



LPJmL4 – a dynamic global vegetation model with managed land – Part 2: Model evaluation

Sibyll Schaphoff¹, Matthias Forkel², Christoph Müller¹, Jürgen Knauer³, Werner von Bloh¹, Dieter Gerten^{1,4},
Jonas Jägermeyr¹, Wolfgang Lucht^{1,4}, Anja Rammig⁵, Kirsten Thonicke¹, and Katharina Waha^{1,6}

¹Potsdam Institute for Climate Impact Research, Telegraphenberg, P.O. Box 60 12 03, 14412 Potsdam, Germany

²TU Wien, Climate and Environmental Remote Sensing Group, Department of Geodesy and Geoinformation,
Gusshausstraße 25–29, 1040 Vienna, Austria

³Max Planck Institute for Biogeochemistry, Hans-Knöll-Str. 10, 07745 Jena, Germany

⁴Humboldt Universität zu Berlin, Department of Geography, Unter den Linden 6, 10099 Berlin, Germany

⁵Technical University of Munich, School of Life Sciences Weihenstephan, 85354 Freising, Germany

⁶CSIRO Agriculture & Food, 306 Carmody Rd, St Lucia QLD 4067, Australia

Correspondence: Sibyll Schaphoff (sibyll.schaphoff@pik-potsdam.de)

Received: 21 June 2017 – Discussion started: 27 July 2017

Revised: 26 February 2018 – Accepted: 5 March 2018 – Published: 12 April 2018

Abstract. The dynamic global vegetation model LPJmL4 is a process-based model that simulates climate and land use change impacts on the terrestrial biosphere, agricultural production, and the water and carbon cycle. Different versions of the model have been developed and applied to evaluate the role of natural and managed ecosystems in the Earth system and the potential impacts of global environmental change. A comprehensive model description of the new model version, LPJmL4, is provided in a companion paper (Schaphoff et al., 2018c). Here, we provide a full picture of the model performance, going beyond standard benchmark procedures and give hints on the strengths and shortcomings of the model to identify the need for further model improvement. Specifically, we evaluate LPJmL4 against various datasets from in situ measurement sites, satellite observations, and agricultural yield statistics. We apply a range of metrics to evaluate the quality of the model to simulate stocks and flows of carbon and water in natural and managed ecosystems at different temporal and spatial scales. We show that an advanced phenology scheme improves the simulation of seasonal fluctuations in the atmospheric CO₂ concentration, while the permafrost scheme improves estimates of carbon stocks. The full LPJmL4 code including the new developments will be supplied open source through <https://gitlab.pik-potsdam.de/lpjml/LPJmL>. We hope that this will lead to new model developments and applications that improve the model perfor-

mance and possibly build up a new understanding of the terrestrial biosphere.

1 Introduction

The terrestrial biosphere is a central element in the Earth system supporting ecosystem functioning and also providing food to human societies. Dynamic global vegetation models (DGVMs) have been developed and used to study biosphere dynamics under climate and land use change. LPJmL4 is a DGVM with managed land that has been developed to investigate the potential impacts of climate change on the terrestrial biosphere, including natural and managed ecosystems, and is now described in full detail in the companion paper (Schaphoff et al., 2018c). LPJmL and its predecessors were originally benchmarked against ecosystem carbon and water fluxes and global maps of vegetation distribution (Sitch et al., 2003), run-off (Gerten et al., 2004), agricultural yield statistics (Bondeau et al., 2007), satellite observations of fire activity (Thonicke et al., 2001, 2010), permafrost distribution and active layer thickness (Schaphoff et al., 2013), satellite observations of fraction of absorbed photosynthetically active radiation (FAPAR) and albedo (Forkel et al., 2014, 2015), and atmospheric CO₂ concentrations (Forkel et al., 2016). These previous evaluation studies focussed on single

processes or components of the model. Here we present a comprehensive multi-sectoral evaluation to demonstrate that LPJmL4 can consistently represent multiple aspects of biosphere dynamics.

LPJmL4 spans a wide range of processes (from biogeochemical to ecological aspects, from leaf-level photosynthesis to biome composition) and combines natural ecosystems, terrestrial water cycling, and managed ecosystems in one consistent framework. As such, it is increasingly applied for cross-sectoral studies, such as the quantification of planetary boundaries (Steffen et al., 2015) and SDG interactions (Jägermeyr et al., 2017), and the multidimensional impacts of climate and land use change (e.g. Gerten et al., 2013; Ostberg et al., 2015; Warszawski et al., 2014; Zscheischler et al., 2014; Müller et al., 2016). With this complexity, its evaluation against historical observations along multiple dimensions is essential (Harrison et al., 2016). For such a purpose, standardized benchmarking systems have been proposed (Luo et al., 2012; Kelley et al., 2013; Abramowitz, 2005) and iLAMB (<https://www.ilamb.org/>), the international land model benchmarking project, has been established. In the present evaluation of a broad range of fundamental features of the LPJmL4 model, we basically follow the benchmarking procedures, variables, performance metrics, and diagnostic plots suggested by Luo et al. (2012) and Kelley et al. (2013). Thus the presented evaluation goes well beyond earlier evaluations of DGVMs and LPJmL (and its predecessors) itself. We pay special attention to LPJmL4's capability to reproduce observed seasonal and inter-annual dynamics and patterns of key biogeochemical, hydrological, and agricultural processes at various spatial scales. In doing so, we highlight the model's unique feature of representing the interaction of processes for both natural and agricultural ecosystems in a single, internally consistent framework.

2 Model benchmark

In the following we describe in detail the model benchmarking scheme employed here, which allows for a consistent evaluation of processes simulated by LPJmL4 at seasonal and annual resolution and at spatial scales from site level (using e.g. eddy flux measurements for comparison) to global level (using e.g. remote sensing products). The evaluation spans the time period from 1901 to 2011. The benchmarking analysis also considers results from different model setups and previous model versions in order to demonstrate advancements achieved with the current LPJmL4 version and the sensitivity of results to individual new modules.

2.1 Model set-up and simulation experiments

As described in Schaphoff et al. (2018c), we drive the model simulations with observation-based monthly input data on daily mean temperatures from the Climatic Research

Unit (CRU TS version 3.23, University of East Anglia Climatic Research Unit, 2015; Harris et al., 2014) and precipitation provided by the Global Precipitation Climatology Centre (GPCC Full Data Reanalysis version 7.0; Becker et al., 2013). Shortwave downward radiation and net downward longwave radiation are reanalysis data from ERA-Interim (Dee et al., 2011). Monthly average wind speeds are based on the National Centers for Environmental Prediction (NCEP) reanalysis data and were regridded to CRU (NOAA-CIRES Climate Diagnostics Center, Boulder, Colorado, USA; Kalnay et al., 1996b). The number of wet days per month, which is used to allocate monthly precipitation data to individual days of the corresponding months, is derived synthetically as suggested by New et al. (2000). Dew-point temperature is approximated from daily minimum temperature (Thonicke et al., 2010). Global annual values for atmospheric carbon dioxide concentration are taken from the Mauna Loa station (Tans and Keeling, 2015). The spatial resolution of all input data is 0.5° and the model simulations are conducted at this spatial resolution. All model simulations are based on a 5000-year spin-up simulation after initializing all pools to zero. A second spin-up simulation of 390 years is conducted, in which human land use is introduced in 1700, using the data of Fader et al. (2010). In addition to the original dataset description of Fader et al. (2010), sugar cane is now represented explicitly. Cropping intensity as calibrated following Fader et al. (2010) is kept static in the simulations, whereas sowing dates are computed dynamically as a function of climatic conditions until 1971 following Waha et al. (2012) and kept static afterwards. Soil texture is given by the Harmonized World Soil Database (HWSD) version 1 (FAO/IASA/ISRIC/ISSCAS/JRC, 2012; Nachtergaele et al., 2008) and parameterized based on the relationships between texture and hydraulic properties from Cosby et al. (1984). The river-routing scheme is from the simulated Topological Network (STN-30) drainage direction map (Vorosmarty and Fekete, 2011). Reservoir parameters are taken from Biemans et al. (2011), and locations are obtained from the GRanD database (Lehner et al., 2011). We test the influence of specific processes that have been implemented or improved in LPJmL4 (specifically, permafrost, phenology, and fire) on overall model performance by conducting the following factorial experiments.

- LPJmL4-GSI-GlobFIRM is a simulation with all standard model features enabled as used in Schaphoff et al. (2018c), i.e. with land use, permafrost dynamics, the growing season index (GSI) phenology scheme, and the simplified fire model (GlobFIRM). This model experiment is the default LPJmL4 model experiment.
- LPJmL4-GSI-GlobFIRE-PNV is the same, but for potential natural vegetation (PNV) to evaluate the role of managed land in global pattern and processes. This model experiment mimics the original LPJ model (i.e. without agriculture) but with improved phenology.

- LPJmL4-NOGSI-GlobFIRM is a simulation with land use, permafrost dynamics, and the simplified fire model, but without the GSI phenology for testing the sole effect of the GSI phenology. Instead of the GSI phenology, here we use the original phenology model (Sitch et al., 2003) that is based on a growing-degree day approach. This experiment mimics the LPJmL 3.5 version (including the LPJ core, agriculture, and permafrost) as described in Schaphoff et al. (2013).
- LPJmL4-NOGSI-NOPERM-GlobFIRM is a simulation with land use and the simplified fire model but without permafrost and without the GSI phenology. This model experiment mimics the original LPJmL 3.0 model with the LPJ core (Sitch et al., 2003) and the agricultural modules (Bondeau et al., 2007).
- LPJmL4-GSI-SPITFIRE is a simulation set-up as LPJmL4-GSI-GlobFIRM but with the process-based fire model (SPITFIRE; Thonicke et al., 2010). This experiment is an LPJmL4 model run with an alternative fire module.

2.2 Evaluation datasets

Following Kelley et al. (2013) we compare LPJmL4 simulations against independent data for vegetation cover, atmospheric CO₂ concentrations, carbon stocks and fluxes, fractional burnt area, river discharge, and FAPAR. Beyond these suggestions of Kelley et al. (2013), we extend the benchmarking system to datasets of eddy flux tower measurements of evapotranspiration and net ecosystem exchange rate (NEE). Ecosystem respiration (R_e) is evaluated against both eddy flux measurements and operational remote sensing data. Crop yields are evaluated against FAOSTAT data (FAO-AQUASTAT, 2014). For FAPAR, we use not just one but three different reference datasets to account for uncertainties from multiple satellite datasets (see Sect. 2.2.6). We also compare LPJmL4 results against data that are not fully independent of other models (mostly empirical, data-driven modelling concepts), acknowledging the limitations of these data in a benchmark system. However, this allows for the assessment of LPJmL4's performance in additional aspects for which fully data-based products are not available. These data comprise global gridded datasets of vegetation or aboveground biomass carbon (Carvalhais et al., 2014; Liu et al., 2015), cropping calendars (Portmann et al., 2010), global gross primary production (GPP) (Jung et al., 2011), R_e (Jägermeyr et al., 2014), soil carbon (Carvalhais et al., 2014), and evapotranspiration (Jung et al., 2011). We use both site-level and global gridded data because they provide complementary information but have different advantages for the comparison with simulated data like those from LPJmL4. Site-level data are fully independent from model estimates and assumptions, but typically only represent a specific ecosystem with a certain vegetation and soil type

and a specific site history. Thus site-level data have only a limited representativeness for 0.5° grid cells. On the other hand, global gridded data of GPP (Beer et al., 2010; Jung et al., 2011) and R_e (Jägermeyr et al., 2014) are available at the same scale and thus can be directly compared to simulation outputs of DGVMs. However, global gridded datasets usually rely on empirical modelling approaches and ancillary data to upscale and extrapolate site-level data to large regions. Nevertheless, specific site conditions like forest management affecting site age, biomass, and carbon fluxes can hardly be re-simulated for a large number of global sites within a DGVM. Although Kelley et al. (2013) reject the use of such datasets for model benchmarking because they depend on modelling approaches, we accept the additional use of such datasets because they prevent the scale mismatch between site-level data and global DGVM simulations.

2.2.1 Vegetation cover

We compare simulated vegetation cover to the ISLSCP II vegetation continuous fields of Defries and Hansen (2009) as suggested by Kelley et al. (2013). This dataset is a gridded snapshot of vegetation cover for the years 1992–1993 from remote sensing data and distinguishes bare soil, herbaceous, and tree cover fractions aggregated to 0.5° resolution (Defries and Hansen, 2009; Kelley et al., 2013). Tree cover fractions are further distinguished into evergreen vs. deciduous and into broadleaved vs. needle-leaved tree types, respectively. The herbaceous vegetation class includes woody vegetation that is less than 5 m tall. Data uncertainties increase in regions where tree cover is < 20 % due to understorey vegetation and soil disturbing the signal, as well as above 80 % due to signal saturation (Defries and Hansen, 2009; Kelley et al., 2013). To test if the simulated land cover of LPJmL4 performs better than a randomly generated land cover distribution we compare the performance of LPJmL4 to the random model as suggested by Kelley et al. (2013, Sect. 2.3.5), whereas in the original dataset ISLSCP II vegetation continuous fields were randomly resampled.

2.2.2 Atmospheric CO₂ concentration

To evaluate the model's capacity to capture global-scale, intra- and inter-annual fluctuations of atmospheric CO₂ concentrations as driven by the uptake activity of the terrestrial biosphere, we compare simulated CO₂ concentrations with those recorded continuously at two remote measurements at Mauna Loa (MLO; 19.53° N, 155.58° W) and Point Barrow (BRW; 71.32° N, 156.60° W; see Rödenbeck, 2005 for further details on these measurements). We use monthly CO₂ concentrations from flasks and continuous measurements from 1980 to 2010 for the comparison with LPJmL4 simulations. CO₂ observations were temporally smoothed and interpolated using a standard method (Thoning et al., 1989). The atmospheric transport model (TM3; Rödenbeck et al., 2003)

in Jacobian representation (Kaminski et al., 1999) simulates the global CO₂ transport using estimates of net biome production (NBP; here simulated by LPJmL4; see Forkel et al., 2016), estimated net ocean CO₂ fluxes from the Global Carbon Project (Le Quéré et al., 2015) and fossil fuel emissions from the Carbon Dioxide Information Analysis Center (CDIAC; Boden et al., 2013). Atmospheric transport in TM3 is driven by wind fields of the NCEP reanalysis (Kalnay et al., 1996a) at a spatial resolution of 4° × 5°.

2.2.3 Terrestrial carbon stocks and fluxes

Model-independent reference data for carbon stocks and fluxes are available from Luyssaert et al. (2007) for various sites globally distributed. This dataset is comprised of vegetation carbon, aboveground biomass, GPP, and net primary production (NPP). GPP flux data from Luyssaert et al. (2007) are based on eddy flux measurements and are subject to the uncertainties reported in Luyssaert et al. (2007, Table 2). Contrastingly, NPP data are derived from direct measurements of continuous leaf-litter collection, allometry-based estimates of stem and branch NPP from basal measurements, root NPP estimates from soil cores, mini-rhizotrons or soil respiration, and destructive understorey harvest. Estimates here are subject to uncertainties, depending on the sampling methods (Luyssaert et al., 2007). Several individual sites of this dataset can be located within one simulation unit of a 0.5° grid cell and we thus compare simulated values to the range of site measurements in that grid cell. Alternatively to the site-based GPP data from Luyssaert et al. (2007), we also compare spatial patterns and grid-cell-specific GPP simulations to the GPP dataset of Jung et al. (2011), as also suggested by Kelley et al. (2013). This global dataset is based on a larger set of eddy flux tower measurements than the dataset of Luyssaert et al. (2007), but uses additional satellite and climate data and empirical modelling for extrapolation to full global coverage. R_e is evaluated for the time period 2000 to 2009 directly against plot-scale FLUXNET (<http://fluxnet.fluxdata.org/data/la-thuile-dataset/>) measurements (ORNL DAAC, 2011), but also against large-scale R_e estimates from an empirical model based on operational remote sensing data by the Moderate Resolution Imaging Spectroradiometer (MODIS) with a resolution of 1 km and 8 days (Jägermeyr et al., 2014). In addition to GPP, R_e , and NPP, we also compare simulated NEE fluxes with eddy flux tower measurements directly. We use 70 time series of estimated NEE from eddy flux tower sites that measure the exchanges of carbon and water fluxes continuously over a broad range of climate and biome types (ORNL DAAC, 2011). Nevertheless, eddy flux tower sites are not well distributed across the globe and sites in the temperate and boreal zone are better represented than the tropical zone. For the global comparison of the soil and vegetation carbon stocks we use the data compiled by Carvalhais et al. (2014). The soil organic carbon (SOC) estimations are based on the Harmonized World Soil

Database (HWSD) (Nachtergaele et al., 2008). Carvalhais et al. (2014) used an empirical model to calculate SOC stocks (kg m⁻²) from soil organic content (%), layer thickness (m; here for the first 3 m), gravel content (% vol), and bulk density (kg m⁻³). They pointed out that regions such as North America and northern Eurasia are less reliable as HWSD was a work in progress at that time. The vegetation carbon data of Carvalhais et al. (2014) are based on a forest biomass map for temperate and boreal forests from microwave satellite observations (Turner et al., 2014), a biomass map for tropical forests based on lidar observations (Saatchi et al., 2011), and an additional estimate of grassland biomass. Uncertainties in biomass are in most regions between 30 and 40 % and are strongly related to uncertainties in belowground biomass. We also compare simulated aboveground biomass to the estimates of Liu et al. (2015), which is also based on satellite-based passive microwave data. This comparison requires additional assumptions on the separation of aboveground and belowground biomass in LPJmL4 simulations. Liu et al. (2015) estimate for 2000 a global aboveground biomass of 362 Pg C with a 90 % confidence interval of 310–422 Pg C.

2.2.4 Terrestrial water fluxes

River discharge measurements are taken from the ArcticNET (<http://www.r-arcticnet.sr.unh.edu/v4.0/index.html>) and UNH/GRDC (<http://www.grdc.sr.unh.edu/index.html>) datasets for 287 gauges (Vörösmarty et al., 1996). From this database, we only selected river gauges with catchment areas ≥ 10 000 km² as the model set-up and resolution are not suitable for comparison with smaller catchments. We also only selected river gauge records with a temporal coverage of more than 95 % of the observation period and an observation period longer than 2 years at a monthly resolution.

Evapotranspiration fluxes are taken from the FLUXNET database (<http://fluxnet.fluxdata.org/data/la-thuile-dataset/>) and comprise 126 sites, of which we selected sites ($n = 99$) with at least 3 years of data available (ORNL DAAC, 2011). Additional to site-level data, we used global gridded ET data from Jung et al. (2011), which is based on an upscaling of site-level eddy covariance observations with satellite and climate data using a machine learning approach.

The irrigation withdrawal and consumption data that we compare to are from other modelling approaches. Nonetheless, human water use for irrigation is an important component in the terrestrial water cycle and we discuss modelled LPJmL4 estimates in comparison to other model-based estimates, acknowledging the limitation of this comparison and addressing different sources of uncertainty.

2.2.5 Permafrost

For the evaluation of simulated permafrost dynamics, we use the measured thaw depth data from 131 stations of the Cir-

cumpolar Active Layer Monitoring (CALM) station dataset (<https://www2.gwu.edu/~calm/>; Brown et al., 2000) and the International Permafrost Association (IPA) Circum-Arctic Map of Permafrost (<http://nsidc.org/data/ggd318>; Brown et al., 1998). The distribution of permafrost is based on regional elevation, physiography, and surface geology. The permafrost extent represents four classes which categorize the percentage of the ground underlain by permafrost (continuous, 90–100 %; discontinuous, 50–90 %; sporadic, 10–50 %; isolated patches of permafrost, 0–10 %).

2.2.6 Fractional area burnt

For the evaluation of simulated fire dynamics, we employ data on fractional area burnt from the Global Fire Emissions Database GFED4 version 4 (GFED4; <http://www.globalfiredata.org/>; Giglio et al., 2013) for the period 1995 to 2014 and Climate Change Initiative (CCI) Fire version 4.1 (<http://cci.esa.int/data>; Chuvieco et al., 2016) for the period 2005 to 2011. Mean annual burnt area was computed for both datasets for the overlapping period (2005–2011). Both datasets are derived from satellite data. Active fire data were used in GFED4 to prolong the dataset prior to the MODIS period (i.e. for 1995–2000).

2.2.7 Fraction of absorbed photosynthetic active radiation and albedo

Data on the fraction of absorbed photosynthetically active radiation (FAPAR) are derived from three different satellite datasets to account for differences between datasets for model evaluation (see Table 4, Forkel et al., 2015): the MODIS (USGS, 2001) FAPAR (Knyazikhin et al., 1999), the Geoland2 BioPar (GEOV1) FAPAR dataset (Baret et al., 2013) (hereafter called VGT2 FAPAR), and the GIMMS3g FAPAR dataset (Zhu et al., 2013). The MODIS FAPAR dataset is taken from the MOD15A2 product with a temporal resolution of 8 days at a spatial resolution of 1 km, covering the period 2001 to 2011. VGT2 is based on SPOT VGT with a temporal resolution of 10 days and 0.05° spatial resolution (Baret et al., 2013), covering the period 2003 to 2011. The GIMMS3g dataset has a 15-day temporal resolution and 1/12° spatial resolution and covers the period from 1982 to 2011. Data on FAPAR are also subject to uncertainties from the processing of the remotely sensed data and are not available continuously for all areas. We compare the spatial patterns of the peak FAPAR, the temporal dynamics of FAPAR in each grid cell, and seasonal variations in FAPAR averaged for Köppen–Geiger climate zones for the three different FAPAR datasets. The aggregated FAPAR represents the average monthly time series for all grid cells that belong to a certain Köppen–Geiger climate zone (see also Forkel et al., 2015). For the Köppen–Geiger climate zones, FAPAR time series are averaged over all grid cells that belong to the same Köppen–Geiger climate zone (see also Forkel et al.,

2015). For the evaluation of the reflectance of the Earth's surface we used the MODIS C5 albedo time series dataset (https://lpdaac.usgs.gov/dataset_discovery/modis/modis_products_table/mcd43c3) from 2000–2010 (Lucht et al., 2000; Schaaf et al., 2002), which we also aggregated to Köppen–Geiger climate zones for the evaluation here.

2.2.8 Agricultural productivity

Detailed data on crop growth and productivity are available for individual sentinel sites (Rosenzweig et al., 2014). For global-scale or regional simulations, reference data are available only for crop yields and in (sub-)national aggregations (e.g. FAO-AQUASTAT, 2014) or as processed and interpolated gridded products (Iizumi et al., 2014). In all yield data statistics outside of well-controlled field experiments, yield levels and inter-annual variability are not only affected by variability in weather, but also by variance in management conditions, such as sowing dates, variety choices, cropping areas, fertilizer inputs, and pest control (Schauberger et al., 2016). Consequently, it is difficult to evaluate model performance from a comparison of simulated yields with static assumptions on most management aspects with yield statistics in which the contribution of weather variability to yield variability is unknown. Müller et al. (2017) propose a combination of global gridded crop model simulations and different observation-based yield datasets to establish a benchmark for global crop model evaluation. Generally, global gridded crop models perform well in most regions for which statistical models can detect a significant influence of weather on crop yield variability (Ray et al., 2015). We here evaluate LPJmL4 by comparing the simulated and observed yield variability of the 10 top-producing countries of the respective crop (FAO-AQUASTAT, 2014). We refrain from comparing to individual sentinel sites, but refer to the evaluation of LPJmL crop simulations at global, national, and grid cell scale in the global gridded crop model evaluation framework (Müller et al., 2017). As in Müller et al. (2017), we aggregate simulated grid-cell-level yield time series to average national yield time series using the MIRCA2000 dataset for spatial aggregation (Porwollik et al., 2016) and removing trends in observations and simulations with a moving-window average (see Müller et al., 2017, for details). The productivity of biomass plantations is evaluated with data from experimental sites for *Miscanthus*, switchgrass, poplar, willow, and eucalyptus production using the data collection of Heck et al. (2016). Data on biomass productivity typically report a data range. These are site-specific management differences and reflect the diverse drivers of reported productivity, such as variation in plant species, fertilizer use and irrigation management, crop spacing, or sapling size. We average the minimum and maximum values to derive the mean productivity per site.

2.2.9 Sowing dates

To evaluate the accuracy of the simulated rain-fed sowing dates, we use the global dataset of growing areas and growing periods, MIRCA2000 (Portmann et al., 2008, 2010), at a spatial resolution of 0.5° and a temporal resolution of 1 month, as proposed by Waha et al. (2012). Monthly data in MIRCA2000 were converted to daily data by assuming that the growing period starts on the first day of the month following Portmann et al. (2010). MIRCA2000 reports several growing periods in a year for some administrative units for the crops wheat, rapeseed, rice, cassava, and maize. For comparison we select the best corresponding growing period so that a close agreement indicates that simulated sowing dates are reasonable, but not necessarily the most frequently chosen by farmers. We do not compare simulated sowing dates for sugar cane (see Fig. S94 in the Supplement) to observed sowing dates, as MIRCA2000 assumes it is grown all year-round as a perennial crop.

2.3 Evaluation metrics

We employ Taylor diagrams (Taylor, 2001) to compare the correlation, differences in standard deviation, and the centred root mean squared error (CRMS) between simulated and observed carbon and water fluxes at FLUXNET sites (ORNL DAAC, 2011) and at gauge stations from ArcticNET and UNH/GRDC. The standard deviations of the reference datasets have been normalized to 1.0 so that multiple sites can be displayed in one figure.

For global gridded reference datasets, such as for carbon stocks, we show spatial patterns in maps and aggregations as latitudinal means and quantify overall differences as a spatial correlation analysis over all grid cells (see Table 4). As suggested by Kelley et al. (2013) we use the normalized mean square error (NMSE) to describe differences between model simulation and reference datasets. The NMSE is zero for perfect agreement, 1.0 if the model is as good as using the data mean as a predictor, and larger 1.0 if the model performs less well than that. The squared error term puts stronger emphasis on large deviations between simulations and observations and is thus stricter than the normalized mean error (see Table 1 for equations). Kelley et al. (2013) also suggest using the normalized mean error (NME) as a more robust metric than NMSE. NME is based on absolute residuals (NMSE on squared residuals) and thus is especially better suited for variables that can have very large values and residuals. Additionally, we use the Manhattan metric (MM) proposed by Kelley et al. (2013) for evaluation of vegetation cover. Values for MM less than 1 reflect the fact that the model performs better than the mean value. Additionally, we show the random model, which was generated by bootstrap resampling of the observations as proposed by Kelley et al. (2013, Table 4). The random model was used for the evaluation of vegetation distribution. Table 2 gives an overview of variables evaluated at

Table 1. Evaluation metrics used in this study.

Metric	Equation	Reference
NMSE	$\text{NMSE} = \frac{\sum_{i=1}^N (y_i - x_i)^2}{\sum_{i=1}^N (x_i - \bar{x})^2}$	Kelley et al. (2013)
NME	$\text{NME} = \frac{\sum_{i=1}^N y_i - x_i }{\sum_{i=1}^N x_i - \bar{x} }$	Kelley et al. (2013)
ME	$\text{ME} = \frac{\sum_{i=1}^N y_i - x_i \cdot A_i}{\sum_{i=1}^N A_i}$	
W	$\text{W} = 1 - \frac{\sum_{i=1}^N (y_i - x_i)^2 \cdot A_i}{\sum_{i=1}^N (y_i - \bar{x} + x_i - \bar{x})^2 \cdot A_i}$	Willmott (1982)
MM	$\text{MM} = \frac{\sum_{i=1}^N q_{i,j} - p_{i,j} }{N}$	Kelley et al. (2013)

Note: y_i is the simulated and x_i the observed value in grid cell i , \bar{x} the mean observed value, A_i the area weight in grid cell i , and N the number of grid cells or sites; $q_{i,j}$ is the simulated and $p_{i,j}$ is the observed fraction of item j in grid cell i . Normalized mean square error – NMSE, normalized mean error – NME, ME – mean absolute error, W – Willmott coefficient of agreement, MM – Manhattan metric.

the local scale and the measures that were used for the evaluation of time series for crop yields. We employ a simple time series correlation analysis after removing trends with a moving-window detrending method. For comparison with point measurements, we extract the time series from corresponding 0.5° grid cells. These simulated time series may differ in terms of weather and soil conditions from the actual site as the simulations are based on gridded global dataset inputs. The time period is given by the respective measurements, which differ for each observation point.

To envisage the degree of agreement between simulated (LPJmL4) and observed (MIRCA2000) sowing dates, we follow Waha et al. (2012) and compute two different metrics: the Willmott coefficient of agreement (W) (Willmott, 1982) and the mean absolute error (ME), both weighted by the crop-specific cultivated area according to Portmann et al. (2010). For an overview of all metrics used, see Table 1.

3 Results and discussion

In the following we compare the standard version LPJmL4, which refers to the experiment LPJmL4-GSI-GlobFIRM. In the case of the other experiments we refer to the names defined in Sect. 2.1.

3.1 Vegetation cover

LPJmL4 reproduces the observed vegetation distribution better than the random model (Table 3). LPJmL4 can best reproduce the distinction between bare soil and vegetated areas (MM = 0.22) and between tree-covered areas and areas without trees (MM = 0.31), but with considerably better scores than the random model (MM = 0.56 and 0.54, respectively). Moreover, LPJmL4 simulation results reach the lowest MM scores for the distinction of evergreen vs. deciduous trees (MM = 0.52) and for the distribution and composition of life forms (trees vs. herbaceous vs. bare soil; MM = 0.45); these are substantially better than the random model (MM = 0.87

Table 2. Overview of variables and measures used for the evaluation of LPJmL4 local scale.

Variable	Measure			Reference to figures	Reference	
	CRMSE	Standard deviation	Correlation		Data	Citation
CO ₂			x	Figs. 1, 2	Atmospheric transport ¹	Rödenbeck (2005)
NEE	x	x	x	Fig. 3	FLUXNET ²	ORNL DAAC (2011)
ET	x	x	x	Fig. 7	FLUXNET ²	ORNL DAAC (2011)
NPP			x	Fig. 4d		Luyssaert et al. (2007)
GPP			x	Fig. 4c		Luyssaert et al. (2007)
BIOMASS			x	Fig. 4a, b		Luyssaert et al. (2007)
DISCHARGE	x	x	x	Figs. 8, S19–S66	ArcticNET ³ , UNH/GRDC ⁴	Vörösmarty et al. (1996)

Centred root mean square error (CRMSE). ¹ <http://pubman.mpdl.mpg.de/pubman/faces/viewItemOverviewPage.jsp?itemId=escidoc:1691952>;

² <http://fluxnet.fluxdata.org/data/la-thuile-dataset/>; ³ <http://www.r-arcticnet.sr.unh.edu/v4.0/index.html>; ⁴ <http://www.grdc.sr.unh.edu/index.html>.

Table 3. Comparison metric scores for LPJmL4 simulations against observations of fractional vegetation cover data from International Satellite Land-Surface Climatology Project (ISLSCP) II vegetation continuous field (VCF) (Defries and Hansen, 2009).

Vegetation cover	Manhattan metric (MM)	
	LPJmL4	Random model*
Life forms	0.45	0.88
Tree vs. non-tree	0.31	0.54
Herb vs. non-herb	0.42	0.66
Bare vs. covered ground	0.22	0.56
Evergreen vs. deciduous	0.52	0.87
Broadleaf vs. needle-leaf	0.37	0.94

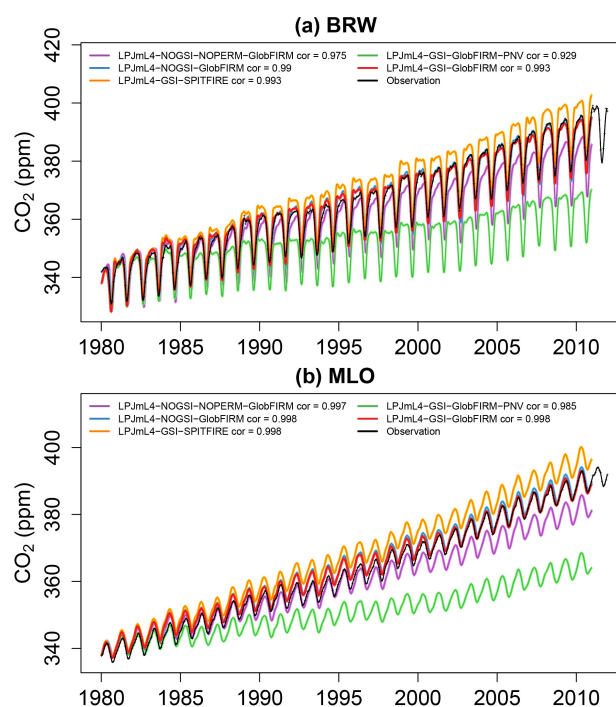
MM suggested by Kelley et al. (2013), * values taken from Kelley et al. (2013, Table 4).

and 0.88, respectively). The largest improvement in LPJmL4 simulations over the random model is found for the patterns of broadleaved vs. needle-leaved trees (MM = 0.37 for LPJmL4 vs. 0.94 for the random model; see Table 3).

3.2 Atmospheric CO₂ concentration and NEE

3.2.1 Comparison of simulated NBP to atmospheric CO₂ concentration at MLO and BRW

LPJmL4 reproduces the observed long-term and seasonal dynamics of atmospheric CO₂ well (Figs. 1 and 2). The long-term trend of atmospheric CO₂ is reproduced well in all the different model set-ups (Fig. 1), except for the set-up with natural vegetation only (LPJmL4-GSI-GlobFIRM-PNV). The experiment with all processes included (LPJmL4-GSI-GlobFIRM) gives the best correlation and trend reproduction, which suggests that an integral representation of the LPJmL4 features is required to match observations best. Next to land use dynamics, the inclusion of permafrost dynamics has the strongest effects on the sim-

**Figure 1.** Comparison of the atmospheric CO₂ concentrations at Point Barrow (BRW; panel a) and Mauna Loa (MLO; panel b) for the different LPJmL4 experiments.

ulated trend (LPJmL4-NOGSI-NOPERM-GlobFIRM vs. LPJmL4-NOGSI-GlobFIRM). The use of the process-based fire model SPITFIRE leads to a small overestimation of the trend in atmospheric CO₂ concentrations compared to the other model set-ups, especially at MLO. Seasonal variations in atmospheric CO₂ can be reproduced well by LPJmL4, especially by the standard set-up (LPJmL4-GSI-GlobFIRM; Fig. 2). The simulation of seasonal variations in atmospheric CO₂ content are especially improved by the GSI phenol-

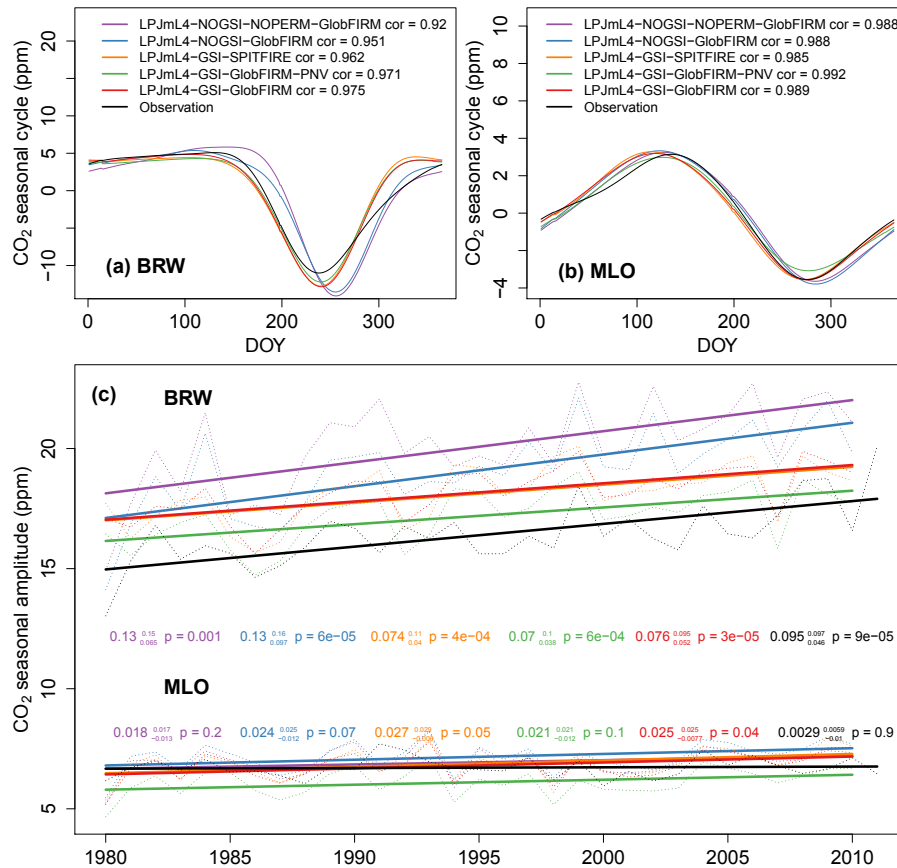


Figure 2. Comparison of the atmospheric CO₂ concentration at Mauna Loa (MLO) and Point Barrow (BRW) simulated in the different LPJmL4 experiments. (a, b) Seasonal cycle, (c) trend of the seasonal amplitude, and slopes are given for the different LPJmL4 experiments.

ogy scheme (LPJmL4-NOGSI-GlobFIRM vs. LPJmL4-GSI-GlobFIRM; Fig. 2a, b). All model set-ups (except LPJmL4-GSI-SPITFIRE) can reproduce the observed strong significant increase in the seasonal CO₂ amplitude at BRW and the weak (and insignificant) increase at MLO (Fig. 2c). These results are in agreement with a previous evaluation of simulated seasonal CO₂ changes in LPJmL (Forkel et al., 2016).

Further analysis shows that the standard set-up (LPJmL4-GSI-GlobFIRM) can best produce the mean seasonal cycle in MLO, whereas the version that omits land use (LPJmL4-GSI-GlobFIRM-PNV) performs slightly better than this in BRW (Fig. 2). The standard set-up (LPJmL4-GSI-GlobFIRM) can also best reproduce the increase in the seasonal amplitude at BRW, whereas it is the only set-up that produces a statistically significant but still very small increase in the seasonal amplitude at MLO where observations also do not show a statistically significant increase.

3.2.2 Comparison of simulated NEE to eddy flux measurements

We evaluate the model performance of simulated NEE from LPJmL4 for temporal and spatial variation in NEE data from eddy flux measurements using Taylor diagrams (Taylor, 2001). Stations are sorted from north to south (see Fig. 3) for all NEE measurements available for > 3 years. The model is able to reproduce the mid-latitudes best (represented by yellow over green to light blue colours), with correlation coefficients mostly between 0.4 and 0.9 and standard deviations often within $\pm 30\%$ of the reference data. The northernmost regions are reproduced well at some flux towers, but often with higher standard deviation than in the flux tower data, which means that the simulated time series are largely in phase but more variable than the observations. In contrast, the evaluation is comparatively poor for tropical regions, especially the station at Santarém with strong negative correlations ($r < -0.6$) but realistic standard deviations. For this site, however, Saleska et al. (2003) have already pointed out that the eddy flux measurements show the opposite sign compared to tree

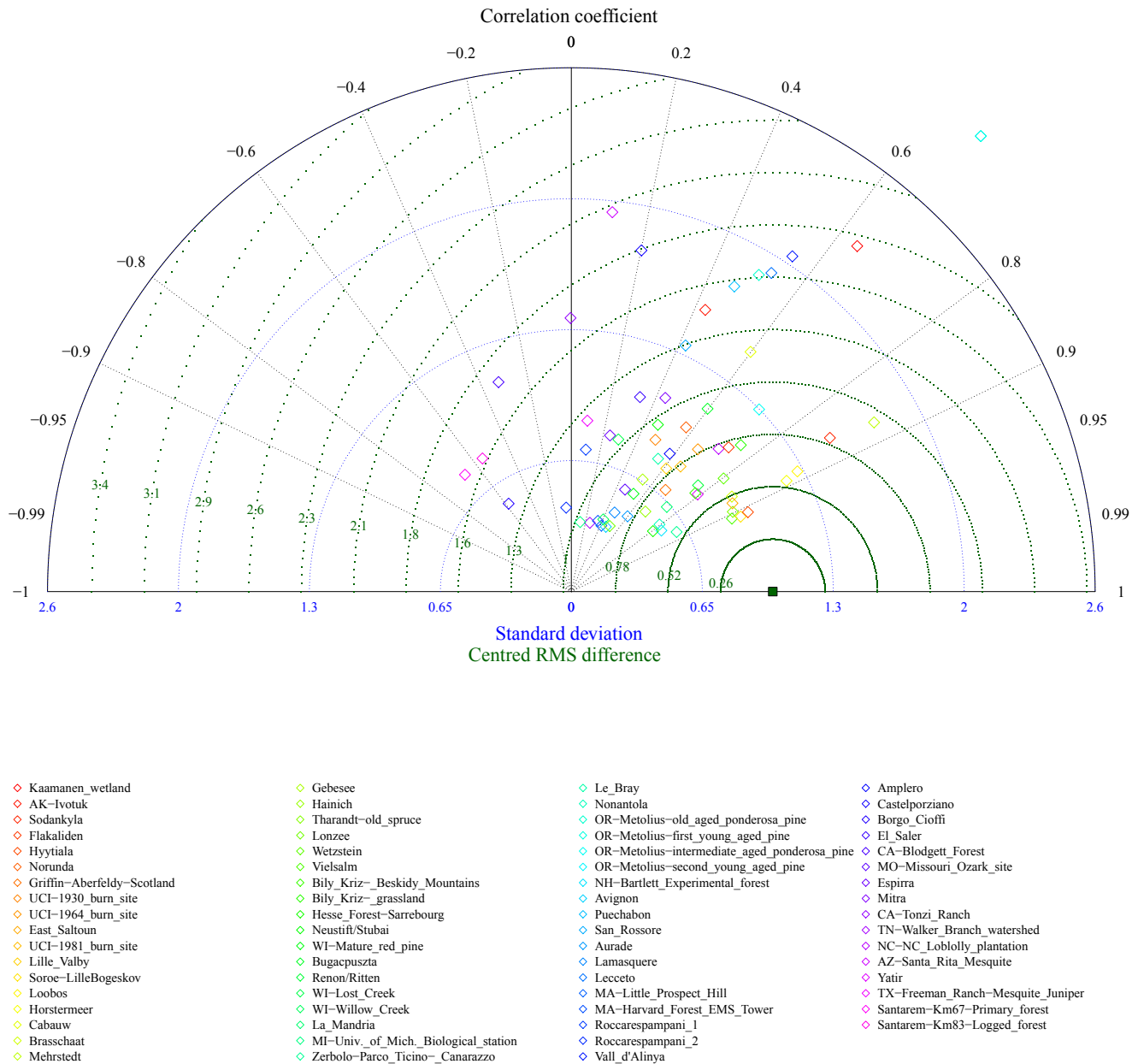


Figure 3. Net ecosystem exchange rate measured at eddy flux towers: ORNL DAAC (2011). Available online at FLUXNET (<http://fluxnet.fluxdata.org/data/la-thuile-dataset/>). Sites (colours) are ordered from north to south.

growth observations and model predictions, which is also the case for LPJmL4. We stress that this evaluation is done for a standard LPJmL4 run and standard input (the LPJmL4-GSI-GlobFIRM as described in Schaphoff et al., 2018c); i.e. we did not calibrate the model to site-specific conditions and also drive the model with gridded input data rather than the observed soil and weather data at individual stations. More detail for comparisons with eddy flux tower measurements

for individual locations is supplied in the Supplement (see Figs. S1–S7). Additionally, we have simulated NEE by conducting simulations with station-specific meteorological observations (see Fig. S17). The results are similar to simulations driven by global climate data.

Table 4. Overview of variables evaluating LPJmL4 showing measures and references at the global scale.

Variable	Measure					Reference	
	NME	NMSE	Spatial correlation	Temporal correlation	Visual comparison	Data	Citation
GPP – Av	0.20	0.13	0.87		Figs. 5, S68	GPP ¹	Jung et al. (2011)
R _e – Av	0.67	0.55	0.67		Figs. 6, S70		Jägermeyr et al. (2014)
SoilC – Av	0.48	0.75	0.29		Fig. S67	Soil carbon stocks ¹	Carvalhais et al. (2014)
VegC – Av	0.33	0.36	0.84		Fig. S69a Fig. S69b	Total biomass ¹ AGB	Carvalhais et al. (2014) Liu et al. (2015)
FAPAR – I-aMv	0.17	0.13	0.63	Fig. 10a		MODIS FAPAR ²	Knyazikhin et al. (1999)
FAPAR – I-aMv	0.18	0.15	0.59	Fig. 10b		GIMMS3g FAPAR ³	Zhu et al. (2013)
FAPAR – I-aMv	0.21	0.20	0.69	Fig. 10c		VGT2 FAPAR ⁴	Baret et al. (2013)
ET	1E-6	0.07	0.84		Fig. S71	Latent heat flux ¹	Jung et al. (2011)
fBA					Fig. S72		GFED4 & CCI Fire (4.1)
Albedo					Fig. S72	MODIS C5	Lucht et al. (2000)
Discharge						ArcticNET ⁵ , UNH/GRDC ⁶	Vörösmarty et al. (1996)
Ov				$R^2 = 0.90$			
Mav				$R^2 = 0.92$			
I-av				$R^2 = 0.97$			

Normalized mean error (NME) and normalized mean square error (NMSE) as suggested by Kelley et al. (2013); Av – annual average; I-aMv – inter-annual monthly variability; overall variability – Ov; monthly average variability – Mav; inter-annual variability – I-av; vegetation carbon – VegC; aboveground biomass – AGB; soil carbon – SoilC; fBA – fractional burnt area.

¹ <https://www.bgc-jena.mpg.de/geodb/BGI/Home>; ² https://pdaac.usgs.gov/dataset_discovery/modis/modis_products_table/mod15a2;

³ <http://cliveg.bu.edu/modismisr/lai3g-fpar3g.html>; ⁴ http://cordis.europa.eu/result/rcn/140496_en.html; ⁵ <http://www.r-arcticnet.sr.unh.edu/v4.0/index.html>;

⁶ <http://www.grdc.sr.unh.edu/index.html>.

3.3 Vegetation and soil carbon stocks and vegetation productivity

3.3.1 Soil carbon and vegetation carbon stocks

The spatial correlation between simulated and observation-based estimates of SOC by Carvalhais et al. (2014) is weak ($r = 0.29$; Table 4) with disagreements in the subtropics where LPJmL4 simulations substantially underestimate soil carbon stocks, whereas LPJmL4 reports much higher soil carbon in the high northern latitudes ($> 50^\circ \text{N}$) and lower values for the tropical and temperate zone compared to Carvalhais et al. (2014) (see Fig. S67). Other estimates by Tarnocai et al. (2009) show much higher carbon content for the permafrost-affected areas than the dataset of Carvalhais et al. (2014). We thus assume that the disagreement between simulations and the Carvalhais et al. (2014) data may also result from an underestimation of carbon stocks in the Carvalhais et al. (2014) data. However, the estimation of global soil carbon is less in LPJmL4 (1869 Pg C) than estimated by Carvalhais et al. (2014) ($2352 \pm 400 \text{ Pg C}$).

The comparison of simulated and observation-based assessments of vegetation carbon show a good spatial correlation ($r = 0.84$; Table 4). Globally, Carvalhais et al. (2014) estimate slightly lower biomass ($445 \pm 8 \text{ Pg C}$) as simulated by LPJmL4 (507 Pg C). The spatial patterns of vegetation car-

bon stocks are shown in Fig. S69a for simulations and the data product of Carvalhais et al. (2014). While the broad geographical patterns are in overall agreement with the evaluation data, the absolute values differ in some regions. Specifically, LPJmL4 simulates much higher biomass (see the latitudinal pattern in Fig. S69) for the tropics and lower biomass between 20 and 40° in the Northern and Southern Hemisphere where Carvalhais et al. (2014) show higher values compared to LPJmL4. This is probably due to an overestimation of vegetation carbon in agricultural regions by Carvalhais et al. (2014), as Liu et al. (2015) show similar aboveground biomass estimates there (see Fig. S69b). The subtropical region where biomass carbon is underestimated also corresponds to the region where LPJmL4 simulations underestimate soil carbon stocks compared to Carvalhais et al. (2014). Also, the comparison of aboveground biomass estimates with the dataset of Liu et al. (2015) shows a similar spatial pattern of overestimation of vegetation biomass with too-high values in boreal and tropical areas. The comparison is complicated by uncertainties in the estimation of belowground biomass (Saatchi et al., 2011) and the assumed distribution between aboveground and belowground biomass in LPJmL4 simulations, in which LPJmL4 assumes that belowground biomass consists of all fine root biomass and one-third of all sapwood biomass. The simulation experiments without permafrost dynamics (LPJmL4-NOGSI-NOPERM-GlobFIRM)

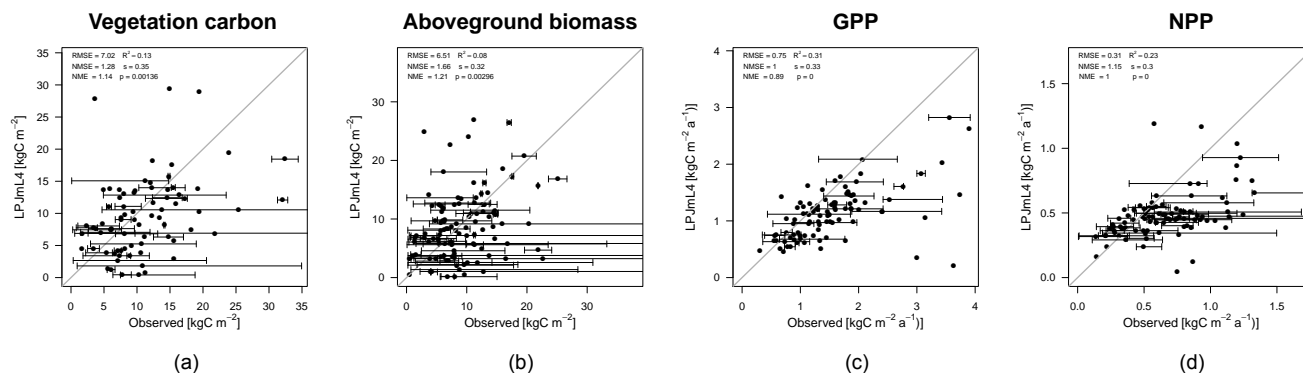


Figure 4. Evaluation of vegetation carbon (a), aboveground biomass (b), GPP (c), and NPP (d). Observed data are provided by Luysaert et al. (2007). Bars give the minimum and maximum of the estimation within one 0.5° cell simulated by LPJmL4.

show a high overestimation of biomass in the high latitudes. Similarly, the inclusion of the GSI phenology substantially reduces the biomass overestimation in comparison to Carvalho et al. (2014) and Liu et al. (2015), which is consistent with the finding of Forkel et al. (2014). The consideration of human land use in the simulations improves carbon stock simulations in the temperate zones (Fig. S69). This clearly demonstrates the importance of permafrost, human land use, and the GSI phenology for the simulation of the terrestrial carbon cycle, even though the remaining discrepancies warrant further model improvement.

Figure 4a and b compare site data estimation with the representative LPJmL4 grid cell estimation with an uncertainty range which comes from the different measurements within one 0.5° grid cell. Both vegetation and aboveground carbon are slightly overestimated in some cases but also strongly underestimated in others. As LPJmL4 calculates a representative mean value of a 0.5° grid cell for all benchmarks, the simulated values should match the mean values. However, it can be assumed that measurements are not evenly distributed through the age classes within one grid cell or forest, and it remains unclear how representative the measurements are for a 0.5° grid cell area.

3.3.2 Gross and net primary production (GPP and NPP)

The global estimation of $123.7 \text{ Pg C a}^{-1}$ GPP from LPJmL4 (see Fig. 5) matches the estimates from Beer et al. (2010) and Jung et al. (2011) of 123 ± 8 and $119 \pm 6 \text{ Pg C a}^{-1}$, respectively, for the years 1982–2005, whereas the highest divergence can be observed in the tropics where LPJmL4 estimates much lower values despite the higher biomass estimations (see Sect. 3.3). LPJmL4 simulated higher GPP for the temperate and boreal zones than reported by Jung et al. (2011). The different model experiments show similar patterns except for LPJmL4-GSI-GlobFIRM-PNV, which shows lower GPP in the Mediterranean (see Fig. 5). Carvalho et al. (2014) estimate global NPP at $54 \pm 10 \text{ Pg C a}^{-1}$

and LPJmL4 at 57 Pg C a^{-1} for the mean of the years 1982–2011.

The site data comparison to Luysaert et al. (2007) shows a good agreement between site measurements and simulated GPP (see Fig. 4c) and NPP (see Fig. 4d). The overestimation of simulated biomass and the good agreement of NPP and GPP leads to the conclusion that LPJmL4 underestimates mortality. This warrants further investigation of why LPJmL4 seems to overestimate global GPP but shows good agreement with site data. The comparison of LPJmL4 against MTE data (Jung et al., 2011) on the local scale for the same points as given by Luysaert et al. (2007) shows a good agreement, especially if outliers are excluded (Fig. S68b, c). Figure S68a compares plot data against the global data.

3.3.3 Ecosystem respiration (R_e)

Comparison of satellite-derived ecosystem respiration with that simulated by LPJmL4 reveals similar spatial patterns (Figs. 6 and S70). However, LPJmL4 shows higher temperature sensitivities (Fig. 6a) and consistently simulates higher R_e values in high-latitude and subtropical regions (Fig. S70). Since satellite-derived ecosystem respiration is calibrated for FLUXNET data and hence exhibits marginal cross-latitude bias, the discrepancies to LPJmL4 are likely associated either with LPJmL4 parameterization or with systematic errors in the FLUXNET processing technique. Additional details and figures are presented in Jägermeyr et al. (2014).

3.4 Water fluxes

3.4.1 Evapotranspiration

The spatial distribution of evapotranspiration in LPJmL4 shows a very similar pattern to that estimated by Jung et al. (2011) (Table 4, Fig. S71). It indicates a general underestimation of ET, especially in the tropics and subtropics, but in most cases within the uncertainty range. This is consistent with the underestimation of GPP in the tropics (Fig. 5), but

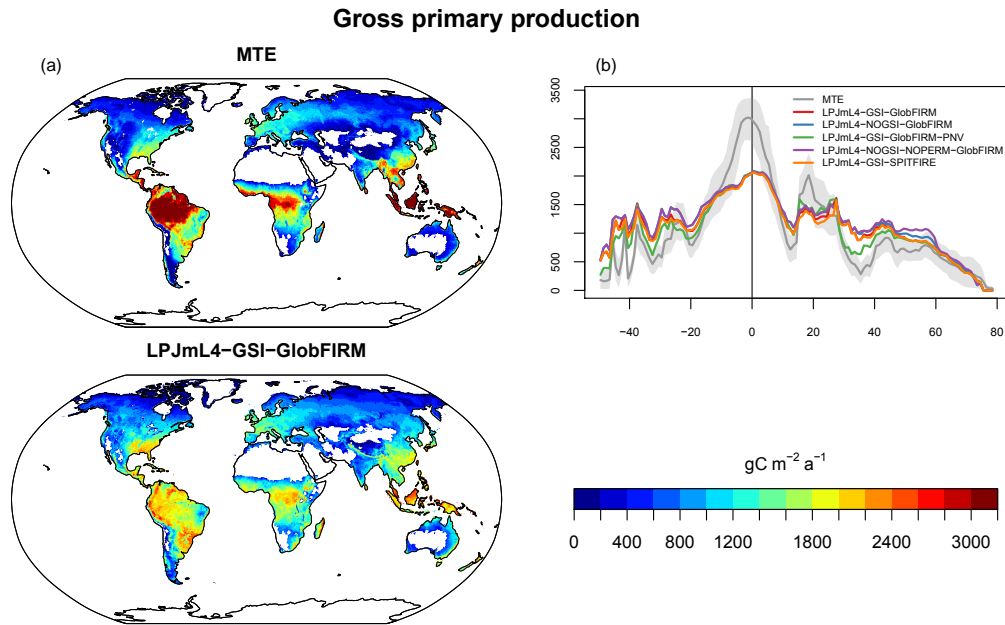


Figure 5. The maps (a) show the spatial pattern of gross primary production (GPP; $\text{g C m}^{-2} \text{a}^{-1}$) distribution from the standard LPJmL4 simulation against the MTE data (Jung et al., 2011). The graph in (b) shows the latitudinal pattern of GPP distribution simulated by the different versions of LPJmL4 against data from Jung et al. (2011).

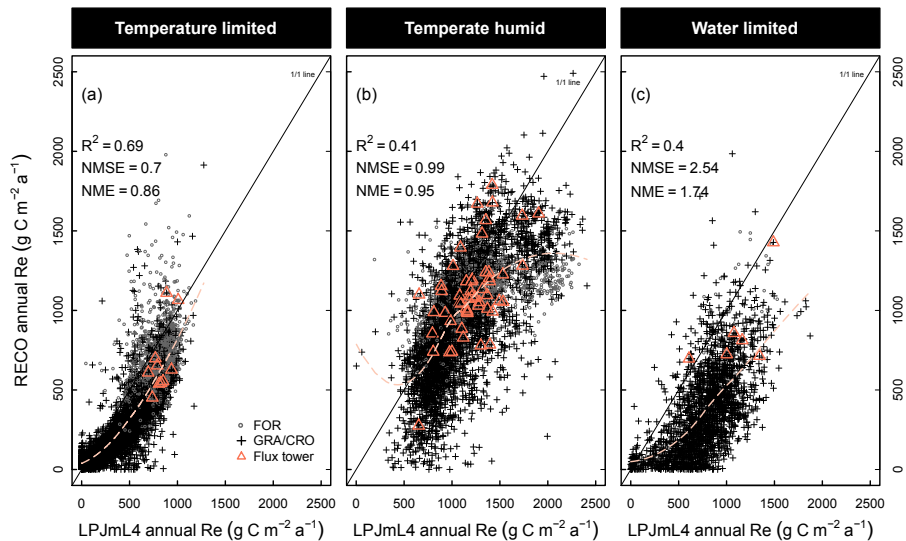


Figure 6. Ecosystem respiration (R_e) evaluation of standard LPJmL4 simulations with satellite-derived estimations from Jägermeyr et al. (2014). Annual R_e sums for all pixels from the displayed extent in Fig. S70 are compared and separated by climate type (a)–(c). Dashed lines indicate a polynomial bias curve. Chart symbols are separated for forest (FOR) and grassland–cropland (GRA–CRO) land cover classes.

not with the general overestimation of vegetation biomass (Fig. S69). The different experiments show nearly no effects on the simulated evapotranspiration. At site level, the evapotranspiration fluxes show a good agreement with eddy flux tower measurements (Fig. 7). LPJmL4 shows good performance in most regions, with correlation coefficients often larger than 0.6. The northern and temperate stations (red

to light blue symbols) show especially high correlation with low CRMS. Simulations of tropical and subtropical ET (dark blue to purple symbols) show weak or even negative correlations coupled with a high CRMS for some stations. We also provide more detailed time series analyses for the evapotranspiration fluxes of individual sites in the Supplement (Figs. S8–S16).

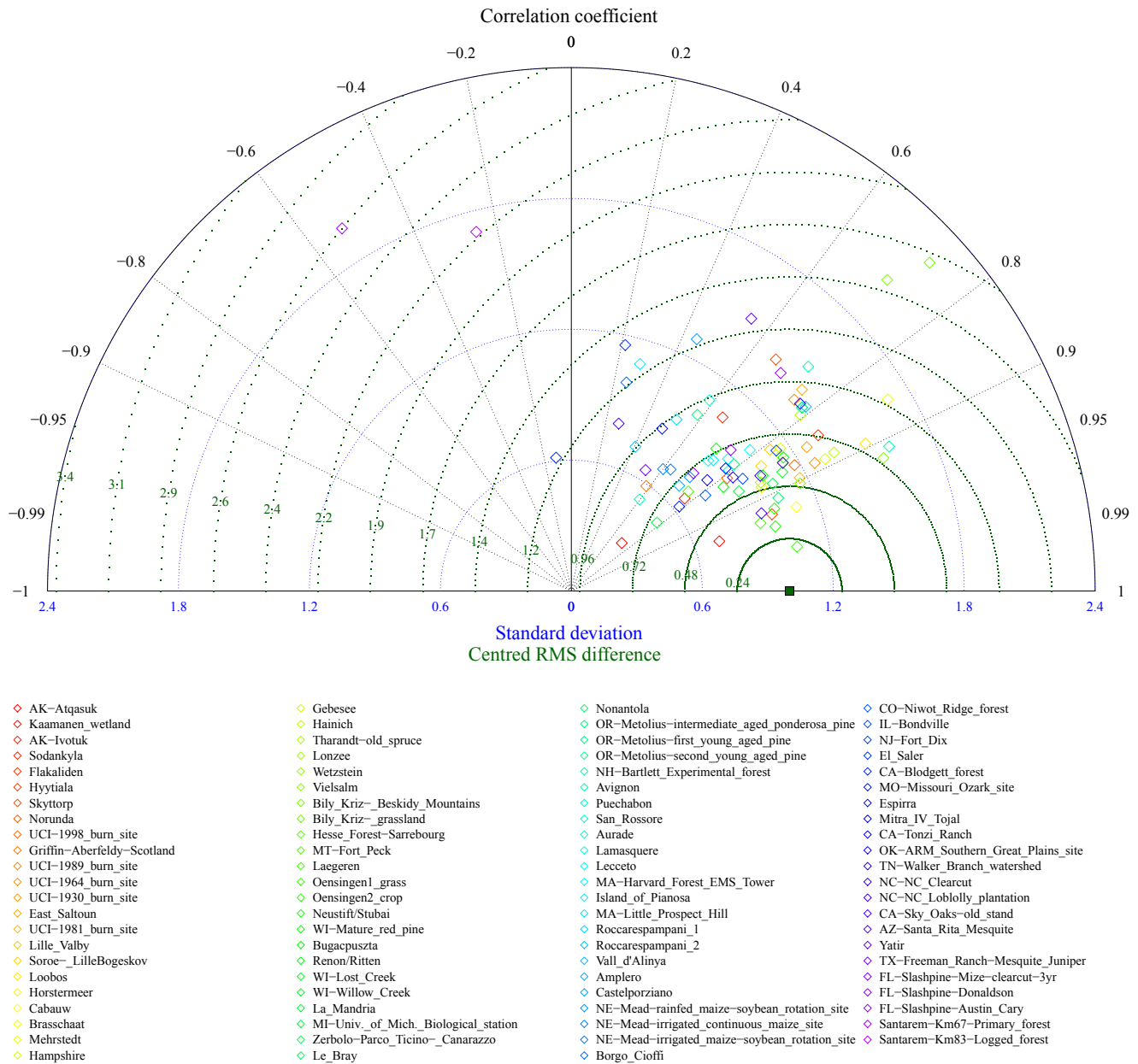


Figure 7. Evaporation rate measured at eddy flux towers: ORNL DAAC (2011). Available online at FLUXNET (<http://fluxnet.fluxdata.org/data/la-thuille-dataset/>). Site locations are ordered from north to south.

3.4.2 River discharge stations evaluation

Discharge simulated by earlier LPJmL versions was previously evaluated in several studies, also in comparison with other global hydrological and land surface models (Haddeland et al., 2011). River discharge was evaluated for major catchments globally, also accounting for the effects of different precipitation datasets (Biemans et al., 2009) and regionally for the Amazon basin (Langerwisch et al., 2013) and the Ganges (Siderius et al., 2013). Figure 8 shows the com-

parison of simulated LPJmL4 and observed river discharge values for all gauges with a basin area $\geq 10\,000\text{ km}^2$. Here, the most northern (blue) and also the most southern (purple) gauges show good agreement, but overall the picture is mixed with respect to correlation coefficients and standard deviation. For further insights, we provide comparisons for all considered gauges in the Supplement (Figs. S19–S66). For many gauges, the simulated seasonal timing of river discharge (peaks) has improved (see Figs. S19–S22) compared to the previous model evaluation of river discharge

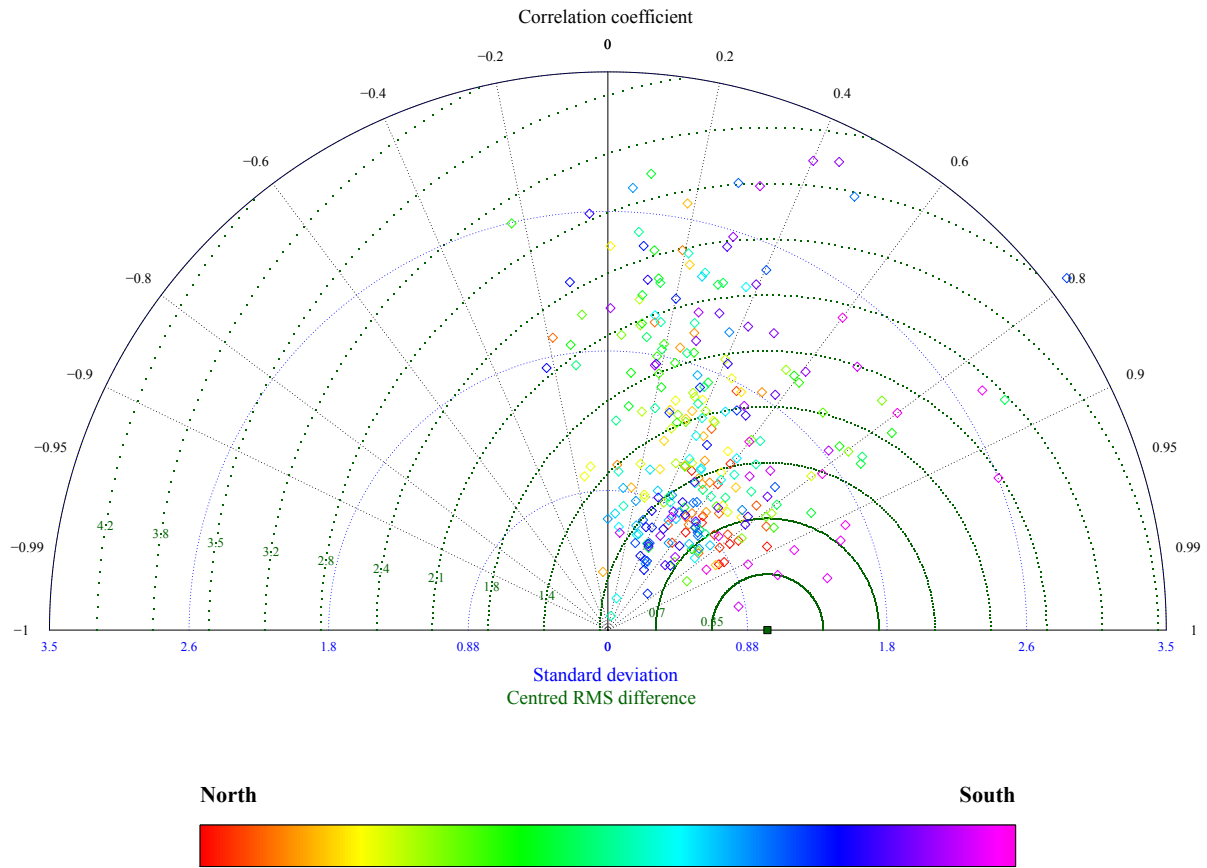


Figure 8. Comparison of simulated discharge with 287 gauges provided by ArcticNET (<http://www.r-arcticnet.sr.unh.edu/v4.0/index.html>) and UNH/GRDC (<http://www.grdc.sr.unh.edu/index.html>). Stations with basin area $\geq 10\,000\text{ km}^2$ are taken into account. Gauges are ordered from north to south (see legend colour).

(Schaphoff et al., 2013), which is mainly a result of the newly implemented GSI phenology scheme (Forkel et al., 2014). The discharge spring peaks in permafrost areas are especially affected by this improvement. At many gauges, LPJmL4 can reproduce the variability for the whole time series and specially the seasonality, with a high R^2 and NME and NMSE, which implies a better performance than the mean model. The dynamics at gauges in the temperate zone (Figs. S49–S50, S61) are not well reproduced in the simulations, and the NME and NMSE also show high values in contrast to gauges in the subtropics and tropics (Figs. S64–S66), which typically show high R^2 and low NME and NMSE.

The evaluation at the global aggregation (computed for all stations and then averaged) shows very high agreement between observed and modelled discharge (see Table 4). Both the explained variance (R^2) and the NME–NMSE contribute to the good performance of the simulated discharge. The constant flow velocity in all rivers, as assumed in LPJmL4 simulations, could be varied by river for further model improvement, especially for the timing in flat areas where wetland dynamics may play an important role.

3.4.3 Irrigation withdrawal and consumption

Global estimates of irrigation water withdrawal (W_d : $2545\text{ km}^3\text{ a}^{-1}$) and consumption (W_c : $1292\text{ km}^3\text{ a}^{-1}$) agree well with previous studies. Reported W_d values for the period 1998–2012 are $2722\text{ km}^3\text{ a}^{-1}$ (FAO-AQUASTAT, 2014), and modelling results range from 2217 to $3185\text{ km}^3\text{ a}^{-1}$ (Döll et al., 2014, 2012; Wada and Bierkens, 2014; Alexandros and Bruinsma, 2012; Wada et al., 2011; Siebert and Döll, 2010). W_c estimations range between 927 and $1530\text{ km}^3\text{ a}^{-1}$ (Chaturvedi et al., 2015; Döll et al., 2014; Hoff et al., 2010). Döll et al. (2012) find that $1179\text{ km}^3\text{ a}^{-1}$ ($1098\text{ km}^3\text{ a}^{-1}$ in Wada and Bierkens, 2014) relates to surface water with an additional $257\text{ km}^3\text{ a}^{-1}$ from groundwater resources. LPJmL4 does not account for fossil groundwater extraction nor desalination. However, previous studies show that 80 % of groundwater withdrawals are recharged by return flows (Döll et al., 2012). It is thus plausible that studies accounting for (fossil) groundwater reach W_d estimates somewhat higher than in LPJmL4. Naturally, irrigation water estimates are associated with uncertainties in the precipita-

tion input employed (Biemans et al., 2009). A representation of multiple cropping systems in LPJmL4 (Waha et al., 2013) and the corresponding growing seasons (Waha et al., 2012) could also help to improve water withdrawal and consumption estimates and eventually river discharge, especially in tropical areas.

Simulated irrigation efficiencies are difficult to compare with observations due to inhomogeneous definitions and field measurement problems. Yet, in Table S1 in the Supplement we relate our results to comparable literature. Our simulations meet the indicative estimates of Brouwer et al. (1989) at the global level. Sauer et al. (2010) provide another independent estimate of field efficiency with global average values of 42, 78, and 89 % for the three irrigation types, respectively. Our estimates agree well with these numbers globally and regionally, even though there are some regional patterns that are not represented in our results. Sauer et al. (2010), for instance, find lower surface irrigation efficiencies in the Middle East, North Africa (MENA), and sub-Saharan Africa (SSA). We simulate above-average efficiencies in MENA and particularly low ones in South Asia, which are both supported by Rosegrant et al. (2002) and Döll and Siebert (2002). Overall, the evaluation of the irrigation model in LPJmL4 demonstrates that it is well in line with reported patterns, and yet it comes with much more detailed depths with respect to process representation and spatio-temporal resolution than these.

3.5 Permafrost distribution and active layer thickness

The current permafrost distribution and the active layer thickness (Fig. 9) is well represented by the LPJmL4 model compared to independent studies (Brown et al., 1998, 2000). LPJmL4 is able to reproduce the distribution of permafrost and the measured active layer thickness in most grid cells. The continuous permafrost zone is characterized by a thawing depth equal to or less than 1 m in LPJmL4, while the model simulates for sporadic permafrost and isolated patches a thawing depth of more than 3 m. The spatial distribution of greater thaw depth from north to south is simulated well by the model. CALM station data show a similar thawing depth as simulated by LPJmL4 (Fig. 9b), but CALM station data also indicate that thawing depth can be different for the same grid cell, as other processes (e.g. exposition) not represented by LPJmL4 can play an important role.

3.6 Fire

3.6.1 Burnt area

Simulated fractional area burnt is largest in the seasonal dry tropics and temperate regions in all model versions and smallest in cold or wet environments (Fig. S72). However, maximum fractional burnt area does not exceed 0.0625 in tropical and subtropical savanna and shrubland areas when

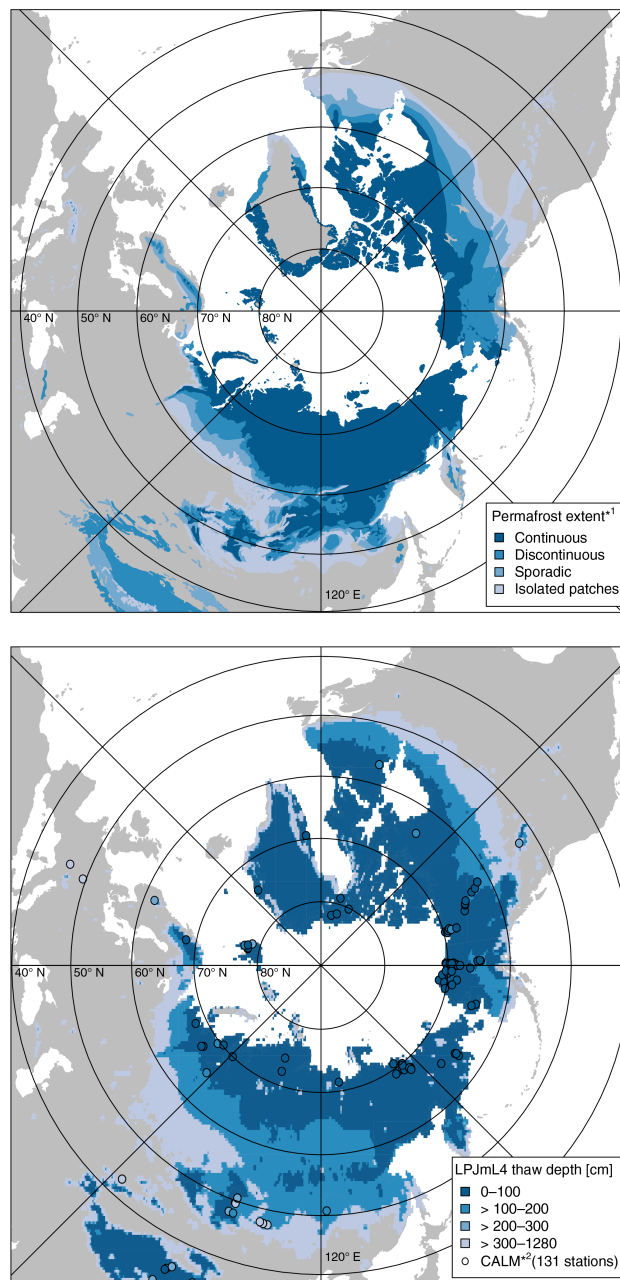


Figure 9. Observed and simulated permafrost distribution and active layer thickness. (a) Contemporary permafrost extent according to the IPA Circum-Arctic Map of Permafrost (*¹Brown et al., 1998). (b) LPJmL4-simulated active layer thickness compared to the *²CALM station data means both for the observation time 1991–2009 (<https://www2.gwu.edu/~calm/>; Brown et al., 2000). The colour scheme used at the bottom is the same for simulated thaw depth and Circumpolar Active Layer Monitoring (CALM) data.

the GlobFIRM model is applied. It is comparable to GFED4 and CCI estimates only in South America, while in other tropical regions GFED4 (Giglio et al., 2013) and CCI report fractional burnt area between 0.125 and 0.75 (Fig. S72). In these regions, the fractional burnt area simulated by the SPITFIRE model is overestimated with values between 0.25 and 1, specifically in Southern Hemispheric Africa and northern Australia. SPITFIRE is very sensitive to vegetation, and thus fuel composition in areas with homogeneous C_4 grasslands can lead to an overestimation of simulated area burnt, which is specifically the case for seasonally dry South America and the Indian subcontinent. LPJmL4-GSI-SPITFIRE captures the distribution of fractional burnt area much better than LPJmL4-GSI-GlobFIRM, which is too homogeneous in its response. In contrast, LPJmL4-GSI-SPITFIRE better captures the very small fractions reported for the wet tropical forests, which is better comparable to GFED4. Here, the approach of simulating fire risk based on the climatic fire danger index instead of deriving a fire probability from the topsoil soil moisture is of great advantage in these regions. While LPJmL4-GSI-GlobFIRM simulates a relatively homogeneous spatial distribution of fractional burnt area in temperate and boreal forest regions, LPJmL4-GSI-SPITFIRE underestimates fractional burnt area in these biomes. LPJmL4-GSI-GlobFIRM underestimates fractional burnt area in the temperate steppe regions, whereas LPJmL4-GSI-SPITFIRE manages to spatially capture the burning conditions in these biomes, even though the total amount is overestimated. The phenology module in LPJmL4 has no effect on the fractional burnt area simulated by LPJmL4-GSI-GlobFIRM, whereas including permafrost increases burnt area in the circum-boreal region, specifically in Siberia, even though the spatial effect is too homogeneous.

3.6.2 Fire effects on biomass and vegetation distribution

Both fire model approaches simulate a comparable latitudinal distribution of biomass starting from the wet tropics towards dry and colder areas in the north and south. Both model versions simulate comparable values in the wet tropics around the Equator and capture the gradient to seasonal dry tropics in the north (until 10° N) and south (until 20° S). The overestimation of burnt area in tropical savannas around 20° N in LPJmL4-GSI-SPITFIRE leads to an underestimation in simulated biomass compared to the other LPJmL4 experiments. The consideration of permafrost and fire dynamics is required to reproduce observed vegetation biomass values in boreal regions.

3.6.3 Global biomass burning

The modelling errors in fractional area burnt compensate in different ways in each fire model. SPITFIRE simulates global biomass burning values of 2.7 Pg C a^{-1} on average between

1996 and 2005, which is comparable to the 2.33 Pg C a^{-1} suggested by Randerson et al. (2015). Here, overestimations of burnt area in tropical savannas and underestimations in boreal forests compensate for each other. GlobFIRM simulates more fires in boreal regions that are less spatially pronounced than in GFED4, but underestimates fractional burnt area in the subtropics and tropics. GlobFIRM therefore estimates global biomass burning by 2.8 Pg C a^{-1} , which is similar to SPITFIRE.

3.7 Fraction of absorbed photosynthetically active radiation (FAPAR) and albedo

Evaluations against multiple satellite datasets of FAPAR have already shown that LPJmL-GSI can reproduce the seasonality of FAPAR and the inter-annual variability and trends well at the start and end of the growing season within observational uncertainties (Forkel et al., 2015). LPJmL4 shows a high spatial correlation with correlation coefficients between 0.6 and 0.71 for PEAK-FAPAR. It shows also a good agreement with the temporal variations (Fig. 10a–c). Large parts of the wet tropics display a negative correlation between simulated and observed FAPAR, which may explain the phase offset in the dynamics of NEE at the station Santarém. However, in these regions the differences between datasets are also large, which is caused by the limitations of optical satellite observations in regions with permanent cloud cover (Forkel et al., 2015).

LPJmL4 reproduces the global patterns of annual peak FAPAR (Fig. 11) well. In northern latitudes and in the tropics, LPJmL4 is within the range of the FAPAR datasets. However, LPJmL4 overestimates peak FAPAR, especially in middle and low latitudes, which originates from an overestimation of FAPAR in semi-arid regions. LPJmL4 reproduces the temporal dynamic of FAPAR well in most climate regions with very high correlations between simulated and observed FAPAR in temperate and boreal climates (climate regions Cf and D*) and with medium to high correlations in semi-arid climate regions (e.g. Am, As, Aw, Bsh, Bsk, Cs in Fig. S73). LPJmL4 and the observational datasets show low correlations in wet tropics (Af) and in winter-dry temperate climates (Cw).

LPJmL4 overestimates albedo in all regions (Fig. S74). The temporal dynamic of snow-free albedo was reproduced well in cold steppes (climate region BSk) and in boreal regions (climate regions D*). The correlation between simulated and observed albedo is poor in tropical semi-arid and temperate climates (e.g. As, Aw, Cs, Cf). This is likely caused by soil-moisture-induced changes in soil and background albedo, which has a great effect on soil reflectance (Lobell and Asner, 2002) outside the vegetation season. Such changes are not considered in LPJmL4.

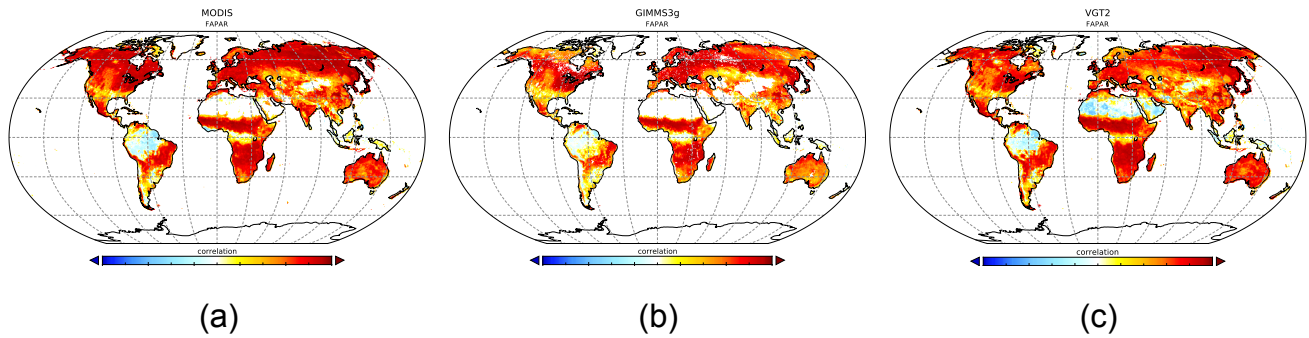


Figure 10. Evaluation of FAPAR for different data sources: MODIS (a), GIMMS (b), and VGT2 (c).

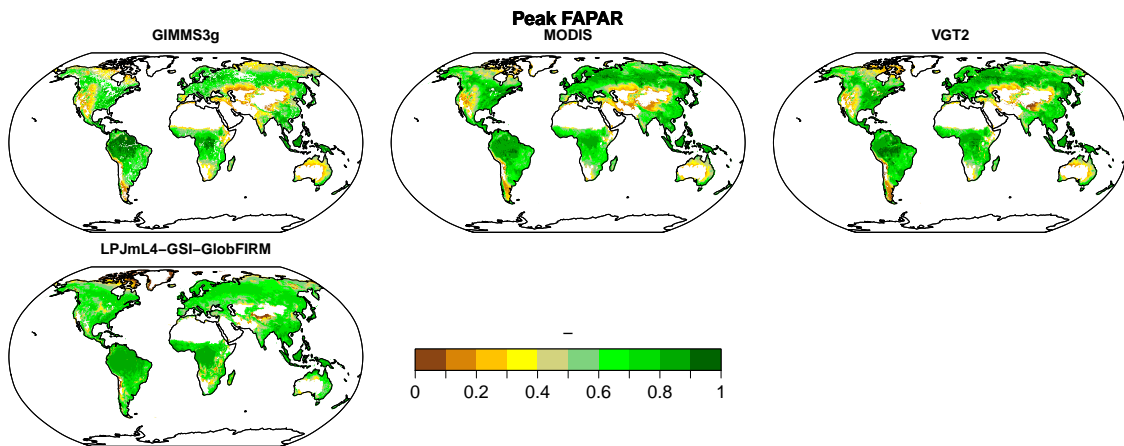


Figure 11. FAPAR mean annual peak comparison with three different remote sensing products.

3.8 Agriculture

3.8.1 Crop yields variability

The evaluation of simulated crop growth and yield can be assessed at individual sites if the model is used as a point model as in different model intercomparison simulations (Asseng et al., 2013, 2015; Bassu et al., 2014; Kollas et al., 2015) in which reference data are available for end-of-season properties (most importantly, crop yield) and within-season dynamics (e.g. development of leaf area index, LAI). The crop yield simulations of LPJmL were evaluated in the framework of the Agricultural Model Intercomparison and Improvement Project (AgMIP) for wheat, maize, rice, and soybean by Müller et al. (2017). They find that the performance of LPJmL is similar to that of the other gridded crop models in that model ensemble ($n = 14$). We supplement the model evaluation with time series correlation analyses for the 10 top-producing countries for all crops implemented in LPJmL4 (Schaphoff et al., 2018c). Results are portrayed in Fig. 12, except for field peas for which no spatial data on crop-specific harvested areas exist for aggregation to national yield time series (Porwollik et al., 2016). As national yield

levels are roughly calibrated in standard LPJmL simulations (Fader et al., 2010), a comparison of the mean bias does not provide insights on model performance. As management intensity is assumed to be static in the simulations (Sect. 2.1), yield trends cannot be reproduced so that simulated and reported national yield time series have been detrended with a running mean approach (Müller et al., 2017) prior to comparison. For a more comprehensive evaluation of LPJmL's performance in yield simulations, see Müller et al. (2017).

The agreement between simulated and observed yields is not only dependent on model performance, but also on the aggregation mask used (Porwollik et al., 2016), assumptions on management and model parameterization (Folberth et al., 2016a), soil parameters (Folberth et al., 2016b), and weather data inputs (Ruane et al., 2016). LPJmL4 yield simulations are typically correlated with national yield statistics (FAO-AQUASTAT, 2014) for some of the 10 top-producing countries for each crop, but only for one country in the case of cassava (Brazil) and sugar cane (China; Figs. 12 and S75–S83 for the other crops).

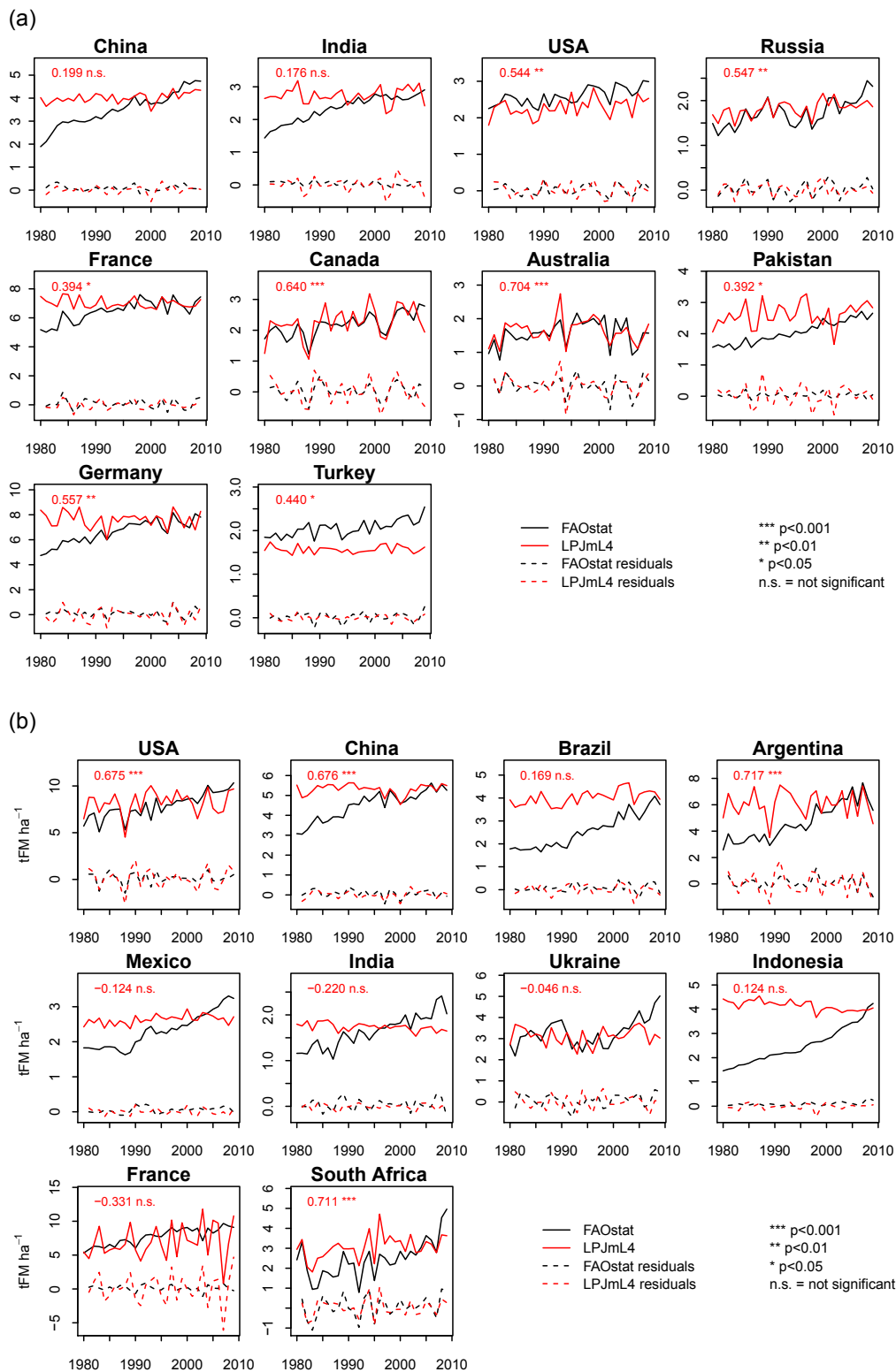


Figure 12. Evaluation of simulated yield variability for wheat (a) and maize (b) in comparison to FAO data (<http://www.fao.org/faostat/en/#data/QC>).

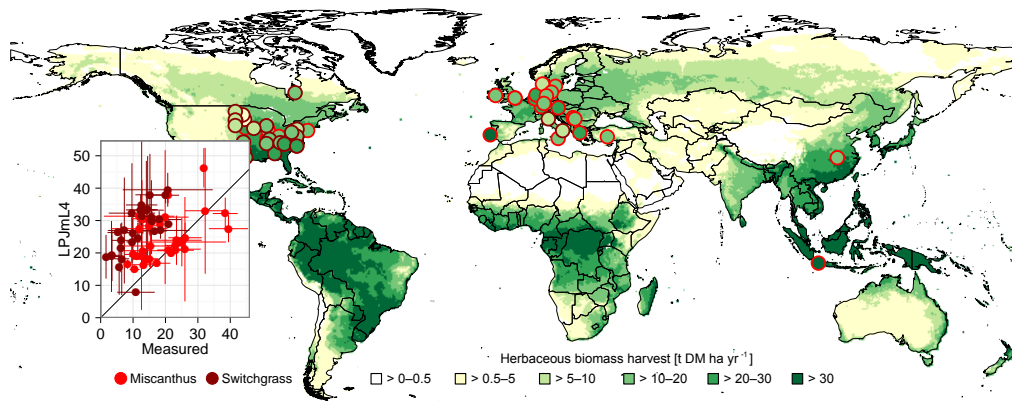
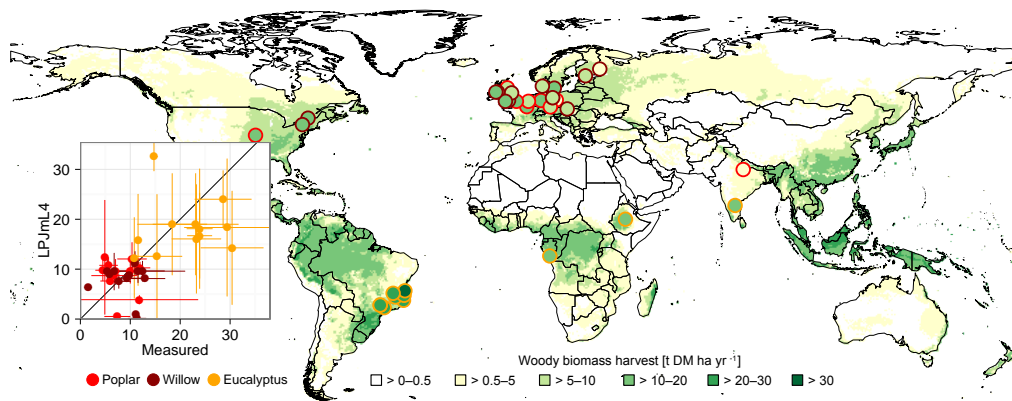
(a) Herbaceous biomass yields [t DM ha⁻¹ a⁻¹](b) Woody biomass yields [t DM ha⁻¹ a⁻¹]

Figure 13. Map of simulated biomass yields by LPJmL4 from rain-fed herbaceous (a) and woody (b) BFTs (averages 1994–2009). Dots indicate the location of the experimental sites and measured yield, with colours scaled to map colours. Scatterplots compare observed and simulated yields in the respective grid cells. Model uncertainty is derived from simulations with and without irrigation. Observation uncertainty reflects dependencies on plantation management (adapted from Heck et al., 2016).

3.8.2 Biomass yield

For the purpose of this evaluation, irrigated and rain-fed biomass plants were simulated to grow globally wherever biophysical conditions allow for sustained growth. The averaged simulated yields for the 16-year period (1994–2009) were compared to reported biomass yields of switchgrass, *Miscanthus*, poplar, willow, and eucalyptus plantations on experimental test sites located in the respective grid cell (Fig. 13). Simulated yields are mostly within the range of observations for *Miscanthus*, poplar, willow, and eucalyptus, but mostly overestimate switchgrass productivity. Management options for BFTs implemented in LPJmL4 are lim-

ited to irrigation management (rain-fed and fully irrigated) because plant species and plantation characteristics (e.g. sapling size and crop spacing) are parameterized as a constant scenario setting and were not varied here. The differences between rain-fed and irrigated biomass yield simulations are depicted as vertical error bars in Fig. 13. The range of rain-fed vs. fully irrigated biomass yields represents an approximation of management uncertainty because simulated yields depend strongly on water availability. Nevertheless, the simulated yield range is likely to represent optimal field management for rain-fed and irrigated plantations as nutrient limitations are not taken into account in these simulations.

Table 5. Indices of agreement between simulated (LPJmL4) and observed (MIRCA2000) sowing dates.

Crop	All cells			Precipitation seasonality			Temperature seasonality		
	W (–)	ME (days)	N	W (–)	ME (days)	N (%)	W (–)	ME (days)	N (%)
Wheat	0.87	44	13 962	0.86	40	15	0.87	44	85
Rice	0.90	25	4995	0.90	24	82	0.87	28	18
Maize	0.88	37	16 333	0.89	37	48	0.85	36	52
Millet	0.89	17	7851	0.92	16	63	0.89	31	37
Pulses	0.63	69	14 712	0.61	80	48	0.84	37	52
Sugar beet	0.37	19	2918	0.24			0.37	19	100
Cassava	0.93	51	6082	0.93	51	83	0.95	57	17
Sunflower	0.92	25	5876	0.87	45	22	0.93	22	78
Soybean	0.94	36	8259	0.94	35	31	0.92	36	69
Groundnut	0.77	34	5642	0.71	36	81	0.96	20	19
Rapeseed	0.86	49	5680	0.36	135	13	0.92	37	87
Wheat (excl. Russia)	0.94	30	11511	0.86	40	18	0.94	29	82

Mean absolute error (ME) and the Willmott coefficient of agreement (W).

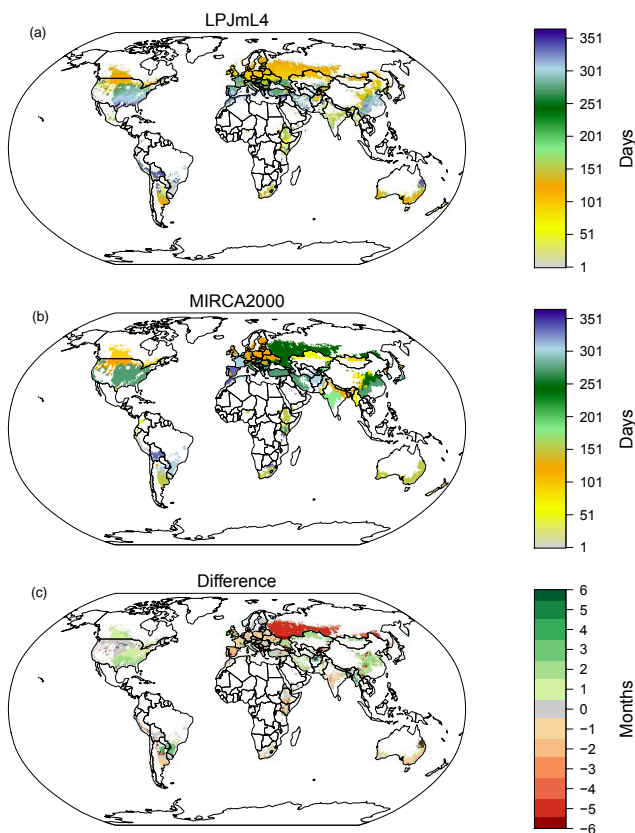


Figure 14. Evaluation of sowing dates for wheat. From (a) to (c): simulated (LPJmL4) sowing date, observed (MIRCA2000) sowing date, and difference between simulated and observed sowing date. Green colours (red colours) in the difference map indicate that simulated sowing dates are too late (too early) compared to observations. White indicates crop area with less than 0.001 % of the grid cell area. Regions without seasonality are not shown.

3.8.3 Month of sowing

The average mean error (ME) for all crops globally is smaller than 2 months, with the exception of pulses (Table 5). For wheat (excl. Russia), millet, rice, sunflower, and sugar beet, the agreement between the simulated and observed timing of sowing is higher, with a difference of about 1 month. The Willmott coefficients (W) are high, indicating good agreement between observations and simulations ($W > 0.85$) for all crops except pulses, sugar beet, and groundnut. Both measures indicate closer agreement for pulses, groundnut, sunflower, and rapeseed in temperate regions (Waha et al., 2012). Poor agreement, with differences between simulated and observed sowing dates of more than 5 months, is found for maize and cassava in South-east Asia and China (for maize in East Africa), for wheat in Russia, for pulses in South-east Asia, India, West and East Africa, the south-east region of Brazil, and southern Australia, for groundnut in India and Indonesia, and for rapeseed in southern Australia and southern Europe (for wheat, Fig. 14; for the other crops, Figs. S84–S93). Divergences are also substantial for crops growing in the southern part of the Democratic Republic of Congo, in South-east Asia, and in tropical climates.

There are several reasons for these disagreements between sowing dates simulated solely using climate data and the global crop calendar; please see Waha et al. (2012) for a more detailed discussion. Firstly the crop varieties in the crop calendar and those simulated here differ, i.e. spring and winter varieties of wheat and rapeseed in temperate regions (e.g. in Russia). Secondly, multiple cropping in tropical regions with high cropping intensity and complex cropping systems is not considered here. Thirdly, we use only one global temperature threshold for simulating sowing temperatures, which is known to vary between regions, and lastly, there are other uncertainties in our method of simulating sowing dates and

in the global crop calendar we use for comparison. We also neglect important factors such as the availability of labour and machinery, social customs, markets and prizes, and the demand for certain agricultural products at certain times in the year.

The comparison to the global crop calendar, however, shows that close agreement between simulated and observed sowing dates can be achieved with purely climate-driven rules for large parts of the Earth for wheat, rice, maize, millet, soybean, and sunflower, as well as for pulses and groundnut in temperate regions. For about 75 % of the global cropping area the difference between simulated and observed sowing dates is 2 months; with the exception of cassava and rapeseed, 80 % of the crop area displays a difference of only 1 month, which is the minimum possible difference as the crop calendar reports monthly sowing dates.

4 Conclusions

This article provides a comprehensive evaluation of the now launched version 4.0 of the LPJmL DGVM that includes an operational representation of agriculture. Unique in its combination of features, the LPJmL4 model enables the simulation of carbon and water fluxes linked to the dynamics of both natural and agricultural vegetation in a single, internally consistent framework. We show that the model has great strength in reproducing carbon fluxes, especially for NBP on the global scale and NEE on the local scale. But we are also able to show that water fluxes match well with other estimates. Both carbon and water fluxes are the link to many ecosystem processes that the model represents and therefore are very important for the understanding of its interrelation. In the agriculture sector we conclude that in regions with a strong weather signal the model is able to match annual yield variability. Nevertheless, in highly managed countries yield variability is not well reproduced by the LPJmL4 model. This can be explained by the absence of a management module in the model. By following suggestions for objective inter-comparative benchmarking systems of multiple models with dedicated software (Abramowitz, 2012; Kelley et al., 2013; Luo et al., 2012), the evaluation takes into account a number of performance metrics, diagnostic plots, and a broad range of fundamental model features. This work thus goes well beyond earlier evaluations of DGVMs (see Kelley et al., 2013) and of model evaluations published for earlier versions LPJmL or its modules.

Pending major model improvements – anticipated as part of forthcoming LPJmL versions – are the incorporation of a scheme for calculating groundwater recharge and storage, the representation of nitrogen cycling for both natural and agricultural landscapes, consideration of ozone effects on plants (Schauberger et al., 2018) and of soil degradation, representation of wetlands with associated methane emissions, the continuous refinement of crop parameterization includ-

ing multi-cropping and other management forms, and possibly a revised implementation of soil moisture (following e.g. Evaristo et al., 2015) and stomatal conductance (following e.g. Lin et al., 2015). As such improvements are expected to have significant effects on plant production and carbon and water fluxes, thus influencing the overall model performance, any future LPJmL version will routinely be subjected to the evaluation protocol used here and, if applicable, tested against other standardized inter-model benchmarks (including participation in model inter-comparisons with evaluation of single components such as in Hattermann et al., 2017). Such continued model maintenance and benchmarking shall also keep pace with recent developments in observational and experimental data, ideally supporting the identification of key uncertainties in model performance (see Medlyn et al., 2015; Smith et al., 2016). Besides identifying features for future model improvement, we demonstrate the adequate performance of the LPJmL4 DGVM in terms of the simulation of long-term averages and also the temporal dynamics across biogeochemical, hydrological, and agricultural processes. This unique capacity renders the LPJmL4 model suitable for process-based analyses of biosphere dynamics including assessments of multi-sectoral impacts of climate change or other anthropogenic Earth system interferences.

Code and data availability. The model code of LPJmL4 is publicly available through PIK's gitlab server at <https://gitlab.pik-potsdam.de/lpjml/LPJmL>, and an exact version of the code described here is archived under <https://doi.org/10.5880/pik.2018.002> and should be referenced as Schaphoff et al. (2018b). The output data from the model simulations described here are available at the research data repository <http://dataservices.gfz-potsdam.de/portal/> under <https://doi.org/10.5880/pik.2017.009> and can be referenced as Schaphoff et al. (2018a).

The Supplement related to this article is available online at <https://doi.org/10.5194/gmd-11-1377-2018-supplement>.

Competing interests. The authors declare that they have no conflict of interest.

Acknowledgements. This study was supported by the German Federal Ministry of Education and Research (BMBF) project “PalMod 2.3 Methankreislauf, Teilprojekt 2 Modellierung der Methanemissionen von Feucht- und Permafrostgebieten mit Hilfe von LPJmL” (FKZ 01LP1507C). Matthias Forkel was funded by the TU Wien Wissenschaftspreis 2015 awarded to Wouter Dorigo. This work used eddy covariance data acquired and shared by the FLUXNET community, including these networks: AmeriFlux, AfriFlux, AsiaFlux, CarboAfrica, CarboEuropeIP, CarboItaly, CarboMont, ChinaFlux, Fluxnet-Canada, GreenGrass, ICOS,

KoFlux, LBA, NECC, OzFlux-TERN, TCOS-Siberia, and USCCC. The ERA-Interim reanalysis data are provided by ECMWF and processed by LSCE. The FLUXNET eddy covariance data processing and harmonization were carried out by the European Fluxes Database Cluster, AmeriFlux Management Project, and Fluxdata project of FLUXNET, with the support of CDIAC and ICOS Ecosystem Thematic Center, and the OzFlux, ChinaFlux, and AsiaFlux offices. We thank the coordinators of the Circumpolar Active Layer Monitoring (CALM) programme for providing thaw depth, the National Snow & Ice Data Center for providing the Circum-Arctic Map of Permafrost, and the R-ArcticNET for discharge data. Furthermore, we thank Jena-BGI data for providing GPP, latent heat flux, total biomass, and soil carbon data. We also thank the providers of the evaluation data of FAPAR (GIMMS3g FAPAR, VGT2 FAPAR, MODIS FAPAR) and the remote sensing data of GFED4 and CCI for evaluating fractional burnt area. MODIS C5 albedo time series data product was retrieved from the online data pool courtesy of the NASA Land Processes Distributed Active Archive Center (LP DAAC), USGS/Earth Resources Observation and Science (EROS) Center, Sioux Falls, South Dakota. We thank the Climatic Research Unit for providing global gridded temperature input, the Global Precipitation Climatology Centre for providing precipitation input, and the coordinators of ERA-Interim for providing shortwave downward radiation and net downward longwave radiation. Furthermore, we thank the coordinators of MIRCA2000 for providing land use input. Finally, we thank Kirsten Elger for her great support in archiving data and the LPJmL4 code and two anonymous reviewers for their helpful comments on earlier versions of the paper.

Edited by: Julia Hargreaves

Reviewed by: two anonymous referees

References

- Abramowitz, G.: Towards a benchmark for land surface models, *Geophys. Res. Lett.*, 32, L22702, <https://doi.org/10.1029/2005GL024419>, 2005.
- Abramowitz, G.: Towards a public, standardized, diagnostic benchmarking system for land surface models, *Geosci. Model Dev.*, 5, 819–827, <https://doi.org/10.5194/gmd-5-819-2012>, 2012.
- Alexandratos, N. and Bruinsma, J.: World agriculture towards 2030/2050: the 2012 revision, Tech. Rep. 12, FAO, Rome, FAO, 2012.
- Asseng, S., Brisson, N., Basso, B., Martre, P., Aggarwal, P. K., Angulo, C., Bertuzzi, P., Biernath, C., Challinor, A. J., Doltra, J., Gayler, S., Goldberg, R., Grant, R., Heng, L., Hooker, J., Hunt, L. A., Ingwersen, J., Izaurralde, R. C., Kersebaum, K. C., Müller, C., Kumar, S. N., Nendel, C., Leary, G. O., Olesen, J. E., Osborne, T. M., Palosuo, T., Priesack, E., Ripoche, D., Semenov, M. A., Shcherbak, I., Steduto, P., Stöckle, C., Stratonovitch, P., Streck, T., Supit, I., Tao, F., Travasso, M., Waha, K., Wallach, D., Williams, J. R., and Wolf, J.: Uncertainty in simulating wheat yields under climate change, *Nat. Clim. Change*, 3, 827–832, <https://doi.org/10.1038/NCLIMATE1916>, 2013.
- Asseng, S., Ewert, F., Martre, P., Rotter, R. P., Lobell, D. B., Cammarano, D., Kimball, B. A., Ottman, M. J., Wall, G. W., White, J. W., Reynolds, M. P., Alderman, P. D., Prasad, P. V. V., Aggarwal, P. K., Anothai, J., Basso, B., Biernath, C., Challinor, A. J., De Sanctis, G., Doltra, J., Fereres, E., Garcia-Vila, M., Gayler, S., Hoogenboom, G., Hunt, L. A., Izaurralde, R. C., Jabloun, M., Jones, C. D., Kersebaum, K. C., Koehler, A.-K., Muller, C., Naresh Kumar, S., Nendel, C., O’Leary, G., Olesen, J. E., Palosuo, T., Priesack, E., Eyshi Rezaei, E., Ruane, A. C., Semenov, M. A., Shcherbak, I., Stockle, C., Stratonovitch, P., Streck, T., Supit, I., Tao, F., Thorburn, P. J., Waha, K., Wang, E., Wallach, D., Wolf, J., Zhao, Z., and Zhu, Y.: Rising temperatures reduce global wheat production, *Nat. Clim. Change*, 5, 143–147, <https://doi.org/10.1038/nclimate2470>, 2015.
- Baret, F., Weiss, M., Lacaze, R., Camacho, F., Makhmara, H., Pacholczyk, P., and Smets, B.: GEOV1: LAI and FAPAR essential climate variables and FCOVER global time series capitalizing over existing products, Part1: Principles of development and production, *Remote Sens. Environ.*, 137, 299–309, <https://doi.org/10.1016/j.rse.2012.12.027>, 2013.
- Bassu, S., Brisson, N., Durand, J.-L., Boote, K., Lizaso, J., Jones, J. W., Rosenzweig, C., Ruane, A. C., Adam, M., Baron, C., Basso, B., Biernath, C., Boogaard, H., Conijn, S., Corbeels, M., Deryng, D., De Sanctis, G., Gayler, S., Grassini, P., Hatfield, J., Hoek, S., Izaurralde, C., Jongschaap, R., Kemanian, A. R., Kersebaum, K. C., Kim, S.-H., Kumar, N. S., Makowski, D., Müller, C., Nendel, C., Priesack, E., Pravia, M. V., Sau, F., Shcherbak, I., Tao, F., Teixeira, E., Timlin, D., and Waha, K.: How do various maize crop models vary in their responses to climate change factors?, *Glob. Change Biol.*, 20, 2301–2320, <https://doi.org/10.1111/gcb.12520>, 2014.
- Becker, A., Finger, P., Meyer-Christoffer, A., Rudolf, B., Schamm, K., Schneider, U., and Ziese, M.: A description of the global land-surface precipitation data products of the Global Precipitation Climatology Centre with sample applications including centennial (trend) analysis from 1901–present, *Earth Syst. Sci. Data*, 5, 71–99, <https://doi.org/10.5194/essd-5-71-2013>, 2013.
- 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., Woodward, F. I., and Papale, D.: Terrestrial gross carbon dioxide uptake: global distribution and covariation with climate, *Science*, 329, 834–8, <https://doi.org/10.1126/science.1184984>, 2010.
- Biemans, H., Hutjes, R. W. a., Kabat, P., Strengers, B. J., Gerten, D., and Rost, S.: Effects of Precipitation Uncertainty on Discharge Calculations for Main River Basins, *J. Hydrometeor.*, 10, 1011–1025, <https://doi.org/10.1175/2008JHM1067.1>, 2009.
- Biemans, H., Haddeland, I., Kabat, P., Ludwig, F., Hutjes, R. W. A., Heinke, J., von Bloh, W., and Gerten, D.: Impact of reservoirs on river discharge and irrigation water supply during the 20th century, *Water Resour. Res.*, 47, W03509, <https://doi.org/10.1029/2009WR008929>, 2011.
- Boden, T., Marland, G., and Andres, R.: Global, Regional, and National Fossil-Fuel CO₂ Emissions, Carbon Dioxide Information Analysis Center (CDIAC), Oak Ridge National Laboratory, US Department of Energy, Oak Ridge, available at: <http://cdiac.ornl.gov/trends/emis/overview.html>, 2013.
- Bondeau, A., Smith, P., Zaehle, S., Schaphoff, S., Lucht, W., Cramer, W., Gerten, D., Lotze-Campen, Hermann, Müller, C., Reichstein, M., and Smith, B.: Modelling the role of agri-

- culture for the 20th century global terrestrial carbon balance, *Glob. Change Biol.*, 13, 679–706, <https://doi.org/10.1111/j.1365-2486.2006.01305.x>, 2007.
- Brouwer, C., Prins, K., and Heibloem, M.: Irrigation Water Management: Irrigation Scheduling, Training manual no. 4, Tech. Rep. 4, FAO Land and Water Development Division, Rome, Italy, available at: <http://www.fao.org/docrep/t7202e/t7202e00.htm>, 1989.
- Brown, J., Ferrians, O. J. J., Heginbottom, J. A., and Melnikov, E. S.: Circum-Arctic map of permafrost and ground-ice conditions, Boulder, CO: National Snow and Ice Data Center/World Data Center for Glaciology, available at: <http://nsidc.org/data/ggd318.html>, 1998.
- Brown, J., Hinkel, K. M., and Nelson, F. E.: The circumpolar active layer monitoring (calm) program: Research designs and initial results, *Polar Geogr.*, 24, 166–258, <https://doi.org/10.1080/10889370009377698>, 2000.
- Carvalhais, N., Forkel, M., Khomik, M., Bellarby, J., Jung, M., Migliavacca, M., Mu, M., Saatchi, S., Santoro, M., Thurner, M., Weber, U., Ahrens, B., Beer, C., Cescatti, A., Randerson, J. T., and Reichstein, M.: Global covariation of carbon turnover times with climate in terrestrial ecosystems, *Nature*, 514, 213–217, <https://doi.org/10.1038/nature13731>, 2014.
- Chaturvedi, V., Hejazi, M., Edmonds, J., Clarke, L., Kyle, P., Davies, E., and Wise, M.: Climate mitigation policy implications for global irrigation water demand, *Mitig. Adapt. Strat. Gl.*, 20, 389–407, <https://doi.org/10.1007/s11027-013-9497-4>, 2015.
- Chuvieco, E., Yue, C., Heil, A., Mouillot, F., Alonso-Canas, I., Padilla, M., Pereira, J. M., Oom, D., and Tansey, K.: A new global burned area product for climate assessment of fire impacts, *Global Ecol. Biogeogr.*, 25, 619–629, <https://doi.org/10.1111/geb.12440>, 2016.
- Cosby, B. J., Hornberger, G. M., Clapp, R. B., and Ginn, T. R.: A Statistical Exploration of the Relationships of Soil Moisture Characteristics to the Physical Properties of Soils, *Water Resour. Res.*, 20, 682–690, <https://doi.org/10.1029/WR020i006p00682>, 1984.
- Dee, D. P., Uppala, S. M., Simmons, A. J., Berrisford, P., Poli, P., Kobayashi, S., Andrae, U., Balmaseda, M. A., Balsamo, G., Bauer, P., Bechtold, P., Beljaars, A. C. M., van de Berg, I., Biblot, J., Bormann, N., Delsol, C., Dragani, R., Fuentes, M., Greer, A. J., Haimberger, L., Healy, S. B., Hersbach, H., Holm, E. V., Isaksen, L., Kallberg, P., Kohler, M., Matricardi, M., McNally, A. P., Mong-Sanz, B. M., Morcette, J.-J., Park, B.-K., Peubey, C., de Rosnay, P., Tavolato, C., Thepaut, J. N., and Vitart, F.: The ERA-Interim reanalysis: Configuration and performance of the data assimilation system, *Q. J. Roy. Meteorol. Soc.*, 137, 553–597, <https://doi.org/10.1002/qj.828>, 2011.
- Defries, R. and Hansen, M.: ISLSCP II Continuous Fields of Vegetation Cover, 1992–1993, ORNL Distributed Active Archive Center, <https://doi.org/10.3334/ORNLDAAC/931>, 2009.
- Döll, P. and Siebert, S.: Global modeling of irrigation water requirements, *Water Resour. Res.*, 38, 1037, <https://doi.org/10.1029/2001WR000355>, 2002.
- Döll, P., Hoffmann-Dobrev, H., Portmann, F., Siebert, S., Eicker, A., Rodell, M., Strassberg, G., and Scanlon, B.: Impact of water withdrawals from groundwater and surface water on continental water storage variations, *J. Geodyn.*, 59–60, 143–156, <https://doi.org/10.1016/j.jog.2011.05.001>, 2012.
- Döll, P., Müller Schmied, H., Schuh, C., Portmann, F. T., and Eicker, A.: Global-scale assessment of groundwater depletion and related groundwater abstractions: Combining hydrological modeling with information from well observations and GRACE satellites, *Water Resour. Res.*, 50, 5698–5720, <https://doi.org/10.1002/2014WR015595>, 2014.
- Evaristo, J., Jasechko, S., and McDonnell, J. J.: Global separation of plant transpiration from groundwater and streamflow, *Nature*, 525, 91–94, <https://doi.org/10.1038/nature14983>, 2015.
- Fader, M., Rost, S., Müller, C., Bondeau, A., and Gerten, D.: Virtual water content of temperate cereals and maize: Present and potential future patterns, *J. Hydrol.*, 384, 218–231, <https://doi.org/10.1016/j.jhydrol.2009.12.011>, 2010.
- FAO-AQUASTAT: AQUASTAT database – Food and Agriculture Organization of the United Nations (FAO), available at: <http://www.fao.org/nr/water/aquastat/data/query/index.html?lang=en>, 2014.
- FAO/IIASA/ISRIC/ISSCAS/JRC: Harmonized World Soil Database (version 1.2), available at: <http://www.iiasa.ac.at/Research/LUC/External-World-soil-database/HTML/>, 2012.
- Folberth, C., Elliott, J., Müller, C., Balkovic, J., Chrystanthacopoulos, J., Izaurralde, R. C., Jones, C. D., Khabarov, N., Liu, W., Reddy, A., Schmid, E., Skalský, R., Yang, H., Arneth, A., Ciais, P., Deryng, D., Lawrence, P. J., Olin, S., Pugh, T. A. M., Ruane, A. C., and Wang, X.: Uncertainties in global crop model frameworks: effects of cultivar distribution, crop management and soil handling on crop yield estimates, *Biogeosciences Discuss.*, <https://doi.org/10.5194/bg-2016-527>, 2016a.
- Folberth, C., Skalský, R., Moltchanova, E., Balkovič, J., Azevedo, L. B., Obersteiner, M., and van der Velde, M.: Uncertainty in soil data can outweigh climate impact signals in global crop yield simulations, *Nat. Comm.*, 7, 11872, <https://doi.org/10.1038/ncomms11872>, 2016b.
- Forkel, M., Carvalhais, N., Schaphoff, S., v. Bloh, W., Migliavacca, M., Thurner, M., and Thonicke, K.: Identifying environmental controls on vegetation greenness phenology through model–data integration, *Biogeosciences*, 11, 7025–7050, <https://doi.org/10.5194/bg-11-7025-2014>, 2014.
- Forkel, M., Migliavacca, M., Thonicke, K., Reichstein, M., Schaphoff, S., Weber, U., and Carvalhais, N.: Codominant water control on global interannual variability and trends in land surface phenology and greenness, *Glob. Change Biol.*, 21, 3414–3435, <https://doi.org/10.1111/gcb.12950>, 2015.
- Forkel, M., Carvalhais, N., Rödenbeck, C., Keeling, R., Heimann, M., Thonicke, K., Zaehle, S., and Reichstein, M.: Enhanced seasonal CO₂ exchange caused by amplified plant productivity in northern ecosystems, *Science*, 351, 696–699, <https://doi.org/10.1126/science.aac4971>, 2016.
- Gerten, D., Schaphoff, S., Haberlandt, U., Lucht, W., and Sitch, S.: Terrestrial vegetation and water balance–hydrological evaluation of a dynamic global vegetation model, *J. Hydrol.*, 286, 249–270, <https://doi.org/10.1016/j.jhydrol.2003.09.029>, 2004.
- Gerten, D., Lucht, W., Ostberg, S., Heinke, J., Kowarsch, M., Kreft, H., Kundzewicz, Z. W., Rastgooy, J., Warren, R., and Schellnhuber, H. J.: Asynchronous exposure to global warming: freshwater resources and terrestrial ecosystems, *Environ. Res. Lett.*, 8, 034032, <https://doi.org/10.1088/1748-9326/8/3/034032>, 2013.
- Giglio, L., Randerson, J. T., and van der Werf, G. R.: Analysis of daily, monthly, and annual burned area using the fourth-

- generation global fire emissions database (GFED4), *J. Geophys. Res.-Biogeo.*, 118, 317–328, <https://doi.org/10.1002/jgrg.20042>, 2013.
- Haddeland, I., Clark, D. B., Franssen, W., Ludwig, F., Voß, F., Arnell, N. W., Bertrand, N., Best, M., Folwell, S., Gerten, D., Gomes, S., Gosling, S. N., Hagemann, S., Hanasaki, N., Harding, R., Heinke, J., Kabat, P., Koirala, S., Oki, T., Polcher, J., Stacke, T., Viterbo, P., Weedon, G. P., and Yeh, P.: Multimodel Estimate of the Global Terrestrial Water Balance: Setup and First Results, *J. Hydrometeor.*, 12, 869–884, <https://doi.org/10.1175/2011JHM1324.1>, 2011.
- Harris, I., Jones, P., Osborn, T., and Lister, D.: Updated high-resolution grids of monthly climatic observations – the CRU TS3.10 Dataset, *Int. J. Climatol.*, 34, 623–642, <https://doi.org/10.1002/joc.3711>, 2014.
- Harrison, P. A., Dunford, R. W., Holman, I. P., and Rounsevell, M. D. A.: Climate change impact modelling needs to include cross-sectoral interactions, *Nat. Clim. Change*, 6, 885–890, <https://doi.org/10.1038/nclimate3039>, 2016.
- Hattermann, F. F., Krysanova, V., Gosling, S. N., Dankers, R., Daggupati, P., Donnelly, C., Flörke, M., Huang, S., Motovilov, Y., Buda, S., Yang, T., Müller, C., Leng, G., Tang, Q., Portmann, F. T., Hagemann, S., Gerten, D., Wada, Y., Masaki, Y., Alemayehu, T., Satoh, Y., and Samaniego, L.: Cross-scale inter-comparison of climate change impacts simulated by regional and global hydrological models in eleven large river basins, *Clim. Change*, 141, 561–576, <https://doi.org/10.1007/s10584-016-1829-4>, 2017.
- Heck, V., Gerten, D., Lucht, W., and Boysen, L. R.: Is extensive terrestrial carbon dioxide removal a “green” form of geoengineering? A global modelling study, *Global Planet. Change*, 137, 123–130, <https://doi.org/10.1016/j.gloplacha.2015.12.008>, 2016.
- Hoff, H., Falkenmark, M., Gerten, D., Gordon, L., Karlberg, L., and Rockström, J.: Greening the global water system, *J. Hydrol.*, 384, 177–186, <https://doi.org/10.1016/j.jhydrol.2009.06.026>, 2010.
- Iizumi, T., Yokozawa, M., Sakurai, G., Travasso, M. I., Romanenkov, V., Oettli, P., Newby, T., Ishigooka, Y., and Furuya, J.: Historical changes in global yields: major cereal and legume crops from 1982 to 2006, *Global Ecol. Biogeogr.*, 23, 346–357, <https://doi.org/10.1111/geb.12120>, 2014.
- Jägermeyr, J., Gerten, D., Lucht, W., Hostert, P., Migliavacca, M., and Nemani, R.: A high-resolution approach to estimating ecosystem respiration at continental scales using operational satellite data, *Glob. Change Biol.*, 20, 1191–1210, <https://doi.org/10.1111/gcb.12443>, 2014.
- Jägermeyr, J., Pastor, A., Biemans, h., and Gerten, D.: Reconciling irrigated food production with environmental flows for Sustainable Development Goals implementation, *Nat. Comm.*, 8, <https://doi.org/10.1038/ncomms15900>, 2017.
- Jung, M., Reichstein, M., Margolis, H. A., Cescatti, A., Richardson, A. D., Arain, M. A., Arneth, A., Bernhofer, C., Bonal, D., Chen, J., Gianelle, D., Gobron, N., Kiely, G., Kutsch, W., Lasslop, G., Law, B. E., Lindroth, A., Merbold, L., Montagnani, L., Moors, E. J., Papale, D., Sottocornola, M., Vaccari, F., and Williams, C.: Global patterns of land-atmosphere fluxes of carbon dioxide, latent heat, and sensible heat derived from eddy covariance, satellite, and meteorological observations, *J. Geophys. Res.-Biogeo.*, 116, G00J07, <https://doi.org/10.1029/2010JG001566>, 2011.
- Kalnay, E., Kanamitsu, M., Kistler, R., Collins, W., Deaven, D., Gandin, L., Iredell, M., Saha, S., White, G., and Woollen, J.: The NCEP/NCAR 40-year reanalysis project, *B. Am. Meteorol. Soc.*, 77, 437–471, [https://doi.org/10.1175/1520-0477\(1996\)077<0437:TNYRP>2.0.CO;2](https://doi.org/10.1175/1520-0477(1996)077<0437:TNYRP>2.0.CO;2), 1996a.
- Kalnay, E., Kanamitsu, M., Kistler, R., Collins, W., Deaven, D., Gandin, L., Iredell, M., Saha, S., White, G., Woollen, J., Zhu, Y., Leetmaa, A., Reynolds, R., Chelliah, M., Ebisuzaki, W., Higgins, W., Janowiak, J., Mo, K. C., Ropelewski, C., Wang, J., Jenne, R., and Joseph, D.: The NCEP/NCAR 40-Year Reanalysis Project, *B. Am. Meteorol. Soc.*, 77, 437–471, [https://doi.org/10.1175/1520-0477\(1996\)077<0437:TNYRP>2.0.CO;2](https://doi.org/10.1175/1520-0477(1996)077<0437:TNYRP>2.0.CO;2), 1996b.
- Kaminski, T., Heimann, M., and Giering, R.: A coarse grid three-dimensional global inverse model of the atmospheric transport: 2. Inversion of the transport of CO₂ in the 1980s, *J. Geophys. Res.-Atmos.*, 104, 18555–18581, <https://doi.org/10.1029/1999JD900146>, 1999.
- Kelley, D. I., Prentice, I. C., Harrison, S. P., Wang, H., Simard, M., Fisher, J. B., and Willis, K. O.: A comprehensive benchmarking system for evaluating global vegetation models, *Biogeosciences*, 10, 3313–3340, <https://doi.org/10.5194/bg-10-3313-2013>, 2013.
- Knyazikhin, Y., Glassy, J., Privette, J. L., Tian, Y., Lotsch, A., Zhang, Y., Wang, Y., Morisette, J. T., Votava, P., Myneni, R. B., and others: MODIS leaf area index (LAI) and fraction of photosynthetically active radiation absorbed by vegetation (FPAR) product (MOD15) algorithm theoretical basis document, Theoretical Basis Document, NASA Goddard Space Flight Center, Greenbelt, MD, 20771, 1999.
- Kollas, C., Kersebaum, K. C., Nendel, C., Manevski, K., Müller, C., Palosuo, T., Armas-Herrera, C. M., Beaudoin, N., Bindi, M., Charfeddine, M., Conradt, T., Constantin, J., Eitzinger, J., Ewert, F., Ferrise, R., Gaiser, T., Cortazar-Atauri, I. G. d., Giglio, L., Hlavinka, P., Hoffmann, H., Hoffmann, M. P., Lounay, M., Manderscheid, R., Mary, B., Mirschel, W., Moriondo, M., Olesen, J. E., Öztürk, I., Pacholski, A., Ripoche-Wachter, D., Roggero, P. P., Roncossek, S., Rötter, R. P., Ruget, F., Sharif, B., Trnka, M., Ventrella, D., Waha, K., Wegehenkel, M., Weigel, H.-J., and Wu, L.: Crop rotation modelling—A European model intercomparison, *Eur. J. Agron.*, 70, 98–111, <https://doi.org/10.1016/j.eja.2015.06.007>, 2015.
- Langerwisch, F., Rost, S., Gerten, D., Poulter, B., Rammig, A., and Cramer, W.: Potential effects of climate change on inundation patterns in the Amazon Basin, *Hydrol. Earth Syst. Sci.*, 17, 2247–2262, <https://doi.org/10.5194/hess-17-2247-2013>, 2013.
- Le Quéré, C., Moriarty, R., Andrew, R. M., Canadell, J. G., Sitch, S., Korsbakken, J. I., Friedlingstein, P., Peters, G. P., Andres, R. J., Boden, T. A., Houghton, R. A., House, J. I., Keeling, R. F., Tans, P., Arneeth, A., Bakker, D. C. E., Barbero, L., Bopp, L., Chang, J., Chevallier, F., Chini, L. P., Ciais, P., Fader, M., Feely, R. A., Gkritzalis, T., Harris, I., Hauck, J., Ilyina, T., Jain, A. K., Kato, E., Kitidis, V., Klein Goldewijk, K., Koven, C., Landschützer, P., Lauvset, S. K., Lefèvre, N., Lenton, A., Lima, I. D., Metzl, N., Millero, F., Munro, D. R., Murata, A., Nabel, J. E. M. S., Nakaoka, S., Nojiri, Y., O’Brien, K., Olsen, A., Ono, T., Pérez, F. F., Pfeil, B., Pierrot, D., Poulter, B., Rehder, G., Rödenbeck, C., Saito, S., Schuster, U., Schwinger, J., Séférian, R., Steinhoff, T., Stocker, B. D., Sutton, A. J., Takahashi, T., Tilbrook, B., van der Laan-Luijkx, I. T., van der Werf, G. R., van Heuven, S., Van-

- demark, D., Viovy, N., Wiltshire, A., Zaehle, S., and Zeng, N.: Global Carbon Budget 2015, *Earth Syst. Sci. Data*, 7, 349–396, <https://doi.org/10.5194/essd-7-349-2015>, 2015.
- Lehner, B., Liermann, C. R., Revenga, C., Vörösmarty, C., Fekete, B., Crouzet, P., Döll, P., Endejan, M., Frenken, K., and Magome, J.: High-resolution mapping of the world's reservoirs and dams for sustainable river-flow management, *Front. Ecol. Environ.*, 9, 494–502, <https://doi.org/10.1890/100125>, 2011.
- Lin, Y., Medlyn, B. E., Duursma, R. E., Prentice, I. C., Wang, H., Baig, S., Eamus, D., Resco de Dios, V., Mitchell, P., Ellsworth, D. S., Op de Beeck, M., Wallin, G., Uddling, J., Tarvainen, L., Linderson, M., Cernusak, L. A., Nippert, J. B., Ocheltree, T. W., Tissue, D. T. and Martin-StPaul, N. K., Rogers, A., Warren, J. M., De Angelis, P., Hikosaka, K., Han, Q., Onoda, Y., Gimeno, T. E., Barton, C. V. M. and Bennie, J., Bonal, J. and Bosc, A., Löw, M., Macinins-Ng, C., Rey, A., Rowland, L., Setterfield, S. A., Tausz-Posch, S., Zaragoza-Castells, J. and Broadmeadow, M. S. J., Drake, J. E., Freeman, M., Ghannoum, O., Hutley, L. B., Kelly, J. W., Kikuzawa, K., Kolari, P., Koyama, K., Limousin, J., Meir, P., Lola da Costa, A. C., Mikkelsen, T. N., Salinas, N., Sun, W., and Wingate, L.: Optimal stomatal behaviour around the world, *Nature Clim. Change*, 5, 459–464, <https://doi.org/10.1038/nclimate2550>, 2015.
- Liu, Y. Y., van Dijk, A. I. J. M., de Jeu, R. A. M., Canadell, J. G., McCabe, M. F., Evans, J. P., and Wang, G.: Recent reversal in loss of global terrestrial biomass, *Nat. Clim. Change*, 5, 470–474, <https://doi.org/10.1038/nclimate2581>, 2015.
- Lobell, D. B. and Asner, G. P.: Moisture Effects on Soil Reflectance, *Soil. Sci. Soc. Am. J.*, 66, 722–727, <https://doi.org/10.2136/sssaj2002.7220>, 2002.
- Lucht, W., Schaaf, C., and Strahler, A.: An algorithm for the retrieval of albedo from space using semiempirical BRDF models, *IEEE T. Geosci. Remote*, 38, 977–998, <https://doi.org/10.1109/36.841980>, 2000.
- Luo, Y. Q., Randerson, J. T., Abramowitz, G., Bacour, C., Blyth, E., Carvalhais, N., Ciais, P., Dalmonech, D., Fisher, J. B., Fisher, R., Friedlingstein, P., Hibbard, K., Hoffman, F., Huntzinger, D., Jones, C. D., Koven, C., Lawrence, D., Li, D. J., Mahecha, M., Niu, S. L., Norby, R., Piao, S. L., Qi, X., Peylin, P., Prentice, I. C., Riley, W., Reichstein, M., Schwalm, C., Wang, Y. P., Xia, J. Y., Zaehle, S., and Zhou, X. H.: A framework for benchmarking land models, *Biogeosciences*, 9, 3857–3874, <https://doi.org/10.5194/bg-9-3857-2012>, 2012.
- Luyssaert, S., Inglima, I., Jung, M., Richardson, A. D., Reichstein, M., Papale, D., Piao, S. L., Schulze, E. D., Wingate, L., Matteucci, G., Aragao, L., Aubinet, M., Beer, C., Bernhofer, C., Black, K. G., Bonal, D., Bonnefond, J. M., Chambers, J., Ciais, P., Cook, B., Davis, K. J., Dolman, A. J., Gielen, B., Goulden, M., Grace, J., Granier, A., Grelle, A., Griffis, T., Grünwald, T., Guidolotti, G., Hanson, P. J., Harding, R., Hollinger, D. Y., Hutyrá, L. R., Kolari, P., Kruijt, B., Kutsch, W., Lagergren, F., Laurila, T., Law, B. E., Le Maire, G., Lindroth, A., Loustau, D., Malhi, Y., Mateus, J., Migliavacca, M., Misson, L., Montagnani, L., Moncrieff, J., Moors, E., Munger, J. W., Nikinmaa, E., Ollinger, S. V., Pita, G., Rebmann, C., Rouspard, O., Saigusa, N., Sanz, M. J., Seufert, G., Sierra, C., Smith, M. L., Tang, J., Valentini, R., Vesala, T., and Janssens, I. A.: CO₂ balance of boreal, temperate, and tropical forests derived from a global database, *Glob. Change Biol.*, 13, 2509–2537, <https://doi.org/10.1111/j.1365-2486.2007.01439.x>, 2007.
- Medlyn, B. E., Zaehle, S., De Kauwe, M. G., Walker, A. P., Dietze, M. C., Hanson, P. J., Hickler, T., Jain, A. K., Luo, Y., Parton, W., Prentice, I. C., Thornton, P. E., Wang, S., Wang, Y.-P., Weng, E., Iversen, C. M., McCarthy, H. R., Warren, J. M., Oren, R., and Norby, R. J.: Using ecosystem experiments to improve vegetation models, *Nat. Clim. Change*, 5, 528–534, <https://doi.org/10.1038/nclimate2621>, 2015.
- Müller, C., Stehfest, E., Minnen, J. G. v., Strengers, B., Bloh, W. v., Beusen, A. H. W., Schaphoff, S., Kram, T., and Lucht, W.: Drivers and patterns of land biosphere carbon balance reversal, *Environ. Res. Lett.*, 11, 044002, <https://doi.org/10.1088/1748-9326/11/4/044002>, 2016.
- Müller, C., Elliott, J., Chryssanthacopoulos, J., Arneith, A., Balkovic, J., Ciais, P., Deryng, D., Folberth, C., Glotter, M., Hoek, S., Iizumi, T., Izaurrealde, R. C., Jones, C., Khabarov, N., Lawrence, P., Liu, W., Olin, S., Pugh, T. A. M., Ray, D. K., Reddy, A., Rosenzweig, C., Ruane, A. C., Sakurai, G., Schmid, E., Skalsky, R., Song, C. X., Wang, X., de Wit, A., and Yang, H.: Global gridded crop model evaluation: benchmarking, skills, deficiencies and implications, *Geosci. Model Dev.*, 10, 1403–1422, <https://doi.org/10.5194/gmd-10-1403-2017>, 2017.
- Nachtergaele, F., van Velthuisen, H., Verelst, L., Batjes, N., Dijkshoorn, K., van Engelen, V., Fischer, G., Jones, A., Montanarella, L., and Petri, M.: Harmonized world soil database, Food and Agriculture Organization of the United Nations, available at: <http://www.fao.org/soils-portal/soil-survey/soil-maps-and-databases/harmonized-world-soil-database-v12/en/>, 2008.
- New, M., Hulme, M., and Jones, P.: Representing Twentieth-Century Space–Time Climate Variability, Part II: Development of 1901–96 Monthly Grids of Terrestrial Surface Climate, *J. Climate*, 13, 2217–2238, [https://doi.org/10.1175/1520-0442\(2000\)013<2217:RTCSTC>2.0.CO;2](https://doi.org/10.1175/1520-0442(2000)013<2217:RTCSTC>2.0.CO;2), 2000.
- ORNL DAAC, Oak Ridge, T. U.: Oak Ridge National Laboratory Distributed Active Archive Center (ORNL DAAC), available at: <http://fluxnet.ornl.gov/>, 2011.
- Ostberg, S., Schaphoff, S., Lucht, W., and Gerten, D.: Three centuries of dual pressure from land use and climate change on the biosphere, *Environ. Res. Lett.*, 10, 44011, <https://doi.org/10.1088/1748-9326/10/4/044011>, 2015.
- Portmann, F. T., Siebert, S., Bauer, C., and Döll, P.: Global dataset of monthly growing areas of 26 irrigated crops, Frankfurt Hydrology Paper, 2008.
- Portmann, F. T., Siebert, S., and Döll, P.: MIRCA2000 – Global monthly irrigated and rainfed crop areas around the year 2000: A new high-resolution data set for agricultural and hydrological modeling, *Global Biogeochem. Cy.*, 24, 1–24, <https://doi.org/10.1029/2008GB003435>, 2010.
- Porwollik, V., Müller, C., Elliott, J., Chryssanthacopoulos, J., Iizumi, T., Ray, D. K., Ruane, A. C., Arneith, A., Balković, J., Ciais, P., Deryng, D., Folberth, C., Izaurrealde, R. C., Jones, C. D., Khabarov, N., Lawrence, P. J., Liu, W., Pugh, T. A., Reddy, A., Sakurai, G., Schmid, E., Wang, X., de Wit, A., and Wu, X.: Spatial and temporal uncertainty of crop yield aggregations, *Eur. J. Agron.*, <https://doi.org/10.1016/j.eja.2016.08.006>, 2016.
- Randerson, J., van der Werf, G. R., Giglio, L., Collatz, G. J., and Kasibhatla, P. S.: Global Fire Emis-

- sions Database, Version 4, (GFEDv4), ORNL DAAC, <https://doi.org/10.3334/ORNLDAAC/1293>, 2015.
- Ray, D. K., Gerber, J. S., MacDonald, G. K., and West, P. C.: Climate variation explains a third of global crop yield variability, *Nat. Comm.*, 6, 5989, <https://doi.org/10.1038/ncomms6989>, 2015.
- Rödenbeck, C.: Estimating CO₂ sources and sinks from atmospheric mixing ratio measurements using a global inversion of atmospheric transport, Technical Reports, Max Planck Institute for Biogeochemistry, available at: <http://pubman.mpdl.mpg.de/pubman/faces/viewItemOverviewPage.jsp?itemId=escidoc:1691952>, 2005.
- Rödenbeck, C., Houweling, S., Gloor, M., and Heimann, M.: CO₂ flux history 1982–2001 inferred from atmospheric data using a global inversion of atmospheric transport, *Atmos. Chem. Phys.*, 3, 1919–1964, <https://doi.org/10.5194/acp-3-1919-2003>, 2003.
- Rosegrant, M. W., Cai, X., and Cline, S. A.: World Water and Food to 2025: Dealing with Scarcity, Tech. rep., International Food Policy Research Institute, Washington, DC, 2002.
- Rosenzweig, C., Elliott, J., Deryng, D., Ruane, A. C., Müller, C., Arneth, A., Boote, K. J., Folberth, C., Glotter, M., and Khabarov, N.: Assessing agricultural risks of climate change in the 21st century in a global gridded crop model intercomparison, *P. Natl. Acad. Sci. USA*, 111, 3268–3273, <https://doi.org/10.1073/pnas.1222463110>, 2014.
- Ruane, A. C., Hudson, N. I., Asseng, S., Camarrano, D., Ewert, F., Martre, P., Boote, K. J., Thorburn, P. J., Aggarwal, P. K., and Angulo, C.: Multi-wheat-model ensemble responses to inter-annual climate variability, *Environ. Model. Softw.*, 81, 86–101, <https://doi.org/10.1016/j.envsoft.2016.03.008>, 2016.
- Saatchi, S. S., Harris, N. L., Brown, S., Lefsky, M., Mitchard, E. T. A., Salas, W., Zutta, B. R., Buermann, W., Lewis, S. L., Hagen, S., Petrova, S., White, L., Silman, M., and Morel, A.: Benchmark map of forest carbon stocks in tropical regions across three continents, *P. Natl. Acad. Sci. USA*, 108, 9899–9904, <https://doi.org/10.1073/pnas.1019576108>, 2011.
- Saleska, S. R., Miller, S. D., Matross, D. M., Goulden, M. L., Wofsy, S. C., Da Rocha, H. R., De Camargo, P. B., Crill, P., Daube, B. C., De Freitas, H. C., Hutya, L., Keller, M., Kirchhoff, V., Menton, M., Munger, J. W., Hammond Pyle, E., Rice, A. H., and Silva, H.: Carbon in Amazon Forests: Unexpected Seasonal Fluxes and Disturbance-Induced Losses, *Science*, 302, 1554–1557, <https://doi.org/10.1126/science.1091165>, 2003.
- Sauer, T., Havlík, P., Schneider, U. a., Schmid, E., Kindermann, G., and Obersteiner, M.: Agriculture and resource availability in a changing world: The role of irrigation, *Water Resour. Res.*, 46, W06503, <https://doi.org/10.1029/2009WR007729>, 2010.
- Schaaf, C. B., Gao, F., Strahler, A. H., Lucht, W., Li, X., Tsang, T., Strugnell, N. C., Zhang, X., Jin, Y., Muller, J.-P., Lewis, P., Barnsley, M., Hobson, P., Disney, M., Roberts, G., Dunderdale, M., Doll, C., d'Entremont, R. P., Hu, B., Liang, S., Privette, J. L., and Roy, D.: First operational BRDF, albedo nadir reflectance products from MODIS, The Moderate Resolution Imaging Spectroradiometer (MODIS): a new generation of Land Surface Monitoring, *Remote Sens. Environ.*, 83, 135–148, [https://doi.org/10.1016/S0034-4257\(02\)00091-3](https://doi.org/10.1016/S0034-4257(02)00091-3), 2002.
- Schaphoff, S., Heyder, U., Ostberg, S., Gerten, D., Heinke, J., and Lucht, W.: Contribution of permafrost soils to the global carbon budget, *Environ. Res. Lett.*, 8, 014026, <https://doi.org/10.1088/1748-9326/8/1/014026>, 2013.
- Schaphoff, S., von Bloh, W., Rammig, A., Thonicke, K., Forkel, M., Biemans, H., Gerten, D., Heinke, J., Jägermyer, J., Knauer, J., Lucht, W., Müller, C., Rolinski, S., and Waha, K.: LPJmL4 model output for the publications in GMD: LPJmL4 – a dynamic global vegetation model with managed land: Part I – Model description and Part II – Model evaluation, <https://doi.org/10.5880/pik.2017.009>, 2018a.
- Schaphoff, S., von Bloh, W., Thonicke, K., Biemans, H., Forkel, M., Heinke, J., Jägermyer, J., Müller, C., Rolinski, S., Waha, K., Stehfest, E., de Waal, L., Heyder, U., Gumpenberger, M., and Beringer, T.: LPJmL4 Model Code, V. 4.0. GFZ Data Services, <https://doi.org/10.5880/pik.2018.002>, 2018b.
- Schaphoff, S., von Bloh, W., Rammig, A., Thonicke, K., Biemans, H., Forkel, M., Gerten, D., Heinke, J., Jägermyer, J., Knauer, J., Langerwisch, F., Lucht, W., Müller, C., Rolinski, S., and Waha, K.: LPJmL4 – a dynamic global vegetation model with managed land – Part 1: Model description, *Geosci. Model Dev.*, 11, 1343–1375, <https://doi.org/10.5194/gmd-11-1343-2018>, 2018c.
- Schauberger, B., Rolinski, S., and Müller, C.: A network-based approach for semi-quantitative knowledge mining and its application to yield variability, *Environ. Res. Lett.*, 11, 123001, <https://doi.org/10.1088/1748-9326/11/12/123001>, 2016.
- Schauberger, B., Rolinski, S., Schaphoff, S., and Müller, C.: Global historical soybean and wheat yield loss estimates from ozone pollution considering water and temperature as modifying effects, *Glob. Change Biol.*, under review, 2018.
- Siderius, C., Biemans, H., Wiltshire, A., Rao, S., Franssen, W. H. P., Kumar, P., Gosain, A. K., van Vliet, M. T. H., and Collins, D. N.: Snowmelt contributions to discharge of the Ganges, *Sci. Total Environ.*, 468, S93–S101, <https://doi.org/10.1016/j.scitotenv.2013.05.084>, 2013.
- Siebert, S. and Döll, P.: Quantifying blue and green virtual water contents in global crop production as well as potential production losses without irrigation, *J. Hydrol.*, 384, 198–217, <https://doi.org/10.1016/j.jhydrol.2009.07.031>, 2010.
- Sitch, S., Smith, B., Prentice, I. C., Arneth, A., Bondeau, A., Cramer, W., Kaplan, J. O., Levis, S., Lucht, W., Sykes, M. T., Thonicke, K., and Venevsky, S.: Evaluation of ecosystem dynamics, plant geography and terrestrial carbon cycling in the LPJ dynamic global vegetation model, *Glob. Change Biol.*, 9, 161–185, <https://doi.org/10.1046/j.1365-2486.2003.00569.x>, 2003.
- Smith, W. K., Reed, S. C., Cleveland, C. C., Ballantyne, A. P., Anderegg, W. R. L., Wieder, W. R., Liu, Y. Y., and Running, S. W.: Large divergence of satellite and Earth system model estimates of global terrestrial CO₂ fertilization, *Nat. Clim. Change*, 6, 306–310, <https://doi.org/10.1038/nclimate2879>, 2016.
- Steffen, W., Richardson, K., Rockström, J., Cornell, S. E., Fetzer, I., Bennett, E. M., Biggs, R., Carpenter, S. R., de Vries, W., de Wit, C. A., Folke, C., Gerten, D., Heinke, J., Mace, G. M., Persson, L. M., Ramanathan, V., Reyers, B., and Sörlin, S.: Planetary boundaries: Guiding human development on a changing planet, *Science*, 347, 1259855, <https://doi.org/10.1126/science.1259855>, 2015.
- Tans, P. and Keeling, R.: Trends in Atmospheric Carbon Dioxide, National Oceanic & Atmospheric Administration, Earth System Research Laboratory (NOAA/ESRL), available at: <http://www.esrl.noaa.gov/gmd/ccgg/trends>, 2015.

- Tarnocai, C., Canadell, J. G., Schuur, E. A. G., Kuhry, P., Mazhitova, G., and Zimov, S.: Soil organic carbon pools in the northern circumpolar permafrost region, *Global Biogeochem. Cy.*, 23, GB2023, <https://doi.org/10.1029/2008GB003327>, 2009.
- Taylor, K. E.: Summarizing multiple aspects of model performance in a single diagram, *J. Geophys. Res.-Atmos.*, 106, 7183–7192, <https://doi.org/10.1029/2000JD900719>, 2001.
- Thonicke, K., Venevsky, S., Sitch, S., and Cramer, W.: The role of fire disturbance for global vegetation dynamics: coupling fire into a Dynamic Global Vegetation Model, *Global Ecol. Biogeogr.*, 10, 661–677, <https://doi.org/10.1046/j.1466-822X.2001.00175.x>, 2001.
- Thonicke, K., Spessa, A., Prentice, I. C., Harrison, S. P., Dong, L., and Carmona-Moreno, C.: The influence of vegetation, fire spread and fire behaviour on biomass burning and trace gas emissions: results from a process-based model, *Biogeosciences*, 7, 1991–2011, <https://doi.org/10.5194/bg-7-1991-2010>, 2010.
- Thoning, K., Tans, P., and Komhyr, W.: Atmospheric carbon dioxide at Mauna Loa Observatory. II – Analysis of the NOAA GMCC data, 1974–1985, *J. Geophys. Res.*, 94, 8549–8565, <https://doi.org/10.1029/JD094iD06p08549>, 1989.
- Turner, M., Beer, C., Santoro, M., Carvalhais, N., Wutzler, T., Schepaschenko, D., Shvidenko, A., Kompter, E., Ahrens, B., Levick, S. R., and Schimmlius, C.: Carbon stock and density of northern boreal and temperate forests, *Global Ecol. Biogeogr.*, 23, 297–310, <https://doi.org/10.1111/geb.12125>, 2014.
- University of East Anglia Climatic Research Unit, Harris, I. C., and Jones, P.: CRU TS3.23: Climatic Research Unit (CRU) Time-Series (TS) Version 3.23 of High Resolution Gridded Data of Month-by-month Variation in Climate (Jan. 1901–Dec. 2014), Centre for Environmental Data Analysis, <https://doi.org/10.5285/4c7fdfa6-f176-4c58-acee-683d5e9d2ed5>, 2015.
- Vorosmarty, C. and Fekete, B.: ISLSCP II River Routing Data (STN-30p), in: ISLSCP Initiative II Collection, Data set, edited by: Hall, F. G., Collatz, G., Meeson, B., Los, S., Brown de Colstoun, E., and Landis, D., ORNL Distributed Active Archive Center, <https://doi.org/10.3334/ORNLDAAAC/1005>, 2011.
- Vörösmarty, C. J., Fekete, B., and Tucker, B.: River Discharge Database, Version 1.0 (RivDIS v1.0), Volumes 0 through 6, A contribution to IHP-V Theme 1. Technical Documents in Hydrology Series, UNESCO, Paris, 1996.
- Wada, Y. and Bierkens, M. F. P.: Sustainability of global water use: past reconstruction and future projections, *Environ. Res. Lett.*, 9, 104003, <https://doi.org/10.1088/1748-9326/9/10/104003>, 2014.
- Wada, Y., van Beek, L. P. H., Viviroli, D., Dürr, H. H., Weingartner, R., and Bierkens, M. F. P.: Global monthly water stress: 2. Water demand and severity of water stress, *Water Resour. Res.*, 47, W07518, <https://doi.org/10.1029/2010WR009792>, 2011.
- Waha, K., van Bussel, L. G. J., Müller, C., and Bondeau, A.: Climate-driven simulation of global crop sowing dates, *Global Ecol. Biogeogr.*, 21, 247–259, <https://doi.org/10.1111/j.1466-8238.2011.00678.x>, 2012.
- Waha, K., Müller, C., Bondeau, a., Dietrich, J., Kurukulasuriya, P., Heinke, J., and Lotze-Campen, H.: Adaptation to climate change through the choice of cropping system and sowing date in sub-Saharan Africa, *Global Environ. Chang.*, 23, 130–143, <https://doi.org/10.1016/j.gloenvcha.2012.11.001>, 2013.
- Warszawski, L., Frieler, K., Huber, V., Piontek, F., Serdeczny, O., and Schewe, J.: The Inter-Sectoral Impact Model Intercomparison Project (ISI-MIP): Project framework, *P. Natl. Acad. Sci. USA*, 111, 3228–3232, <https://doi.org/10.1073/pnas.1312330110>, 2014.
- Willmott, C. J.: Some comments on the evaluation of model performance, *B. Am. Meteorol. Soc.*, 63, 1309–1313, [https://doi.org/10.1175/1520-0477\(1982\)063<1309:SCOTEO>2.0.CO;2](https://doi.org/10.1175/1520-0477(1982)063<1309:SCOTEO>2.0.CO;2), 1982.
- Zhu, Z., Bi, J., Pan, Y., Ganguly, S., Anav, A., Xu, L., Samanta, A., Piao, S., Nemani, R. R., and Myneni, R. B.: Global Data Sets of Vegetation Leaf Area Index (LAI)3g and Fraction of Photosynthetically Active Radiation (FPAR)3g Derived from Global Inventory Modeling and Mapping Studies (GIMMS) Normalized Difference Vegetation Index (NDVI)3g for the Period 1981 to 2011, *Remote Sens.*, 5, 927–948, <https://doi.org/10.3390/rs5020927>, 2013.
- Zscheischler, J., Mahecha, M., Von Buttlar, J., Harmeling, S., Jung, M., Rammig, A., Randerson, J. T., Schölkopf, B., Seneviratne, S. I., Tomelleri, E., Zaehle, S., and Reichstein, M.: Few extreme events dominate global interannual variability in gross primary production, *Environ. Res. Lett.*, 9, 035001, <https://doi.org/10.1088/1748-9326/9/3/035001>, 2014a.