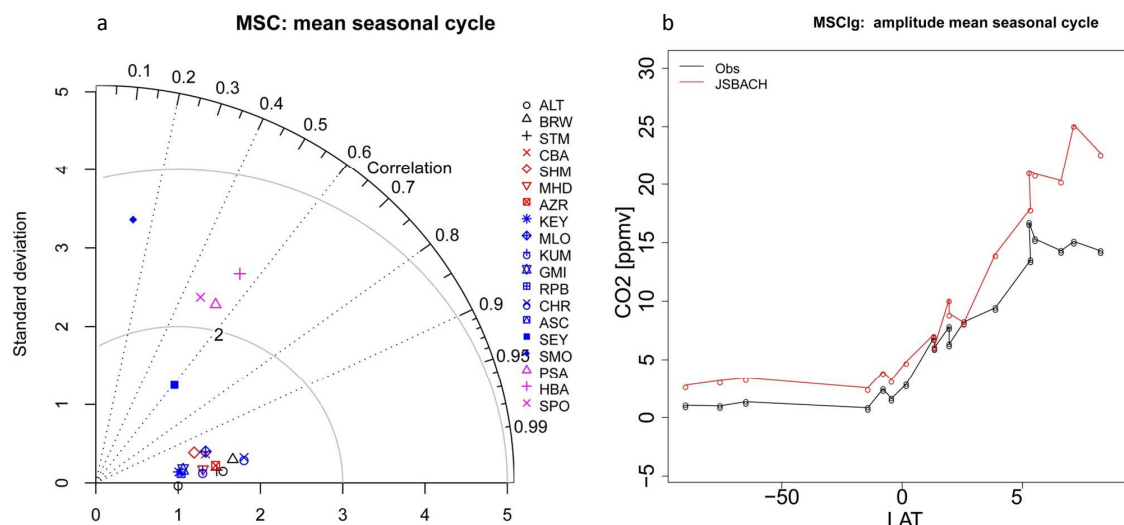
Figure 5b exemplarily shows that JSBACH is able to qualitatively return the seasonal-like shape of the monthly CO<sub>2</sub> trend and the detected land-C uptake weakening, but it is not able to fully explain the observed signal (Fig. 2 and Table 4). Since the selected metric analyses the correspondence of phase of the monthly trend, the non-perfect match could be attributable to divergence in observed and modelled climate sensitivities of photosynthesis and respiration.

### 3.3 Interhemispheric gradient and long-term trend of atmospheric CO<sub>2</sub> (Table 3)

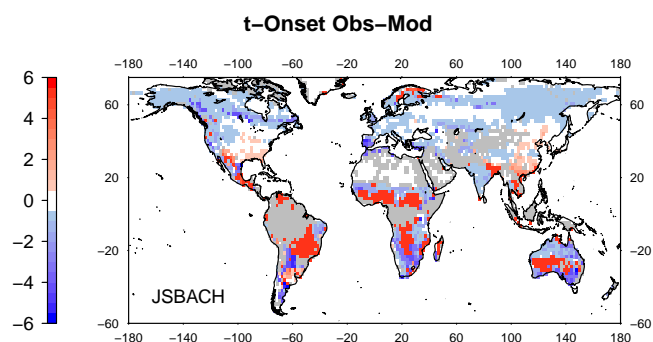
The interhemispheric gradient trait (IHG), which evaluates the interannual variability of the net land–atmosphere C exchange, agrees well between JSBACH and the observations (results not shown, but see Fig. 2). However, the analysis on the long-term C balance trend (C-LTT) shows that JSBACH substantially overestimates the long-term trend compared to observation (Fig. 6a), such that its score is actually lower than the baseline benchmark at all stations (Fig. 2, Table 4, Table A1). Since this detected data–model difference is unlikely to be due to uncertainties in fossil-fuel emissions or ocean net carbon fluxes (le Quere et al., 2009), this result is due to a substantial underestimation of net land-C uptake.

### 3.4 Vegetation activity trend (Table 3)

Figure 6b displays the decadal patterns of the normalized annual vegetation activity time series (GIMMS-NDVI and JSBACH-FAPAR), excluding evergreen tropical forests, glaciers, and desert areas. There appears to be a good qualitative global agreement, suggesting that phenological limitations are not likely the cause for the aforementioned too low increase in land C. However, the good agreement of the global vegetation pattern is partly due to the compensation of errors (Fig. 7). The observed, spatially extensive positive trend in vegetation greenness in the period 1982–1991 is not fully captured by the model because several areas have either no trend or even a negative trend (in parts of the South America, Australia, and South East Asia). During the years 1992–1997, no clear geographical pattern is detected (data not shown). For the years 1998–2006, large areas with an



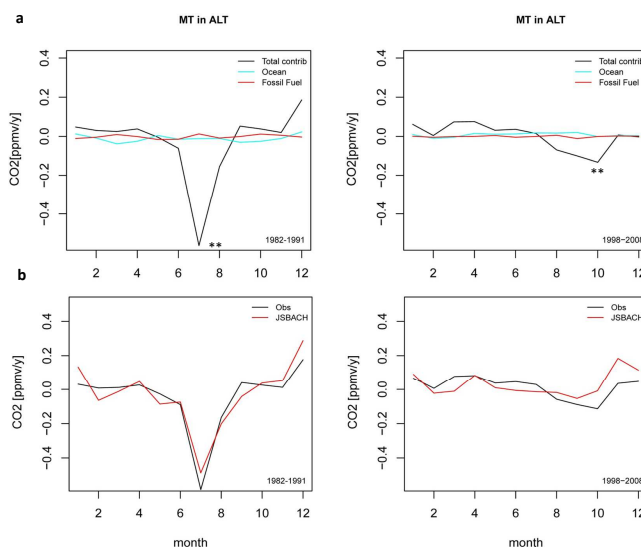
**Fig. 3.** (a) Taylor diagram of the mean seasonal cycle of JSBACH. (b) Latitudinal gradient of the amplitude of the mean seasonal cycle. The x-axis of the Taylor diagram indicates a mismatch in terms of amplitude and the y-axis provides information in term of phase correspondence. Stations in the coloured list are sorted according to latitude: black, 90° N–60° N; red, 60° N–30° N; blue, 30° N–30° S; pink, 30° S–90° S.



**Fig. 4.** Difference between the most frequent month of time of onset for SeaWiFS-FAPAR data and modelled FAPAR (expressed as months) during 1998–2005. Grey areas were masked out from the analysis and indicate missing observations, dominance of tropical evergreen rain-forests, desert or glaciers, or areas with more than one growing season. Red cells indicate missing data from the model.

observed positive trend in the period 1982–1991 appear to have no or even negative trends. This phenomenon is only partly reproduced by JSBACH: in the northern boreal regions and in the Southern Hemisphere, particularly in the South American Temperate region, the negative trends are simulated.

The observed large-scale positive trends in vegetation activity during the period 1982–1991 is consistent with previous results (Myneni et al., 1997; Zhou et al., 2003). However, our analysis underlines that the observed positive warming effect on greening has not been persistent in time, but switched toward a neutral effect in the years 1992–1997 and a localized negative trend in the years 1998–2006. The ob-



**Fig. 5.** Monthly CO<sub>2</sub> trend in the station of Alert (Canada, ALT) for the periods 1982–1991 and 1998–2008. (a) Observations, as well as simulated contribution from fossil-fuel emission and net ocean fluxes (\*\*  $P < 0.01$ ). (b) Monthly record for observations and modelled data. Negative values for a specific month indicate a decrease of seasonal atmospheric CO<sub>2</sub>, indirectly linked to an increase of biosphere C uptake, and vice versa.

served negative pattern in the SH is generally consistent with the trends in evapotranspiration and in particular soil moisture reported in Jung et al. (2010) even though our analyses ends in 2006, while theirs ends in 2008. Several factors might contribute to the observed overall behaviour following the El Niño event in 1997. These include recurrent drought events,

pest outbreaks, and severe fire events over several regions responsible for the detected negative trends in boreal areas and the weakening of the summer C uptake that we reported in Sect. 3.2 (van der Werf et al., 2004; Angert et al., 2005; Goetz et al., 2005).

The low final score of JSBACH in this metric (Fig. 2, Table 4) is in particular the result of the recurrent large-scale negative trends in several areas in the SH and in south-east Asia during the years of 1982–1991 and 1998–2006 (Fig. 7). The non-quantitative nature of this comparison prohibits a too strict interpretation of the mechanisms behind the model–data differences. It is unclear whether these differences are caused by the phenological scheme of the model, land-use change protocol, or other factors such as the drought response or fire processes. However, the disagreement in the sign of the trend can be attributed to model deficiencies, and the ranking metric provides a quantitative measurement of the detected disagreement.

As aforementioned, despite the spatial model–data disagreement, at global scale the errors in the model compensate to return a positive vegetation activity trend. Assuming that vegetation activity is linked to plant productivity, the underestimation of the net land-C uptake in JSBACH (Sect. 3.3) is likely the consequence of a too high soil-C turnover rate.

### 3.5 Terrestrial ecosystems and climate variability

#### 3.5.1 Growing season response to local climate (Table 2)

The timing of the CO<sub>2</sub> drawdown point (C-dd) and the onset of vegetation greening (*t*-onset) represent two independent proxies to measure the effects of land warming on spring phenology (Badeck et al., 2004; Menzel et al., 2006). There is a tendency towards earlier CO<sub>2</sub> drawdown at the stations of STM, BRW, and ALT (Fig. 8a), although this trend is statistically significant only for the latter two stations ( $P < 0.10$ ). Such a negative trend in time is consistent with the advance of spring phenology induced by land surface warming (Fig. 8b): the correlation between climate variability and the timing of vegetation onset is significantly negative with annual temperature. Despite this trait constituting an emerging empirical relationship, the negative correlation mainly in the boreal areas, as clearly shown in Fig. 8b, is consistent to an earlier green-up in warmer years.

JSBACH does not show any discernible trend in any of the three stations (Fig. 8a, example for BRW), despite the fact that it returns a similar correlation pattern at the start of the growing season with local temperature (Fig. 8c), in particular in the extratropical northern areas. The final, global score for this trait is very low, despite the good visual matching, because of the low cell-by-cell correspondence (Fig. 2, Table 4). These two analyses underline that the model, although it realistically simulates the beginning of the growing season (Sect. 3.1), is likely to respond too weakly to land surface temperature anomalies.

#### 3.5.2 Interannual variability of vegetation activity and regional climate (Table 3)

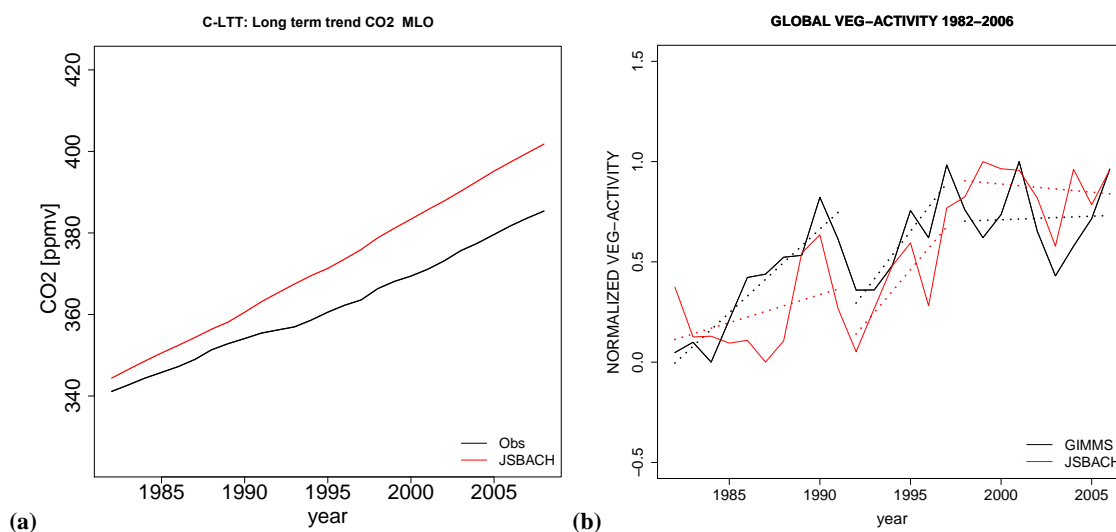
The vegetation activity is analysed separately for each climatic driver. It is not possible to clearly disentangle temperature and precipitation effects. Nonetheless, the analysis suggests that the NDVI at high latitudes is mainly correlated with surface air temperature, where plant growth is mainly limited by temperature. An exception to this pattern is Eurasia Boreal (EAB), which shows a higher co-variation of vegetation activity with precipitation pattern. NDVI in regions with dominance of shrubs/grassland is mainly driven by precipitation anomalies – in agreement with previous studies (Groeneveld and Baugh, 2007).

Figure 9a–b presents exemplary the computed cross correlograms for Eurasia Temperate (EATe) and the North American Boreal (NAB). The pattern returned in NAB, which is common to NATe and EUR, reveals a strong co-variation of vegetation activity and temperature in both data and model. However, the model behaviour suggests a strong correlation with temperature even in areas where the observations suggest a stronger covariation with precipitation (measures as SPI), as for instance in the EATe. One notable feature in these regions is that JSBACH shows a larger delay in the response of vegetation activity to SPI than observed, with differences of the order of 2–3 months (EAB, NA, SA). The final JSBACH score is good for this trait (Fig. 2, Table 4) when considering an average performance over all the regions. Low scores are obtained in precipitation-driven areas mainly due to the different time lag of the response, which corresponds well with the aforementioned too low sensitivity of raingreen vegetation to seasonal drought.

The selected trait underlines the tendency of the system to respond in a specific way to external forcing/climate, or to respond instantaneously or with some lag, and the metric is selected in order to be sensitive to model–data difference of phase rather than absolute differences of climate and vegetation activity. An important aspect emerging from this simple trait is that the detected delay could hide an incorrect representation of the effects of soil drought on vegetation growth or soil hydrology. The same regions in which the model shows a delayed response to precipitation also show a persistent negative trend in vegetation activity (Sect. 3.4, Fig. 7). This pattern is evident in particular in South East Asia, South America Temperate, and Australia, which are mainly dominated by grasslands, shrub lands, or crops. Even if other non-climatic effects at smaller spatial scales (i.e. land degradation and management practices, and fire recurrence) might affect vegetation cover and activity (Foley et al., 2005), the longer lag in the co-variation of vegetation and precipitation might be caused by the same model fault responsible for the mismatch in the vegetation trends. From a biogeophysical point of view, this model feature could also indicate a less reliable capability of the land surface model to return memory effects of the vegetation–precipitation pattern

**Table 3.** List of atmospheric CO<sub>2</sub> and vegetation activity traits used for the analyses at the interannual time scale (higher than annual frequency). A detailed explanation of metrics can be found in Sect. 2 and Appendices B and C).

Interannual time scales					
CO <sub>2</sub> Trait	Label	Test	Metric	Section	
				Methods	Results
Long-term trend	C-LTT	direct comparison of numerical trend	single value comparison metric Eq. (B4)	2.4.3	3.3
Interhemispheric gradient	IHG	variability in time	standard deviation based metric Eq. (B2)	2.4.3	3.3
CO <sub>2</sub> growth rate – regional drivers relationships	C-CL	covariance with time lag	Pearson correlation $r$	2.4.5	3.5.3
Apparent C-land sensitivity to surface temperature	C-CLsens	direct comparison of numerical trend	single-value comparison metric Eq. (B4)	2.4.5	3.5.3
Vegetation Trait	Label	Test	Metric	Methods	Results
Vegetation trend (1982–91/1992–97/1998–2006)	V-LTT	occurrence of positive/negative/no trends	spatial ranking	2.4.4	3.4
Veg. activity ~ regional drivers relationships	V-CL	covariance with time lag	Pearson correlation $r$	2.4.5	3.5.2

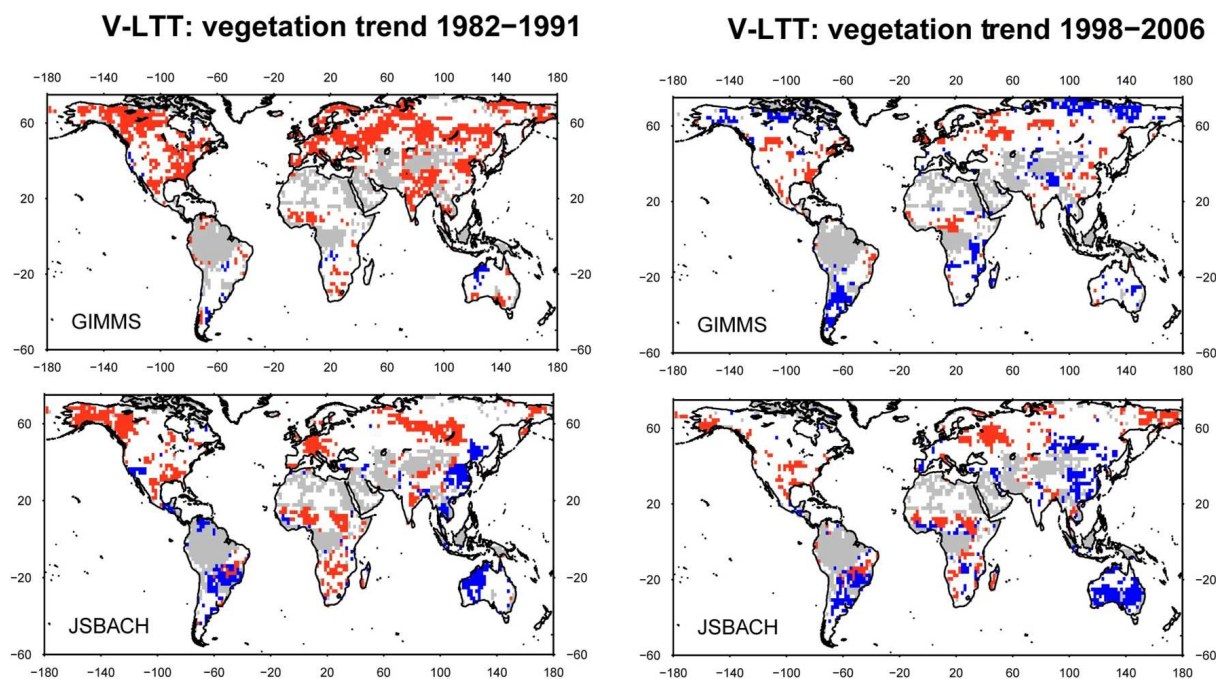
**Fig. 6.** (a) Long-term pattern of atmospheric CO<sub>2</sub> at the station of Mauna Loa (MLO); (b) normalized annual values of vegetation activity (excluding tropical, desert, and ice areas) for GIMMS-NDVI and modelled FAPAR. Period of reference 1982–2006. Dotted lines represent the linear trend computed on the normalized data (qualitative analysis).

emerging in the real Earth system (Alessandri and Navarra, 2008; Hirschi et al., 2010) in a coupled Earth system model setting.

### 3.5.3 Interannual variability of CO<sub>2</sub> growth rate and regional climate (Table 3)

The analysis of the CO<sub>2</sub> growth rate revealed distinctly different behaviour in two latitudinal bands: in tropical latitudes, the correlation structure is similar between observations and model. However, JSBACH performs less well in

particular where the CO<sub>2</sub> growth rate is mainly correlated to temperature anomalies, as for instance in North American Boreal and North American Temperate regions (Fig. 9d). It is noteworthy that this model deficiency occurs despite the good correspondence in terms of vegetation temperature (Fig. 9b). One potential reason for this phenomenon might be modelled temperature sensitivities of ecosystem respiration parameterization, particularly soil-C decomposition – inconsistent with the observations. However, it is also possible that the CO<sub>2</sub> signal at the monitoring station is influenced by net land–atmosphere C fluxes in other extratropical regions,



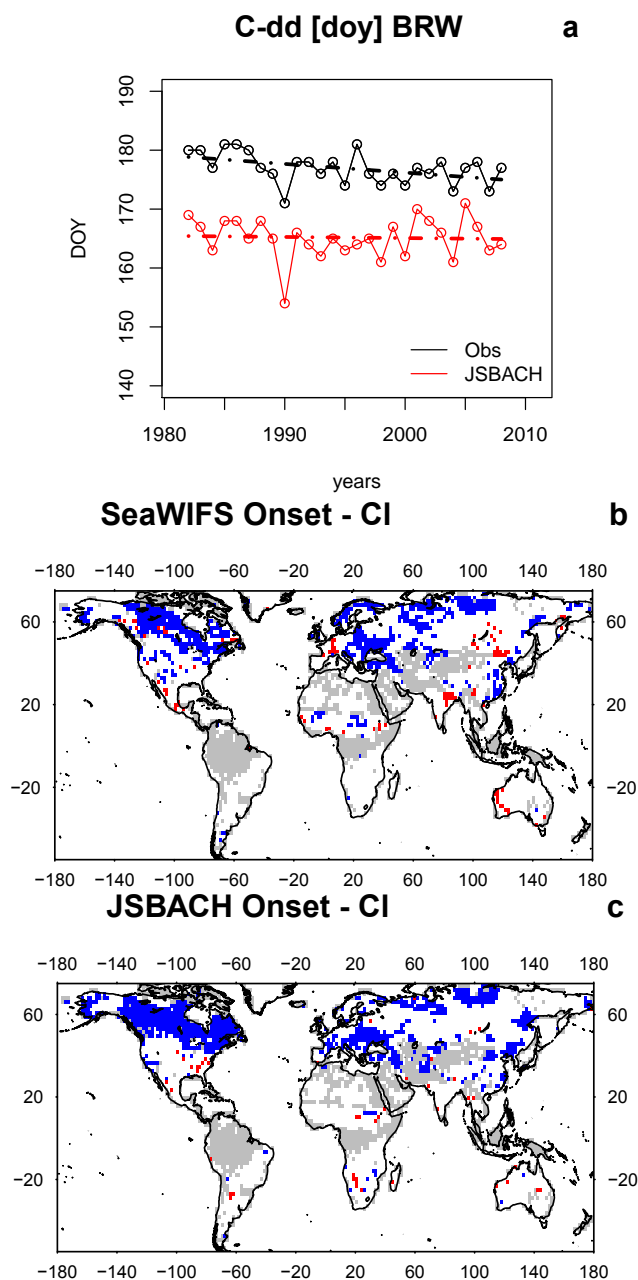
**Fig. 7.** Vegetation activity trend according to the Mann–Kendall statistics for the period of references is reported for GIMMS-NDVI and modelled FAPAR. Red: positive monotonic trend ( $P < 0.10$ ); blue: negative monotonic trend ( $P < 0.1$ ); white: no significant trend; grey: areas masked out from the analysis (grid cells with dominance of tropical forests, dominance of desert and ice).

**Table 4.** Final scores of atmospheric  $\text{CO}_2$  and vegetation activity. Atmospheric  $\text{CO}_2$  scores are reported per latitudinal band. The numerical values prior of the scaling to the baseline benchmark are reported in brackets where they differ from the final scores. For the acronyms refer to Tables 2 and 3.

Atmospheric $\text{CO}_2$ traits		Vegetation activity traits	
MSC 90N60N	0.8 (/)	$t$ -onset	0.65 (/)
MSC 60N30N	0.86 (/)	$t$ -max	0.6 (/)
MSC 30N30S	0.57 (0.64)	Onset-Cl	0.16 (0.42)
MSC 30S90S	0 (0.13)		
MSCc 90N60N	0.97 (/)		
MSCc 60N30N	0.97 (/)		
MSCc 30N30S	0.42(0.46)		
MSCc 30S90S	0 (0.19)		
MSClg	0.43 (/)		
MT 1982–1991	0.81 (0.87)		
MT 1991–1997	0.54 (0.7)		
MT 1998–2006	0.46 (0.51)		
C-dd	0.09 (0.54)		
C-LTT 90N60N	0 (0.56)	V-LTT 1982–1991	0.3 (0.5)
C-LTT 30N30S	0 (0.56)	V-LTT 1992–1997	0.1 (0.35)
C-LTT 30S90S	0 (0.53)	V-LTT 1998–2006	0.7 (0.35)
IHG 90N60N	0.7 (0.98)	V-Cl	0.6 (/)
IHG 30N30S	0.32 (0.75)		
IHG 30S90S	0.9 (0.99)		
C-CL 90N60N	0 (/)		
C-CL 30N30S	0.22 (/)		
C-CL 30S90S	0 (/)		
C-CLsens	0.02 (0.28)		

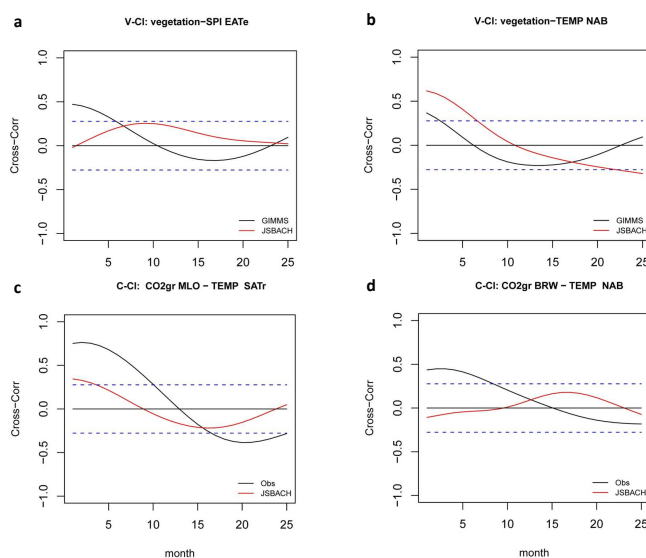
obscuring the local relationship. In general, the observed weak correspondence for the station of BRW is also observed for the station of ALT, while for the stations between  $60^\circ\text{N}$  and  $25^\circ\text{N}$ , no statistically significant co-variations were found in observations (data not shown).

In all stations, where most of the contribution to the observed concentrations is from tropical regions (e.g. South American Tropical, Northern and Southern Africa), the results reveal a good correspondence of the pattern of the covariance. However, in contrast to the observations the modelled correlation is weaker and sometimes not significant (Fig. 9c). A comparison of the time series of atmospheric  $\text{CO}_2$  and land surface climate (data not shown) reveals that the modelled time series exhibits more variability than observed and explained by, for instance, ENSO-related events. The apparent global land-C sensitivity to land surface temperature anomalies (C-Cl<sub>sens</sub>) computed for the model is not significant and very shallow (Fig. 10), in contrast to the observed sensitivity ( $4.2 \text{ Pg C yr}^{-1} \text{ K}^{-1}$ ) ( $P < 0.01$ ). It is not possible to determine to what extent the missing fire module in the current version of the model or the use of a specific transport model contribute to the observed–modelled trait mismatch involving the  $\text{CO}_2$  growth rates. However, the very low sensitivity returned by the model is comparable to the baseline benchmark (assuming a neutral biosphere; see Table 4), suggesting a deficiency in the model rather than a conceptual error in the methodology.



**Fig. 8.** (a) Atmospheric CO<sub>2</sub> drawdown points (C-dd) as computed at the station of Barrow (BRW) for observations and model. (b) and (c) partial correlation between time of onset and mean annual temperature computed for observations and JSBACH, for the period 1998–2005. Red: positive correlations ( $P < 0.1$ ); blue: negative correlations ( $P < 0.1$ ); white: no significant correlations; grey: areas masked out from the analysis (see text).

The CO<sub>2</sub> growth rate is the result of several concurrent biospheric and anthropogenic signals with dominance of land contributions, coming from several areas on the global land. Because of this large contribution from land, and the detailed regional analysis we have performed, the CO<sub>2</sub>-growth-rate-

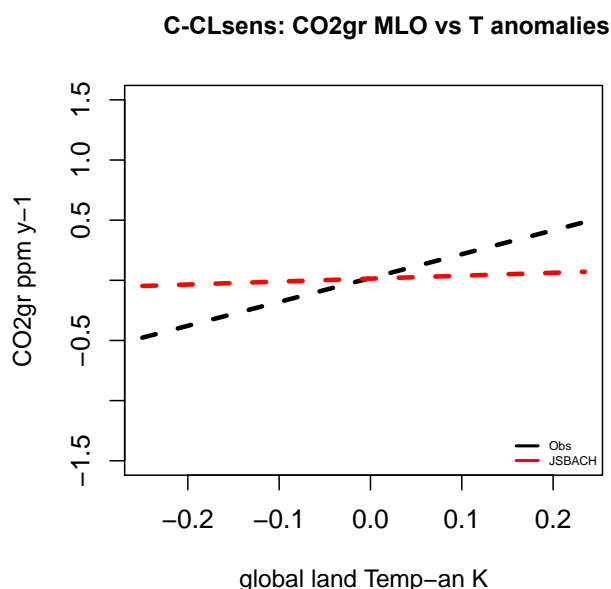


**Fig. 9.** (a) Cross correlation between precipitation pattern (SPI) and vegetation activity in Eurasia Temperate (EATE). (b) Cross correlation between temperature and vegetation activity in North American Boreal (NAB). (c) Atmospheric CO<sub>2</sub> growth rate in the station of Mauna Loa (MLO) and temperature pattern in the South American Tropical (SATr). (d) Atmospheric CO<sub>2</sub> growth rate in Barrow (BRW) and temperature pattern in NAB. Dotted lines are confidence intervals at significant level of  $P < 0.05$  (two-tailed statistics).

based traits are a useful diagnostic to indicate potential conflicts between model and observations that deserve further investigation, even though a process attribution is not possible without the use of further data streams.

As suggested by Rafelski et al. (2009), a similarity in the climate sensitivity of the underlying C processes at interannual and decadal time scales is likely to exist and to be mostly attributable to the land biosphere. This would imply that the poor results obtained from the JSBACH model in the land-C sensitivity trait could also indicate a potential for a model deficiency at longer temporal scales with respect to the net land-C exchange.

Common to most of the evaluation schemes, data and model errors are not considered explicitly in the mathematical formulation of the metrics. This constitutes a major limitation of this and other evaluation frameworks. As stated in the introduction, uncertainties in observations or reference datasets are not always provided or quantified, and this poses challenges for the computation of model–data/model–model differences and the significance of these distances. Structural model errors can only be assessed with a dedicated study investigating the effect of alternative model structures on surface fluxes, which is beyond the scope of a benchmarking scheme. Conversely, the scheme proposed here can help to quantify the model structural error if different model variants are available. Where possible, we have minimized the conceptual difference between model and observation by only



**Fig. 10.** Apparent land-C sensitivity: CO<sub>2</sub> growth rate in Mauna Loa (MLO) versus global land surface temperature. Regression is significant at  $P < 0.01$  for observations. Annual data points are omitted for clarity.

considering those features (traits) that can be robustly compared and have isolated the contribution of the land versus oceanic and anthropogenic influences. The proposed evaluation framework defines bounded metrics that allows stating whether and how much the model adds information to the simulation of carbon cycle trends and thus whether the model lies in the range of acceptable performances. This provides an additional constraint to the model performance.

#### 4 Concluding remarks

Pertinent information on current C-cycle-related processes contained in the atmospheric CO<sub>2</sub> record and the satellite-based records of vegetation activity were compiled and synthesized into easily identifiable traits and a framework of intuitively comparable metrics. The results of the exploratory analysis of C-related processes and climate variability were presented with emphasis on the robustness of the information content of the observations, making use of both atmospheric CO<sub>2</sub> concentration and vegetation activity at the appropriate time- and spatial scale of a global land surface model.

The results show that the simultaneous use of the atmospheric CO<sub>2</sub> record and satellite-based vegetation activity as two independent datasets help to identify the sources of data-model mismatch in terms of regional source of errors, or to detect potential compensation errors. In particular, the separate analysis of the atmospheric CO<sub>2</sub> and vegetation activity circumvent the problem that the atmospheric CO<sub>2</sub> retains the

net effect of both vegetation activity (i.e. photosynthesis) and ecosystem C release response.

The use of a baseline benchmark with a clear ecological meaning was shown to be a valuable approach to provide a more robust and objective quantification of data-model disagreement. In addition, scaling the metric against a reference case allows more independence by the section of a specific metric and avoidance of misleading interpretation of the numerical score.

A key component of the evaluation framework developed here is that it is designed to be suitable and sensitive to evaluate global land surface models both in offline mode – i.e. when driven by observed climate variability – and fully coupled to Earth system models with a different climate and climate variability. Therefore, in addition to providing metrics for key traits that describe climatological mean variables, we use a range of correlational metrics to analyse the climate sensitivity of key carbon cycle traits. We demonstrate that these metrics provide insight into the realism of the carbon cycle simulation that go beyond an evaluation of mean states and trends. In this paper, we described the framework and applied it to an example model. The next step will be the use of this framework to evaluate online and offline versions of JSBACH. Nonetheless, even application of the benchmarking framework for the evaluation of the JSBACH model in offline mode already allows certain conclusions particular to the model:

- The traits at seasonal time scales showed that high-latitude terrestrial ecosystem patterns are a major strength of JSBACH, with good performance both in terms of mean vegetation activity and mean seasonal CO<sub>2</sub> cycle in the high-latitude stations. Lower performance of mean pattern of phenology occurs in the Southern Hemisphere, in particular in shrub-dominated areas and in deciduous broadleaved forests in Southern Africa. A systematic overestimation of the seasonal cycle of CO<sub>2</sub> points to a too high magnitude of the seasonal land gross C fluxes.
- The observed weakening of the positive warming effect in vegetation in the NH and the trend toward a neutral/negative effect in the SH pronounced in last decade are not fully captured by the model, both in CO<sub>2</sub> and vegetation activity traits. The analysis of vegetation–climate covariance revealed that the modelled ecosystem response is primarily driven by temperature anomalies, suggesting that this discrepancy might be associated with an incorrect sensitivity of vegetation to precipitation anomalies at interannual time scales.
- While the analysis of CO<sub>2</sub> growth rate and climate drivers returned a weak covariation of the atmospheric signal with climate on selected regions on land, the model deviates strongly from the observations both in terms of the long-term trend of the atmospheric CO<sub>2</sub>,



**Table A1.** Final scores of atmospheric CO<sub>2</sub> and vegetation activity for the baseline benchmark. Atmospheric CO<sub>2</sub> scores are reported per latitudinal band.

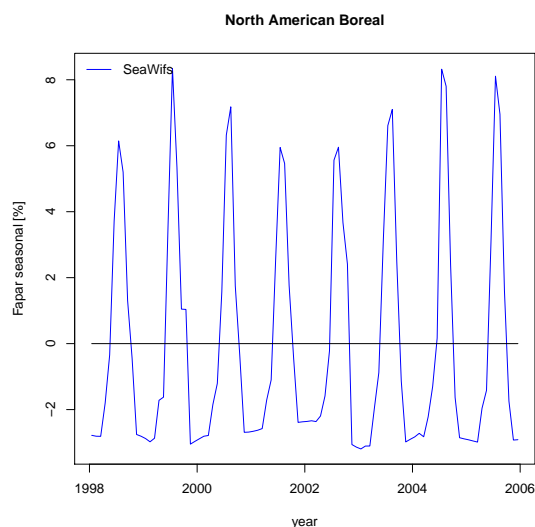
Atmospheric CO <sub>2</sub> traits		Vegetation activity traits	
MSC 90N60N	0	<i>t</i> -onset	0
MSC 60N30N	0	<i>t</i> -max	0
MSC 30N30S	0.23	Onset-Cl	0.3
MSC 30S90S	0.45		
MSCc 90N60N	0.1		
MSCc 60N30N	0.1		
MSCc 30N30S	0.18		
MSCc 30S90S	0.91		
MSClg	0		
MT 1982–1991	0.2		
MT 1992–1997	0.4		
MT 1998–2006	0.25		
C-dd	0.5		
<hr/>			
C-LTT 90N60N	0.56	V-LTT 1982–1991	0.3
C-LTT 30N30S	0.58	V-LTT 1992–1997	0.3
C-LTT 30S90S	0.55	V-LTT 1998–2006	0.3
IHG 90N60N	0.54	V-Cl	0
IHG 30N30S	0.39		
IHG 30S90S	0.18		
C-CL 90N60N	0		
C-CL 30N30S	0		
C-CL 30S90S	0		
C-CLsens	0.26		

and therefore the implied net land-C uptake, and the apparent interannual land-C sensitivity to temperature anomalies. The combined analysis of CO<sub>2</sub> with the vegetation trend analysis suggests that a too high soil-C turnover rate might be responsible for the underestimation of net land-C uptake.

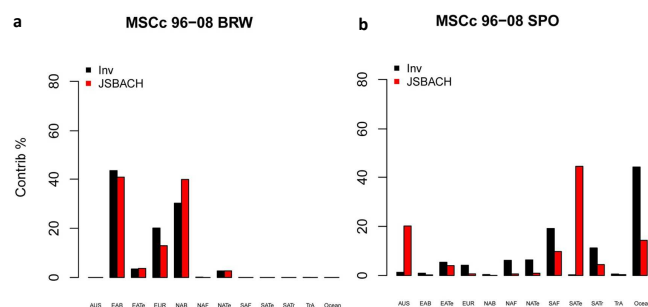
### Appendix A

#### Computation of SPI index

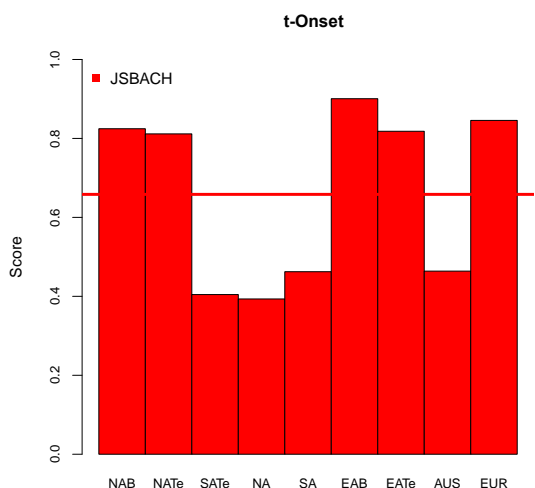
The SPI is the transformation of the precipitation time series into a standardized normal distribution (z distribution). First, a gamma distribution is fitted to the cumulative precipitation frequency distribution. The gamma distribution has been used to fit the empirical frequency of data. Since the gamma distribution is undefined for null values of the variables, the cumulative probability has been corrected according to Lloyd-Hughes and Sanders (2002). Using an equiprobable transformation, the cumulative probability function of the gamma distribution is then transformed into the normal distribution function.



**Fig. A1.** Exemplar of determination of time of onset of the seasonal signal, zero centred, of the vegetation activity (SeaWiFS-FAPAR). Time series extracted from one grid cell located in North American Boreal.



**Fig. A2.** Regional contribution to the mean seasonal cycle in the stations of Barrow (BRW) and South Pole (SPO), expressed as percentage contribution. For the region labels refer to Fig. 1.



**Fig. A3.** Final score of the time of onset trait computed according to Transcom3 land regions. (Tropical Asia and South American Tropical). The horizontal line represents the mean of the regional scores.

## Appendix B

In addition to the classical statistics as the Pearson correlation coefficient ( $r$ ), the squared correlation coefficient ( $r^2$ ), cross correlation, and standard deviation statistics ( $\sigma$ ), metrics were selected as combination of some of the previous statistics and built ad hoc for the specific trait analysed.

As reported in Taylor (2000):

$$\frac{4(1+r)^4}{\hat{\sigma}_f + \frac{1}{\hat{\sigma}_f} (1+R_0)^4} \quad (\text{B1})$$

In this metric, more weight is given to the capability of the model to return the right phase of the trait rather than the amplitude.  $\hat{\sigma}_f = \sigma_m/\sigma_0$  is the ratio between the modelled standard deviation and the observed standard deviation of the trait of interest.  $R_0$  is the maximum correlation achievable and assumed to be 1.

Comparison of variability of the signal via standard deviation:

$$\frac{4}{\left(\hat{\sigma}_f + \frac{1}{\hat{\sigma}_f}\right)^2} \quad (\text{B2})$$

Linear differences metric:

$$\frac{|6 - |O - M||}{6}, \quad (\text{B3})$$

where  $O$  is the observed value and  $M$  is the modelled value. It is applied to the most frequent month of the variable observed (0 when the maximum difference of the variables is six months, 1 when no differences occur).

Single value comparison metric:

$$\frac{1}{(1 + |(O - M)/O|)^2}, \quad (\text{B4})$$

where  $O$  is the observed value and  $M$  is the modelled value.

At exception of the Taylor statistics, all the other metrics are symmetric.

Map cell-by-cell comparison metric:

The ranking metric specifies the number of agreement cells against the total observed cell belonging to a specific class. The final score is the average over the selected classes. Three classes were used in our framework: no statistically significant relationship (i.e. no correlation, no trend detected), positive relationship (i.e. correlation/trend), and negative relationships (i.e. correlation/trend) detected.

In terms of lower benchmark, the case of constant vegetation has been used. This is the equivalent to analysing the returned trend against a null hypothesis of non-changing vegetation. The average score obtained under this setting is equal to 0.3, considering the agreement cell-by-cell to each single class. The score of the model is thereafter scaled to this lower benchmark.

## Appendix C

### Additional constraints for the computation of the final score

Negative correlations between model and dataset the final score to 0, with the exception of the cross-correlation traits. If the modeled signal has no standard deviation (i.e. constant vegetation activity), the score is automatically set to 0, if the observed signal has no standard deviation the score is set to NA. Only grid cells with a valid observed signal were considered in the model–data comparison analysis.

Correlations and cross correlations, trend, and number of growing seasons were tested against random noise ( $t$  two-tailed statistics) at least  $P < 0.1$  significance. For the scores based on cross-correlation statistics with climate drivers, the score is set to NA when observations do not show any statistically significant relationship. If the model does not return any significant relationship, the score is set to 0.

When testing the degree and persistence of the association between temporal series using two-tailed  $t$  test, potential autocorrelations in the temporal series were considered by adjusting the degrees of freedom, hence the number of independent information (Trenberth and Caron, 2000). We assume a number  $N/2$  of independent information, where  $N$  is the total number of months in the record (300 months).

For the atmospheric CO<sub>2</sub> traits, the final score is the average of the scores obtained each individual monitoring station. In terms of comparison to remote sensing data, the scores were first aggregated by vegetation class for each Transcom3 region, and then further aggregated using a weighted average and taking account of the number of grid cells belonging to the specific vegetation class.

## Appendix D

### The Synmap classification vegetation map

The following vegetation classes of the Synmap dataset (Jung et al., 2006) were considered: shrubs, grass, crop, deciduous broad leaved forest (dbf), deciduous needle-leaved forest (dnf), evergreen broadleaved forest (ebf), evergreen needle-leaved forest (enf) along with unvegetated area (i.e. bared soil, ice lands), and water. This way to aggregate the information instead of using the model's vegetation classification helps to maintain flexibility and comparability across different model platforms and thereby creates less uncertainties in the performance evaluation analysis. The most dominant class is computed as the one covering at least the 80% of the total area of each grid cell.

*Acknowledgements.* The research leading to these results has received funding from the Seventh Framework Programme (FP7 2007-2013) under grant agreement no. [238366]. The authors are grateful to Christian Reick, Veronika Gayler, and Reiner Schnur for help with the JSBACH model. The authors furthermore wish to thank Christian Rödenbeck and two anonymous reviewers for helpful comments on the manuscript and constructive discussions.

The service charges for this open access publication have been covered by the Max Planck Society.

Edited by: P. Stoy

## References

- Alessandri, A. and Navarra, A.: On the coupling between vegetation and rainfall inter-annual anomalies: Possible contributions to seasonal rainfall predictability over land areas, *Geophys. Res. Lett.*, 35, L02718, doi:10.1029/2007GL032415, 2008.
- 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*, in press, doi:10.1175/JCLI-D-12-00417.1, 2013.
- Angert, A., Biraud, S., Bonfils, C., Henning, C. C., Buermann, W., Pinzon, J., Tucker, C. J., and Fung, I.: Drier summers cancel out the CO<sub>2</sub> uptake enhancement induced by warmer springs, *P. Natl. Acad. Sci. USA*, 102, 10823–10827, doi:10.1073/pnas.0501647102, 2005.
- Arora, V., Boer, G., Friedlingstein, P., Eby, M., Jones, C., Christian, J., Bonan, G., Bopp, L., Brovkin, V., Cadule, P., Hajima, T., Ilyina, T., Lindsay, K., Tjiputra, J., and Wu, T.: Carbon-concentration and carbon-climate feedbacks in CMIP5 Earth system models, *J. Climate*, in press, doi:10.1175/JCLI-D-12-00494.1, 2013.
- Asner, G. P. and Alencar, A.: Drought impacts on the Amazon forest: the remote sensing perspective, *New Phytol.*, 187, 569–578, doi:10.1111/j.1469-8137.2010.03310.x, 2010.
- Badeck, F.-W., Bondeau, A., Bottcher, K., Doktor, D., Lucht, W., Schaber, J., and Sitch, S.: Responses of spring phenology to climate change, *New Phytol.*, 162, 295–309, doi:10.1111/j.1469-8137.2004.01059.x, 2004.
- 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.
- Beer, C., Reichstein, M., Tomelleri, E., Ciais, P., Jung, M., Carvalhais, N., Rödenbeck, C., Arain, M. A., Baldocchi, D., Bonan, G. B., Bondeau, A., Cescatti, A., Lasslop, G., Lindroth, A., Lomas, M., Luysaert, S., Margolis, H., Oleson, K. W., Rouspard, O., Veenendaal, E., Viovy, N., Williams, C. F. I., and Papale, D.: Terrestrial Gross Carbon Dioxide Uptake: Global Distribution and Covariation with Climate, *Science*, 329, 834–838, doi:10.1126/science.1184984, 2010.
- Blasing, T. J., Broniak, C. T., and Marland, G.: The annual cycle of fossil-fuel carbon dioxide emissions in the United States, *Tellus B*, 57, 107–115, 2005.
- 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.: Forests and climate change: forcings, feedbacks, and the climate benefits of forests, *Science*, 320, 1444–14449, doi:10.1126/science.1155121, 2008.
- Brown, M. E., Pinzón, J. E., Didan, K., Morisette, J. T., and Tucker, C. J.: Evaluation of the Consistency of Long-Term NDVI Time Series Derived From AVHRR, and Landsat ETM + Sensors, Sensors (Peterborough, NH), 44, 1787–1793, 2006.
- Buermann, W., Anderson, B., Tucker, C. J., Dickinson, R. E., Lucht, W., Potter, C., and Myneni, R. B.: Interannual covariability in Northern Hemisphere air temperatures and greenness associated with El Niño–Southern Oscillation and the Arctic Oscillation, *J. Geophys. Res.*, 108, 4396, doi:10.1029/2002JD002630, 2003.
- Buermann, W., Lintner, B. R., Koven, C. D., Angert, A., Pinzon, J. E., Tucker, C. J., and Fung, I. Y.: The changing carbon cycle at Mauna Loa Observatory, *P. Natl. Acad. Sci. USA*, 104, 4249–4254, doi:10.1073/pnas.0611224104, 2007.
- 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 CO<sub>2</sub> measurements, *Global Biogeochem. Cy.*, 24, GB2016, doi:10.1029/2009GB003556, 2010.
- Conway, T. J., Tans, P. P., Waterman, L. S., Thoning, K. W., Kitzis, D. R., Masarie, K. A., and Zhang, N.: Evidence for interannual variability of the carbon cycle from the National Oceanic and Atmospheric Administration/Climate Monitoring and Diagnostics Laboratory Global Air Sampling Network, *J. Geophys. Res.*, 99, 22831–22855, doi:10.1029/94JD01951, 1994.
- Cox, P. M., Betts, R. A., Jones, C. D., Spall, S. A., and Totterdell, I. J.: Acceleration of global warming due to carbon-cycle feedbacks in a coupled climate model, *Nature*, 408, 184–187, doi:10.1038/35041539, 2000.
- Dahlke, C., Loew, A., and Reick, C. H.: 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, 2013.
- Deser, C., Phillips, A., Bourdette, V., and Teng, H.: Uncertainty in climate change projections: the role of internal variability, *Climate Dynam.*, 38, 527–546, doi:10.1007/s00382-010-0977-x, 2010.
- Estrella, N. and Menzel, A.: Responses of leaf colouring in four deciduous tree species to climate and weather in Germany, *Climate Res.*, 32, 253–267, 2006.
- Foley, J. A., Defries, R., Asner, G. P., Barford, C., Bonan, G., Carpenter, S. R., Chapin, F. S., Coe, M. T., Daily, G. C., Gibbs, H. K., Helkowski, J. H., Holloway, T., Howard, E. A., Kucharik, C. J., Monfreda, C., Patz, J. A., Prentice, I. C., Ramankutty, N., and Snyder, P. K.: Global consequences of land use, *Science*, 309, 570–574, doi:10.1126/science.1111772, 2005.
- Friedlingstein, P., Cox, P., Betts, R., Bopp, L., von Bloh, W., Brovkin, V., Cadule, P., Doney, S., Eby, M., Fung, I., Bala, G., John, J., Jones, C., Joos, F., Kato, T., Kawamiya, M., Knorr, W., Lindsay, K., Matthews, H. D., Raddatz, T., Rayner, P., Reick, C., Roeckner, E., Schnitzler, K.-G., Schnur, R., Strassmann, K., Weaver, A. J., Yoshikawa, C., and Zeng, N.: Climate – Carbon

- Cycle Feedback Analysis: Results from the C 4 MIP Model Intercomparison, *J. Climate*, 19, 3337–3353, 2006.
- Gobron, N., Pinty, B., Ausedat, O., Chen, J. M., Cohen, W. B., Fensholt, R., Gond, V., Huemmrich, K. F., Lavergne, T., Mélin, F., Privette, J. L., Sandholt, I., Taberner, M., Turner, D. P., Verstraete, M. M., and Widlowski, J.: Evaluation of fraction of absorbed photosynthetically active radiation products for different canopy radiation transfer regimes: Methodology and results using Joint Research Center products derived from SeaWiFS against ground-based estimations, *J. Geophys. Res.*, 111, D13110, doi:10.1029/2005JD006511, 2006a.
- Gobron, N., Pinty, B., Taberner, M., Mélin, F., Verstraete, M. M., and Widlowski, J.-L.: Monitoring the photosynthetic activity of vegetation from remote sensing data, *Adv. Space Res.*, 38, 2196–2202, doi:10.1016/j.asr.2003.07.079, 2006b.
- Goetz, S. J., Bunn, A. G., Fiske, G. J., and Houghton, R. A.: Satellite-observed photosynthetic trends across boreal North America associated with climate and fire disturbance, *P. Natl. Acad. Sci. USA*, 102, 13521–13525, 2005.
- Groeneveld, D. and Baugh, W.: Correcting satellite data to detect vegetation signal for eco-hydrologic analyses, *J. Hydrol.*, 344, 135–145, doi:10.1016/j.jhydrol.2007.07.001, 2007.
- Gurney, K. R., Law, R. M., Denning, A. S., Rayner, P. J., Baker, D., Bousquet, P., Bruhwiler, L., Chen, Y.-H., Ciais, P., Fan, S., Fung, I. Y., Gloor, M., Heimann M., Higuchi, K., John, J., Maki, T., Maksyutov, S., Masariek, K., Peylin, P., Pratherkk, M., Pakkk, B. C., Randerson, J., Sarmiento, J., Taguchi, S., Takahashi, T., and Yuen, C.: Towards robust regional estimates of CO<sub>2</sub> sources and sinks using atmospheric transport models, *Nature*, 415, 626–630, doi:10.1038/415626a, 2002.
- Gurney, K. R., Law, R. M., Denning, A. S., Rayner, P. J., Baker, D., Bousquet, P., Bruhwiler, L., Chen, Y.-H., Ciais, P., Fan, S., Fung, I. Y., Gloor, M., Heimann, M., Higuchi, K., John, J., Kowalczyk, E., Maki, T., Maksyutov, S., Peylin, P., Prather, M., Pak, B. C., Sarmiento, J., Taguchi, S., Takahashi, T., and Yuen, C.: TransCom3 CO<sub>2</sub> inversion intercomparison: 1. Annual mean control results and sensitivity to transport and prior flux information, *Tellus B*, 55, 555–579, doi:10.1034/j.1600-0889.2003.00049.x, 2003.
- Heimann, M., Esser, G., Haxeltine, A., Kaduk, J., Kicklighter, D. W., Knorr, W., Kohlmaier, G. H., Mcguire, A. D., Melillo, J., Moore III, B., Otto, R. D., Prentice, I. C., Sauf, W., Schloss, A., Sitch, S., Wittenberg, U., and Wurth, G.: Evaluation of terrestrial carbon cycle models through simulations of the seasonal cycle of atmospheric First results of a model intercomparison study, *Global Biogeochem. Cy.*, 12, 1–24, doi:10.1029/97GB01936, 1998.
- Hirschi, M., Seneviratne, S. I., Alexandrov, V., Boberg, F., Boroneant, C., Christensen, O. B., Formayer, H., Orłowski, B., and Stepanek, P.: Observational evidence for soil-moisture impact on hot extremes in southeastern Europe, *Nat. Geosci.*, 4, 17–21, doi:10.1038/ngeo1032, 2010.
- Holben, B. N.: Characteristics of maximum-value composite images from temporal AVHRR data, *Int. J. Remote Sens.*, 7, 1417–1434, 1986.
- Huete, A., Didan, K., Miura, T., Rodriguez, E., Gao, X., and Ferreira, L.: Overview of the radiometric and biophysical performance of the MODIS vegetation indices, *Remote Sens. Environ.*, 83, 195–213, doi:10.1016/S0034-4257(02)00096-2, 2002.
- Hurt, G. C., Frothing, S., Fearon, M. G., Moore, B., Shevliakova, E., Malyshev, S., Pacala, S. W., and Houghton, R. A.: The underpinnings of land-use history: Three centuries of global gridded land-use transitions, wood-harvest activity, and resulting secondary lands, *Glob. Change Biol.*, 1208–1229, doi:10.1111/j.1365-2486.2006.01150.x, 2006.
- Jacobson, A. R., Mikaloff Fletcher, S. E., Gruber, N., Sarmiento, J. L., and Gloor, M.: A joint atmosphere-ocean inversion for surface fluxes of carbon dioxide: 1. Methods and global-scale fluxes, *Global Biogeochem. Cy.*, 21, doi:10.1029/2005GB002556, 2007.
- Jung, M., Henkel, K., Herold, M., and Churkina, G.: Exploiting synergies of global land cover products for carbon cycle modeling, *Remote Sens. Environ.*, 101, 534–553, doi:10.1016/j.rse.2006.01.020, 2006.
- Jung, M., Reichstein, M., Ciais, P., Seneviratne, S. I., Sheffield, J., Goulden, M. L., Bonan, G., Cescatti, A., Chen, J., de Jeu, R., Dolman, A. J., Eugster, W., Gerten, D., Gianelle, D., Gobron, N., Heinke, J., Kimball, J., Law, B. E., Montagnani, L., Mu, Q., Mueller, B., Oleson, K., Papale, D., Richardson, A. D., Rouspard, O., Running, S., Tomelleri, E., Viovy, N., Weber, U., Williams, C., Wood, E., Zaehle, S., and Zhang, K.: Recent decline in the global land evapotranspiration trend due to limited moisture supply, *Nature*, 467, 951–954, doi:10.1038/nature09396, 2010.
- Kalnay, E., Kanamitsu, M., Kistler, R., Collins, W., Deaven, D., Gandin, L., Iredell, M., Saha, S., White, G., Wollen, J., Zhu, Y., Chelliah, M., Ebisuzaki, W., Higgins, W., Janowiak, J., Mo, K. C., Ropelewski, C., Wang, J., Leetmaa, A., Reynolds, R., Jenne, R., and Joseph, D.: The NCEP/NCAR 40-year reanalysis project, *Bull. Am. Meteor. Soc.*, 77, 437–471, 1996.
- Keeling, C. D., Whorf, T. P., Wahlen, M., and van der Plicht, J.: Interannual extremes in the rate of rise of atmospheric carbon dioxide since 1980, *Nature*, 375, 666–670, 1995.
- Keeling, C. D., Chin, J. F. S., and Whorf, T. P.: Increased activity of northern vegetation inferred from atmospheric CO<sub>2</sub> measurements, *Nature*, 382, 146–149, 1996.
- Lloyd-Hughes, B. and Saunders, M. A.: A drought climatology for Europe, *Int. J. Climatol.*, 22, 1571–1592, doi:10.1002/joc.846, 2002.
- Lucht, W., Prentice, I. C., Myneni, R. B., Sitch, S., Friedlingstein, P., Cramer, W., Bousquet, P., Buermann, W., and Smith, B.: Climatic control of the high-latitude vegetation greening trend and Pinatubo effect, *Science*, 296, 1687–1689, doi:10.1126/science.1071828, 2002.
- Maignan, F., Bréon, F. M., Vermote, E., Ciais, P., and Viovy, N.: Mild winter and spring 2007 over western Europe led to a widespread early vegetation onset, *Geophys. Res. Lett.*, 35, L02404, doi:10.1029/2007GL032472, 2008.
- McKee, T. B., Doesken, N. J., and Kleist, J.: The relationship of drought frequency and duration to time scales, *Conference Proceedings, Eighth Conference on Applied Climatology*, 17–22 January, Anaheim, California, 1993.
- Menzel, A., Sparks, T. H., Estrella, N., Koch, E., Aasa, A., Ahas, R., Alm-Kübler, K., Bissolli, P., Braslavská, O., Briede, A., Chmielewski, F. M., Crepinsek, Z., Curnel, Y., Dahl, A., Defila, C., Donnelly, A., Filella, Y., Jatca, K., Mâge, F., Mestre, A., Nordli, Ø., Peñuelas, J., Pirinen, P., Remišova, V., Scheffinger, H., Striz, M., Susni, A., Van Vliet, A. J. H., Wielgolaski, F., Zach, S., and Züst, A.: European phenological response to cli-

- mate change matches the warming pattern, *Glob. Change Biol.*, 12, 1969–1976, doi:10.1111/j.1365-2486.2006.01193.x, 2006.
- Mikaloff Fletcher, S. E., Gruber, N., Jacobson, A. R., Doney, S. C., Dutkiewicz, S., Gerber, M., Follows, M., Joos, F., Lindsay, K., Menemenlis, D., Mouchet, A., Müller, S. A., and Sarmiento, J. L.: Inverse estimates of anthropogenic CO<sub>2</sub> uptake, transport, and storage by the ocean, *Global Biogeochem. Cy.*, 20, GB2002, doi:10.1029/2005GB002530, 2006.
- Mikaloff Fletcher, S. E., Gruber, N., Jacobson, A. R., Gloor, M., Doney, S. C., Dutkiewicz, S., Gerber, M., Follows, M., Joos, F., Lindsay, K., Menemenlis, D., Mouchet, A., Müller, S. A., and Sarmiento, J. L.: Inverse estimates of the oceanic sources and sinks of natural CO<sub>2</sub> and the implied oceanic carbon transport, *Global Biogeochem. Cy.*, 21, GB1010, doi:10.1029/2006GB002751, 2007.
- Myneni, R. B. and Williams, D. L.: On the Relationship between FAPAR and NDVI, *Remote Sens. Environ.*, 49, 200–211, 1994.
- Myneni, R. B., Keeling, C. D., Tucker, C. J., Asrar, G., and Nemani, R. R.: Increased plant growth in the northern high latitudes from 1981 to 1991, *Nature*, 386, 698–702, 1997.
- Notaro, M., Vavrus, S., and Liu, Z.: Global Vegetation and Climate Change due to Future Increases in CO<sub>2</sub> as Projected by a Fully Coupled Model with Dynamic Vegetation\*, *J. Climate*, 20, 70–90, doi:10.1175/JCLI3989.1, 2007.
- Peñuelas, J., Rutishauser, T., and Filella, I.: Ecology. Phenology feedbacks on climate change, *Science*, 324, 887–888, doi:10.1126/science.1173004, 2009.
- Peylin, P., Bousquet, P., Le Quéré, C., Sitch, S., Friedlingstein, P., McKinley, G., Gruber, N., Rayner, P. J. and Ciais, P.: Multiple constraints on regional CO<sub>2</sub> flux variations over land and oceans, *Global Biogeochem. Cy.*, 19, GB1011, doi:10.1029/2003GB002214, 2005.
- Prentice, I. C., Heimann, M., and Sitch, S.: The carbon balance of the terrestrial biosphere: ecosystem models and atmospheric observations, *Ecol. Appl.*, 10, 1553–1573, 2000.
- Le Quere, C., Aumont, O., Bopp, L., Bousquet, P., Ciais, P., Francey, R., Heimann, M., Keeling, C. D., Keeling, R. F., Khesghi, H., Peylin, P., Piper, S. C., and Prentice, I. C.: Two decades of ocean CO<sub>2</sub> sink and variability, *Tellus B*, 55, 649–656, doi:10.1034/j.1600-0889.2003.00043.x, 2003.
- Le Quéré, C., Raupach, M. R., Canadell, J. G., Marland, G., Bopp, L., Ciais, P., Conway, T. J., Doney, S. C., Feely, R. A., Foster, P., Friedlingstein, P., Gurney, K., Houghton, R. A., House, J. I., Huntingford, C., Levy, P. E., Lomas, M. R., Majkut, J., Metzl, N., Ometto, J. P., Peters, G. P., Prentice, I. C., Randerson, J. T., Running, S. W., Sarmiento, J. L., Schuster, U., Sitch, S., Takahashi, T., Viovy, N., van der Werf, G. R., and Woodward, F. I.: Trends in the sources and sinks of carbon dioxide, *Nat. Geosci.*, 2, 831–836, doi:10.1038/ngeo689, 2009.
- Raddatz, T. J., Reick, C. H., Knorr, W., Kattge, J., Roeckner, E., Schnur, R., Schnitzler, K.-G., Wetzell, P., and Jungclaus, J.: Will the tropical land biosphere dominate the climate–carbon cycle feedback during the twenty-first century?, *Climate Dynam.*, 29, 565–574, doi:10.1007/s00382-007-0247-8, 2007.
- Rafelski, L. E., Piper, S. C., and Keeling, R. F.: Climate effects on atmospheric carbon dioxide over the last century, *Tellus B*, 61, 718–731, doi:10.1111/j.1600-0889.2009.00439.x, 2009.
- 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., Curtis, 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.
- Raupach, M. R., Canadell, J. G., and Qu, C. L.: Anthropogenic and biophysical contributions to increasing atmospheric CO<sub>2</sub> growth rate and airborne fraction, *Analysis*, 1991 (June 1991), 1601–1613, 2008.
- Raupach, M. R., Rayner, P. J., Barrett, D. J., Defries, R. S., Heimann, M., Ojima, D. S., Quegan, S., and Schimmler, C. C.: Model–data synthesis in terrestrial carbon observation: methods, data requirements and data uncertainty specifications, *Glob. Change Biol.*, 11, 378–397, doi:10.1111/j.1365-2486.2005.00917.x, 2005.
- Reick, C., Raddatz, T., Brovkin, V., and Gayler, V.: The representation of natural and anthropogenic land cover change in MPI-ESM, *J. Adv. Model. Earth Syst.*, 4, accepted, doi:10.1002/jame.20022, 2013.
- Richardson, A. D., Braswell, B. H., Hollinger, D. Y., Jenkins, J. P., and Ollinger, S. V.: Near-surface remote sensing of spatial and temporal variation in canopy phenology, *Ecol. Appl.*, 19, 1417–1428, 2009.
- Rödenbeck, C., Houweling, S., Gloor, M., and Heimann, M.: CO<sub>2</sub> flux history 1982–2001 inferred from atmospheric data using a global inversion of atmospheric transport, *Atmos. Chem. Phys.*, 3, 1919–1964, doi:10.5194/acp-3-1919-2003, 2003.
- 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.
- Taylor, K. E.: Summarizing multiple aspects of model performance in a single diagram, *J. Geophys. Res.*, 106, 7183–7192, 2001.
- Thoning, K. W., Tans, P. P., and Komhyr, W. D.: Atmospheric Carbon Dioxide at Mauna Loa Observatory 2. Analysis of the NOAA GMCC Data, 1974–1985, *J. Geophys. Res.*, 94, 8549–8565, 1989.
- Tucker, C., Pinzon, J., Brown, M., Slayback, D., Pak, E., Mahoney, R., Vermote, E., and El Saleous, N.: An extended AVHRR 8-km NDVI dataset compatible with MODIS and SPOT vegetation NDVI data, *Int. J. Remote Sens.*, 26, 4485–4498, doi:10.1080/01431160500168686, 2005.
- van der Werf, G. R., Randerson, J. T., Collatz, G. J., Giglio, L., Kasibhatla, P. S., Arellano, A. F., Olsen, S. C., and Kasischke, E. S.: Continental-scale partitioning of fire emissions during the 1997 to 2001 El Niño/La Niña period, *Science*, 303, 73–76, doi:10.1126/science.1090753, 2004.
- Verstraete, M. M., Gobron, N., Auzanedat, O., Robustelli, M., Pinty, B., Widlowski, J.-L., and Taberner, M.: An automatic procedure to identify key vegetation phenology events using the JRC-FAPAR products, *Adv. Space Res.*, 41, 1773–1783, doi:10.1016/j.asr.2007.05.066, 2008.
- Yang, Z., Washenfelder, R. A., Keppel-Aleks, G., Krakauer, N. Y., Randerson, J. T., Tans, P. P., Sweeney, C., and Wennberg, P. O.: New constraints on Northern Hemisphere growing season net flux, *Geophys. Res. Lett.*, 34, L12807, doi:10.1029/2007GL029742, 2007.

Zeng, N., Mariotti, A., and Wetzel, P.: Terrestrial mechanisms of interannual CO<sub>2</sub> variability, *Global Biogeochem. Cy.*, 19, GB1016, doi:10.1029/2004GB002273, 2005.

Zhou, L.: Relation between interannual variations in satellite measures of northern forest greenness and climate between 1982 and 1999, *J. Geophys. Res.*, 108, 4004, doi:10.1029/2002JD002510, 2003.