1 **A new step-wise Carbon Cycle Data Assimilation System**

2 **using multiple data streams to constrain the simulated land**

3 **surface carbon cycle**

4

5 **P. Peylin[1], C. Bacour[2], N. MacBean[1], S. Leonard[1], P. J. Rayner[1,3], S. Kuppel[1], E.**

6 **N. Koffi[1], A. Kane[1], F. Maignan[1], F. Chevallier[1], P. Ciais[1], P. Prunet[2]**

7 [1]{Laboratoire des Sciences du Climat et de l'Environnement, UMR 8212 CEA-CNRS-

8 UVSQ, 91191 Gif-sur-Yvette cedex, France}

9 [2]{Noveltis, Parc Technologique du Canal, 2 avenue de l'Europe, 31520 Ramonville-Saint-

10 Agne, France}

11 [3]{University of Melbourne, 3010, Vic, Melbourne, Australia}

12 Correspondence to: P. Peylin (philippe.peylin@lsce.ipsl.fr)

13

14

Geoscientific
Model Development
Discussions

Open Access

EGU

# 1 **Abstract**

2 Large uncertainties in Land surface models (LSMs) simulations still arise from inaccurate

3 forcing, incorrect model parameter values and incomplete representation of biogeochemical

4 processes. The recent increase in the number and type of carbon cycle related observations,

5 including both in situ and remote sensing measurements, has opened a new road to optimize

6 model parameters via robust statistical model-data integration techniques, in order to reduce

7 the simulated carbon fluxes and stocks uncertainties. In this study we present a Carbon Cycle

8 Data Assimilation System (CCDAS) that assimilates three major data streams, namely

9 MODIS-NDVI observations of vegetation activity, net ecosystem exchange (NEE) and latent

10 heat (LE) flux measurements at more than 70 sites (FLUXNET), and atmospheric $CO_2$

11 concentrations at 53 surface stations, in order to optimize the main parameters of the

12 ORCHIDEE LSM (around 180 parameters in total). The system relies on a step-wise

13 approach that assimilates each data stream in turn, propagating the information gained on the

14 parameters from one step to the next.

15 Overall, the ORCHIDEE model is able to achieve a consistent fit to all three data streams,

16 which suggests that current LSMs have reached the level of development to assimilate these

17 observations. The assimilation of MODIS-NDVI (step 1) reduced the growing season length

18 in ORCHIDEE for temperate and boreal ecosystems, thus decreasing the global mean annual

19 gross primary production (GPP). Using FLUXNET data (step 2) led to large improvements in

20 the seasonal cycle of the NEE and LE fluxes for all ecosystems (i.e., increased amplitude for

21 temperate ecosystems). The assimilation of atmospheric $CO_2$, using the atmospheric transport

22 model LMDz (step 3), provides an overall constraint (i.e., constraint on large scale net $CO_2$

23 fluxes), resulting in an improvement of the fit to the observed atmospheric $CO_2$ growth rate.

24 Thus the optimized model predicts a land C sink of around 2.2 PgC.yr$^{-1}$ (for the 2000-2009

25 period), which is more compatible with current estimates from the Global Carbon Project

26 (GCP) than the prior value. The consistency of the step-wise approach is evaluated with back-

27 compatibility checks. The final optimized model (after step 3) does not significantly degrade

28 the fit to MODIS-NDVI and FLUXNET data that were assimilated in the first two steps,

29 suggesting that a stepwise approach can be used instead of the more "challenging"

30 implementation of a simultaneous optimization in which all data streams are assimilated

31 together. Most parameters, including the scalar of the initial soil carbon pool size, changed

1    during the optimization with a large error reduction. This work opens new perspectives for

2    better predictions of the land carbon budgets.

3

4    **1   Introduction**

5    Atmospheric $CO_2$ concentrations have increased at an unprecedented rate over the last few

6    decades, predominantly due to anthropogenic fossil fuel and cement emissions, as well as

7    land use and land cover change (LULCC). The oceans and the terrestrial biosphere have

8    absorbed $CO_2$, removing on average 50% of anthropogenic emissions from the atmosphere.

9    However, knowledge about the exact location of sources and sinks of carbon (C) and the

10    driving mechanisms is still lacking. Land surface models (LSMs) can be used to improve our

11    understanding of the spatio-temporal patterns of sources and sinks, as well as for attributing

12    changes due to $CO_2$, climate variability and other environmental drivers. However, the spread

13    in the model predictions of terrestrial net C exchange currently has the same order of

14    magnitude as the uncertainty of the terrestrial C budget estimated as the residual of the other

15    components (Le Quéré et al., 2015). In addition to uncertainties in the mean global annual

16    terrestrial C budget and its trend over time (Sitch et al., 2015), there remains strong

17    discrepancies between LSMs in their predictions of regional budgets (Canadell, 2013) at

18    seasonal and inter-annual timescales and in their sensitivity to climate and atmospheric $CO_2$

19    forcing (Piao et al., 2013).

20    Uncertainties in model simulations arise from inaccurate forcing, incorrect model parameter

21    values and/or an inadequate or incomplete representation of biogeochemical processes in the

22    model (for example the impact of nutrient limitation on C fluxes, or C release related to

23    permafrost thawing). Arguably the best way to improve model predictions is to confront

24    simulations with multiple sources of data within an appropriate and rigorous framework

25    (Prentice et al., 2015). In the last two decades significant efforts by the site and satellite

26    observation communities have resulted in a large increase in the number and type of C cycle-

27    related observations. These data contain some information at various spatial and temporal

28    scales and should be combined together to robustly address different aspects of the models.

29    One way in which these data can be used to better quantify and reduce model uncertainty is

30    by optimizing or calibrating the model parameters via robust statistical model-data fusion (or

31    data assimilation – DA) techniques. In particular a Bayesian inference framework allows us to

1    update our prior knowledge of the parameters based on new information contained in the

2    observations.

3    There is a long history of using DA techniques for parameter optimization, particularly in

4    Geophysics (Tarantola, 1987), but the initial studies in the field of global terrestrial C cycle

5    data assimilation started with the initial study of Fung et al. (1987) and a pioneering work by

6    Knorr and Heimann (1995) who used atmospheric $CO_2$ concentration to constrain the Simple

7    Diagnostic Biosphere Model (SDBM). This effort was continued by the original Carbon

8    Cycle Data Assimilation System (CCDAS) described in Rayner et al. (2005) and Kaminski et

9    al. (2012) which used both atmospheric $CO_2$ and satellite-derived Fraction of Absorbed

10   Photosynthetic Radiation (FAPAR) data to optimize vegetation productivity by adjusting the

11   C cycle-related parameters of the Biosphere Energy-Transfer Hydrology (BETHY) model

12   (see a review in Kaminski et al., 2013). Meanwhile substantial efforts have been put into the

13   use of local eddy covariance flux tower measurements of net exchange of $CO_2$ and latent and

14   sensible heat fluxes to optimize photosynthesis, respiration and energy-related parameters of

15   terrestrial ecosystem models, both at individual sites (e.g. Wang et al., 2001, 2007; Williams

16   et al., 2005; Braswell et al., 2005; Knorr and Kattge, 2005; Moore et al., 2008; Ricciuto et al.,

17   2008), and more recently using multiple sites together (hereafter multiple sites) from the

18   global FLUXNET network (e.g. Groenendijk et al., 2011; Kuppel et al., 2012, 2014; Alton,

19   2013; Xiao et al., 2014). Increasingly the focus in carbon cycle data assimilation is moving

20   towards using multiple different data streams as independent constraints, with the aim of

21   bringing more information at different spatial and temporal scales and constraining several

22   processes at once in order to reduce the likelihood of model equifinality (where multiple sets

23   of parameters achieve the same reduction in model-data misfit). Recent examples include the

24   combination of in-situ eddy covariance flux observations and ground-based information on

25   vegetation structure and C stocks (Richardson et al., 2010; Ricciuto et al., 2011; Keenan et al.,

26   2012, 2013; Thum et al., 2015), or in-situ flux data and satellite FAPAR (Kato et al., 2013;

27   Zobitz et al., 2014; Bacour et al., 2015). This is a non-trivial task however, especially when

28   optimizing a complex LSM (see MacBean et al, submitted), which has many parameters

29   acting from local to global scales.

30   When assimilating multiple different data streams we have two options: i) to optimize the

31   model with each data stream in turn, and to propagate the information gained on the

32   parameter values from one step to the next (hereafter referred to as "stepwise" assimilation),

1   or ii) to include all data streams together in the same optimization (hereafter referred to as

2   "simultaneous" assimilation). Kaminski et al. (2012) suggested that it is essential to perform a

3   consistent, simultaneous assimilation that includes all data streams in the same optimization.

4   It is important to note that this is an implementation question. Tarantola (2005) recasts the

5   fundamentals of the approach as the conjunction or multiplication of probability densities.

6   This multiplication is associative so it makes no difference whether it is performed in one step

7   or several. In complex problems such as these, one cannot carry or even describe the full

8   structure of the relevant probability densities so which approach will work best in each case is

9   unclear. In particular, technical difficulties associated with the different number of

10  observations for each data stream and the characterization of error correlations between them,

11  in addition to computational constraints to run global LSMs, might result in the preference for

12  a step-wise assimilation framework. Additionally, it may be more straightforward, to expose a

13  restricted set of parameters to each observation type in a stepwise approach to ensure that

14  each data stream constrains only the most relevant parts of the model. This reduces biases

15  from other poorly-represented processes caused by inadequate model structure. For these

16  reasons we follow the stepwise approach in this paper.

17  We present the first global-scale CCDAS that assimilates three of the main global data

18  streams that have been used to date to understand the terrestrial carbon cycle – atmospheric

19  $CO_2$ concentration, satellite-derived information of vegetation greenness (from the MODIS

20  instrument) and multisite eddy covariance net $CO_2$ and latent heat flux measurements (from

21  FLUXNET) – to optimize the parameters of the Organizing Carbon and Hydrology in

22  Dynamics Ecosystems (ORCHIDEE) process-based LSM (Krinner et al., 2005). The main

23  questions that we aim to answer in this paper are as follows:

24  i) How and to which extend the optimization of the ORCHIDEE model allows to fit the three

25  data streams that are considered?

26  ii) Does the step-wise optimization result in a degradation of the fit to other data streams used

27  in the previous steps?

28  iii) What are the main changes in the optimized parameters when using sequentially these

29  three data streams in a global CCDAS and which processes are constrained?

30  iv) What are the improvements for the land C cycle in terms of net/gross fluxes and stocks as

31  a result of multi-data stream optimization? What preliminary perspectives can we draw that

1   may help us in improving model predictions of trends, variability and the location of

2   terrestrial C sources and sinks?

3   Following these objectives, the paper first describes the new ORCHIDEE-CCDAS including

4   the concept, the observations, the models and the optimization approach. We then present the

5   results, including the fit to the data, consistency checks (question i) above) as well as mean

6   global and regional C cycle budget for the period 2000-2009. The last section discusses issues

7   and perspectives associated with these results.

8

## 2   Methods

### 2.1   ORCHIDEE-CCDAS concept

11  We have designed a CCDAS around the ORCHIDEE land surface model (ORCHIDEE-

12  CCDAS, later also referred to as ORCHIDAS for simplicity) that combines a state-of-the-art

13  description of the driving biogeochemical processes within the model with multiple

14  observational constraints in a robust statistical framework, in order to improve the simulation

15  of land carbon fluxes and stocks. The system allows us to retrieve the best estimate, given the

16  observations and prior information, of selected parameters (see §2.3.3) as well as to evaluate

17  their uncertainty. It relies on a stepwise assimilation of a comprehensive set of three C cycle-

18  related observations that are representative of small (100 m) to large (continental) scales (see

19  §2.2):

20  •   Step 1: Satellite measurements of vegetation activity using the Normalized Difference

21      Vegetation Index (NDVI) from the MODIS instrument over the 2000-2008 period for

22      a randomly selected set of sites for boreal and temperate deciduous vegetation types;

23  •   Step 2: In-situ eddy-covariance net $CO_2$ and water (latent heat) flux measurements

24      from the FLUXNET database for a large set of sites, spanning 7 different vegetation

25      types;

26  •   Step 3: In-situ monthly atmospheric surface $CO_2$ concentration measurements from

27      the GLOBALVIEW-CO2 database over three years (2002-2004).

28  The system relies on two models:

29  •   The ORCHIDEE global LSM, whose main C cycle parameters are optimized (see

30      §2.3)

1    • The atmospheric transport model, LMDz (see §2.3), to relate the surface carbon fluxes

2        to atmospheric $CO_2$ concentrations.

3    The framework combines the different observational data streams within ORCHIDAS in

4    order to optimize selected model parameters using a variational data assimilation system,

5    described in section 2.4. Figure 1 illustrates the structure of the CCDAS and the different

6    components that are involved. Such a framework distinguishes i) the assimilated observations,

7    ii) an ensemble of forcing and input data streams, iii) the models and optimization framework,

8    as well as iv) an evaluation step, where independent datasets are compared to the optimized

9    model stocks and fluxes. As explained in the introduction, a major feature of the current

10   system is the stepwise approach, in which all data streams are assimilated sequentially (i.e.

11   one after the other). The information retrieved at a given step (retrieved optimal parameter

12   values and associated uncertainty) is propagated to the next step (see Fig. 2 and §2.4). Note

13   that for simplicity we did not propagated the error correlations in this first implementation of

14   the system.

15   At each step, the parameter optimization relies on a Bayesian framework that explicitly

16   minimizes the difference between the simulated and observed quantities in addition to

17   minimizing the difference between the optimized model parameters and "a priori" values (see

18   §2.4.2). The dependence of the simulated quantities on the optimized variables is non-linear,

19   which thus necessitates the use of an iterative algorithm. Note that all components of the

20   surface C budget need also to be included in the ORCHIDAS, particularly when using

21   atmospheric $CO_2$ measurements which requires the atmospheric transport model to be

22   prescribed with fossil fuel emissions, $CO_2$ fluxes associated with biomass burning and ocean

23   $CO_2$ fluxes (see §2.5) in addition to net ecosystem exchange (NEE) from ORCHIDEE.

24   **2.2  Assimilated observations**

25   2.2.1 MODIS-NDVI

26   MODIS collection 5 obtained from surface reflectance data (from 2000-2008) in the red (R)

27   and near-infrared (NIR) bands at 5 km resolution (CMG) are used to optimize the phenology-

28   related parameters of ORCHIDEE in the first step. The R and NIR data were processed to

29   correct for directional effects following Vermote et al. (2009) and then used to calculate the

30   NDVI, which is assumed to be linearly related to the model FAPAR. The NDVI are then i)

31   aggregated to the 0.72° spatial resolution of the ERA-Interim meteorological fields that are

1   used to force ORCHIDEE, ii) interpolated to a daily time series and iii) checked for quality

2   (see MacBean et al., 2015 for details). If there is a gap in the observations of more than 15

3   days, no interpolation is done (i.e., no data during the gap are assimilated). Figure 3 displays

4   the location of the sites that were selected (see §2.4.1).

5   ## 2.2.2  Eddy covariance flux data

6   Eddy covariance flux measurements of net surface $CO_2$ flux – hereafter referred to as net

7   ecosystem exchange (NEE) and latent heat flux (LE) from 78 observation sites of a network

8   of regional networks (FLUXNET; see Fig. 3) are used to constrain ecosystem physiology and

9   fast C-related processes at daily to seasonal timescales in ORCHIDEE in the second step. We

10  use quality-checked and gap-filled data from a global synthesis called the La Thuile dataset

11  (Papale, 2006). In order to avoid dealing with the large error correlations in the half-hourly

12  data (see Lasslop et al., 2008), daily mean values of NEE and LE are used in the ORCHIDAS.

13  Days with less than 80% of the half-hourly data are left out of the assimilation. The selection

14  of the sites and the data processing (gap-filling, correction for energy balance closure) are

15  detailed in Kuppel et al. (2014).

16  ## 2.2.3  Atmospheric $CO_2$ concentrations

17  Atmospheric $CO_2$ concentration measurements were taken from an ensemble of selected

18  surface stations around the world (Fig. 3). The spatial concentration gradients relate to the

19  integral of the fluxes over large areas and thus allow the optimization of large-scale global

20  patterns of carbon fluxes. These data were taken from the NOAA Earth System Laboratory

21  (ESRL) GLOBALVIEW-CO2 collaborative product (GLOBALVIEW-CO2, 2013) and

22  averaged to monthly means. We assimilated the monthly values for 53 sites for the 2002-2004

23  period inclusive in the last step of the assimilation system. Such restricted period (3 years

24  only) was chosen for practical reasons (computing resources) while constructing the

25  ORCHIDAS system. The station locations, indicated in Fig. 3, favor the background

26  conditions i.e. the surrounding air masses are only weakly influenced by local continental

27  sources, such as power plants. The choice of monthly mean is related to the use of pre-

28  calculated transport fields with LMDZ (see §2.3.2).

1 **2.3    Models and optimized parameters**

2 2.3.1 ORCHIDEE land surface model

3 In this study we use the ORCHIDEE process-oriented land surface model (Krinner et al.,

4 2005), which computes water, carbon and energy balances at the land surface on a half hourly

5 time step, using a mechanistic description of the physical and biogeochemical processes (see,

6 http://labex.ipsl.fr/orchidee/). The model describes the exchange of carbon and water at the

7 leaf level, the allocation of carbon within plant compartments (leaves, roots, heartwood and

8 sapwood), the autotrophic respiration, the production of litter, the plant mortality and the

9 degradation of soil organic matter (CENTURY model; Parton et al., 1988). The hydrological

10 processes for the soil reservoir rely on a double bucket scheme (Ducoudré et al., 1993). The

11 link between the water and carbon modules is via photosynthesis, which is based on the leaf-

12 scale equations of Farquhar et al., (1980) for C3 plants, and Collatz et al. (1992) for C4 plants,

13 that are then integrated over the canopy by assuming an exponential attenuation of light. The

14 FAPAR by each layer of the canopy is calculated from the leaf area index (LAI) following a

15 Beer-Lambert extinction law (Bacour et al., 2015).

16 ORCHIDEE uses the concept of the plant functional type (PFT) to describe the vegetation

17 distribution, with 13 PFTs (including bare soil) that can co-exist in each grid cell. Except for

18 the phenology (see a recent description in MacBean et al., 2015), the equations governing the

19 different processes are generic, but with specific parameter values for each PFT. Detailed

20 descriptions of model equations can be found in numerous publications (see for instance

21 Krinner et al., 2005). ORCHIDEE can be run at either global scale on a grid, or at site-level

22 using point-scale surface meteorological forcing variables. It is the land surface component of

23 the Institut Pierre Simon Laplace (IPSL) Earth System Model, and the version that we used

24 corresponds to CMIP5 simulations in the IPCC 5[th] Assessment Report (Dufresne et al., 2013).

25 However, in this study the model is run offline using the ERA-Interim 3-hourly near surface

26 meteorological forcing fields (Dee et al., 2011) aggregated at the spatial resolution of the

27 atmospheric transport model for the global simulations (see § 2.3.2). However, when we

28 assimilate in situ flux data in the second step, we force the model with the gap-filled half-

29 hourly meteorological data measured at each site. The global PFT map was derived from the

30 high-resolution IGBP AVHRR land data set (Vérant et al., 2004). The carbon pools are

31 brought to equilibrium (spin-up procedure) for both site and global scale simulations by

32 cycling the available meteorological forcing over several millennia, to ensure that the long-

1  term net carbon flux is close to zero. For the global simulation in third step, we spun-up the

2  model recycling the 1989-1998 meteorology and then used a transient simulation from 1990

3  to 2001 with changing climate (ERA-Interim) and increasing $CO_2$, before starting the

4  optimization with atmospheric data over 2002-2004. For the site simulations (i.e., the

5  assimilation of flux data) we recycled the available in situ meteorological forcing to spin-up

6  the model, with present day $CO_2$.

### 2.3.2 LMDz model

8  The transport model used in this study is version 3 of the General Circulation Model (GCM),

9  LMDz (Hourdin and Armengaud, 1999) with a horizontal resolution of 3.75° (longitude) x

10  2.5° (latitude) and 19 sigma-pressure layers up to 3 hPa. The calculated winds (u and v) are

11  relaxed to the ECMWF ERA-40 meteorological data (Uppala et al. 2005) with a relaxation

12  time of 2.5h (guiding) in order to realistically account for large-scale advection (Hourdin et

13  al., 2000). Deep convection is parameterized according to the scheme of Tiedtke (1989) and

14  the turbulent mixing in the planetary boundary layer is based on a local second-order closure

15  formalism. The LMDz GCM model has been widely used to model climate (IPCC, 2007,

16  2013) and its derived transport model has been used for the simulation of chemistry of gas

17  and particles and greenhouse gases distributions (Hauglustaine et al., 2004; Folberth et al.,

18  2005; Bousquet et al. 2005, 2006; Rivier et al., 2006). For this study, we used pre-calculated

19  transport fields, as described in Peylin et al. (2005), that correspond to the sensitivity of

20  concentration at each atmospheric site and each month to the surface flux of each model grid-

21  cell for each day (often called influence functions). The sensitivities (using inter-annual

22  winds) were calculated with the "retro-transport" formulation implemented in the LMDz

23  transport model (Hourdin et al. 2006). This approach decreases the computing time of the

24  optimization compared to the use of the full forward LMDz model at each iteration, as the

25  transport is replaced by a matrix multiplication with the vector of surface fluxes. Note that the

26  initial 3D state of the atmospheric concentrations was be defined from Chevallier et al. (2010)

### 2.3.3 Parameters optimized

28  The optimized parameters are described in Table 1, and their prior values, uncertainty and

29  range are given in Table 2. In the most recent studies using ORCHIDAS at site scales a large

30  set of ORCHIDEE parameters has been optimized (Kuppel et al., 2014; Santaren et al., 2014;

31  Bacour et al., 2015). In this study a smaller set was chosen, based on a Morris sensitivity

1   analysis (Morris, 1991; results not shown) that determines the sensitivity of the NEE and LE

2   to all model parameters at various FLUXNET sites (for each PFT), in order to reduce the

3   computational cost of the global optimization in step 3 (see §2.5). We considered 9 PFT-

4   dependent and 4 "global" (i.e. non PFT-dependent) parameters that control mostly the fast

5   carbon processes (diurnal to seasonal). In addition, we introduced a new parameter, $K_{soilC}$, to

6   scale the initial values (after spin-up) of the modeled slow and passive soil carbon pools, in

7   order to take account of all the historical effects not accounted for in the model that would

8   result in a disequilibrium of these pools in reality. For the site-specific optimizations with

9   FLUXNET data, we have one $K_{soilC,site}$ parameter per site. For the global scale optimization

10  step, we used 30 $K_{soilC,reg}$ parameters corresponding to 30 regions (see Fig. A2), thus the initial

11  soil carbon pools of all pixels within each region were scaled by the same value. The prior

12  value for all $K_{soilC}$ parameters was set to one, i.e. the default state of soil carbon pools is

13  assumed to be in equilibrium.

14  Overall (including all PFT-dependent parameters), we optimize 16 parameters related to

15  phenology, 36 to photosynthesis, 3 to respiration, 1 to the energy budget, 78 soil C pool

16  scalars (one for each FLUXNET site), and 30 regional soil C pool scalars for the global

17  simulations – a total of 184 parameters. Note that the soil C pool multipliers at the FLUXNET

18  sites are independent from the regional C pool multipliers, as the history of soil carbon over

19  large eco-regions of several millions square kilometers is rather heterogeneous (as it is mainly

20  related to previous land use changes), and most likely, the FLUXNET sites are not

21  representative of larger regions in terms of the soil carbon disequilibrium. The prior standard

22  deviation for each parameter is equal to 40% of the parameter range (lower and higher

23  boundaries) prescribed for each parameter following Kuppel et al. (2012). The parameter

24  ranges were specified following expert judgment of their meaning in the ORCHIDEE

25  equations and based on literature reviews or databases (such as TRY, Kattge et al., 2011).

26  **2.4   System description:  a step-wise approach**

27  2.4.1 Stepwise assimilation of three data streams

28  The ORCHIDAS system relies on a stepwise assimilation of the three data streams described

29  in section 2.2. Figure 2 illustrates the flow of information in this sequential approach:

30  ***Step 1 – Assimilation of MODIS-NDVI:*** Four parameters related to the seasonal cycle of the

31  vegetation (phenology) are optimized for the temperate and boreal deciduous PFTs (TeBD,

1  BoND, BoBD and NC3 – see caption of Table 2). These four deciduous PFTs alone are

2  considered in step 1 in this ORCHIDAS version because the tropical deciduous phenology

3  modules in ORCHIDEE require further modifications to improve the functions that control

4  leaf growth and fall in response to water availability (MacBean et al., 2015). Evergreen PFTs

5  were also not considered, as the there are no phenology modules related to these PFTs in the

6  model. The procedure is similar to that described in detail in MacBean et al. (2015) and

7  therefore only briefly recalled here. A simple linear relationship between the modeled

8  Fraction of Absorbed Photosynthetically Active Radiation (FAPAR) and MODIS-NDVI

9  observations is assumed, based on studies such as Knyazikhin et al. (1998). Following Bacour

10  et al. (2015), we use only the temporal information in the NDVI observations and not the

11  actual values, and thus we normalized both the model FAPAR output and the NDVI

12  observations to their $5^{th}$ and $95^{th}$ percentiles. The model was run for fifteen randomly selected

13  grid cells for each of the four PFTs using the ERA-Interim meteorological forcing. Only grid

14  cells that included vegetation fraction of greater than 60% for the PFT optimized were

15  considered. The fifteen sites for each PFT were included in one optimization for each PFT

16  following a multi-site approach, in which all observations are used simultaneously to optimize

17  the model parameters. The optimized parameters are described in Table 1. They correspond to

18  a scalar on the growing degree days (GDD) threshold for the start of the vegetation ($K_{pheno,crit}$),

19  a parameter controlling the use of carbohydrate reserve during the start of leaf growth

20  ($K_{lai,happy}$), a temperature threshold for the onset of leaf senescence ($CT_{,senes}$) and the critical

21  age for leaves ($L_{agecrit}$).

22  ***Step 2 – Assimilation of FLUXNET data:*** Mean daily NEE and LE flux measurements for 78

23  sites, including up to 10 years worth of data for each site, are used to optimize a set of model

24  parameters controlling the fast carbon and water processes (photosynthesis, respiration,

25  phenology – see Table 1). The site selection and the choice of a daily time step are described

26  in more details in Kuppel et al. (2014). These sites cover 7 of the PFTs in ORCHIDEE (see

27  Table 2). The posterior parameter values of the four phenology parameters derived in step 1,

28  and their associated uncertainties, are input as prior information in step 2. For the additional

29  parameters, the default ORCHIDEE values are used for the prior and the uncertainties are set

30  as described in §2.3.3. A multi-site optimization is performed for each PFT independently as

31  in step 1. Global parameters, i.e. those that are not PFT-dependent, were optimized for each

32  PFT and the mean across all PFTs was then calculated to define the prior parameter vector in

33  step 3 of the assimilation with atmospheric $CO_2$ data (at global scale). Such an approach was

1  chosen to allow us to optimize all PFTs in parallel and therefore to simplify the assimilation

2  process.

3  ***Step 3 – Assimilation of atmospheric CO$_2$ concentrations:*** We use monthly mean CO$_2$

4  concentrations from 53 surface stations over three years (2002-2004) to provide a large-scale

5  constraint to the land surface fluxes (i.e. to match the global CO$_2$ growth rate, mean seasonal

6  cycle and its latitudinal variation, as well as the spatial gradients between stations). We use

7  the LMDz atmospheric transport model (see §2.3.2) to assimilate these observations. The set

8  of parameters optimized in step 2 are included in step 3, except for the albedo scaling

9  parameter ($K_{albedo,veg}$), as the net carbon fluxes are only weakly sensitive to that parameter. We

10 used the posterior parameter distributions from step 2 (parameter optimal values and

11 associated uncertainties) as prior information for step 3, and expanded the parameter vector to

12 include the 30 $K_{soilC}$ parameters that scale the initial soil carbon pools for large "spatially-

13 coherent regions" (see §2.1.2 and Fig. A2). The air-sea fluxes and fossil fuel and biomass

14 burning emissions are also accounted for (but not optimized) in this final step, in order to

15 close the global carbon budget within the atmospheric transport model (see §2.5).

16 ## 2.4.2 Optimization procedure (for all steps):

17 In each step the statistically optimal parameter values are derived with an optimization

18 procedure following the principle of the 4-D variational assimilation systems (developed for

19 numerical weather prediction), using a tangent linear operator (and finite differences for a few

20 parameters, Bacour et al. 2015). Assuming that the errors associated with the parameters, the

21 observations and the model outputs follow Gaussian distributions, the optimal parameter set

22 corresponds to the minimum of a cost function, *J(x),* that measures the mismatch between i)

23 the observations (*y*) and the corresponding model outputs, *H(x),* (where *H* is the model

24 operator), and ii) the a priori (*x$_b$*) and optimized parameters (*x*), weighted by their error

25 covariance matrices (Tarantola, 1987; Eq. (1)):

26
$$J(\boldsymbol{x}) = \frac{1}{2}\left[(H(\boldsymbol{x}) - \boldsymbol{y})^T \, \mathbf{R}^{-1} \, (H(\boldsymbol{x}) - \boldsymbol{y}) \quad + (\boldsymbol{x} - \boldsymbol{x}_b)^T \mathbf{B}^{-1} \, (\boldsymbol{x} - \boldsymbol{x}_b) \quad \right]$$
(1)

27 **R** represents the error variance/covariance matrix associated with the observations and **B** the

28 parameter prior error variance/covariance matrix. At each step a different cost function is

29 defined with the observations and parameters related to that step (see Fig. 2). **R** includes the

30 errors on the measurements, the model structure and the meteorological forcing. Model errors

31 are rather difficult to assess and may be much larger than the measurement error itself.

Geoscientific
Model Development
Discussions

1    Therefore we chose to focus on the structural error and defined the variances in **R** as the mean

2    squared difference between the prior model and the observations for both step 1 and step 2

3    (see Kuppel et al. 2013). For simplicity we assumed that the observation error covariances

4    were independent between the different observations and therefore we kept **R** diagonal (off-

5    diagonal terms set to zero), given the rapid decline of the model error auto-correlation beyond

6    one day (Kuppel et al., 2013). For step 3 we used a different approach, given the large bias in

7    the model a priori concentrations, and therefore followed the methodology of Peylin et al.

8    (2005) based on the observed and modeled temporal concentration variability at each site.

9    Overall, data uncertainties in the optimization procedure are between 0.1 and 0.45 for NDVI

10    (step 1), around 3-6 $gCm^{-2}d^{-1}$ for daily NEE, and 15-30 $Wm^{-2}$ for daily LE (step 2) and

11    between 0.1 ppm at remote oceanic stations and 4 ppm at continental sites (step 3).

12    The determination of the optimal parameter vector that minimizes *J(x)* is performed by

13    successive calls to a "gradient-descent" minimization algorithm L-BFGS-B (Byrd et al.

14    1995), which is specifically dedicated to solving large nonlinear optimization problems that

15    are subject to simple bounds on the parameter values. In order to find the minimum of *J(x)* the

16    algorithm requires the gradient of *J(x)* (Jacobian) with respect to the ORCHIDEE parameters.

17    L-BFGS-B explores each parameter space simultaneously along the gradient of the cost

18    function, and uses an approximation of the Hessian (second derivative) of *J(x),* which is

19    updated at each iteration, to define the size of the step at each iteration.

20    For step 1 and step 2, the model "*H*" simply corresponds to the land surface model: $H = S$,

21    with *S(x)* representing the surface fluxes from the ORCHIDEE model using the parameter

22    vector, *x*. The gradients *dJ(x)/dx* are calculated from the tangent linear model of ORCHIDEE

23    that was automatically generated by the numerical Transformation of Algorithms in Fortran

24    (www.fastopt.de), except for two parameters linked to the model phenology for which the

25    threshold functions prevent the use of a linear approximation. A finite difference approach

26    was used for these parameters.

27    For step 3, the model "*H*" corresponds to the composition of the land surface model with the

28    transport model: $H = T$ o $S$ (see Kaminski et al. (2002) for details), with *T* representing the

29    LMDz transport model. *T* is a linear operator for a non-reactive species: $T(S(x)) = \mathbf{T} \cdot S(\mathbf{x}),$

30    with **T** a matrix representation of the transport operator. It corresponds to the sensitivity of

31    $CO_2$ concentrations at each site and for each month to the daily surface flux of each model

32    grid-cell. It is then combined with the ORCHIDEE surface fluxes (*S(x)*) through a matrix

1  multiplication to derive *H(x)*. **T** has been pre-calculated for all atmospheric stations in order

2  to save computing time during the iterative optimization process (see §2.3.2). For simplicity

3  we use monthly mean values for both the fluxes *S(x)* and the transport sensitivities (**T**) in the

4  computation of the gradients *dJ(x)/dx*.

5  For improved minimization efficiency, the inversion is preconditioned (following Chevallier

6  et al., 2005), which means that L-BFGS-B is fed with the control variable $\mathbf{x'} = \mathbf{B}^{-1/2}(\mathbf{x} -$

7  $\mathbf{x_b})$, rather than with $\mathbf{x}$, as this homogenizes the range of variation of the optimized

8  parameters.

### 2.4.3 Error estimation

10  The posterior parameter error covariance matrix, **A**, can be approximated to the inverse

11  Hessian of the cost function, using the linearity assumption at the minimum of *J(x)*. It can be

12  derived with the Jacobian of the model at the end of the minimization (i.e. the last iteration),

13  $\mathbf{H_\infty}$, following Tarantola (1987):

14
$$\mathbf{A} = [\mathbf{H_\infty^T}.\mathbf{R}^{-1}.\mathbf{H_\infty} + \mathbf{B}^{-1}]^{-1}$$
(4)

15  Note that for step 3, $\mathbf{H_\infty} = \mathbf{T}.\mathbf{S_\infty}$, where $\mathbf{S_\infty}$ is the Jacobian of the ORCHIDEE model at the

16  last iteration. The posterior parameter error covariance, **A**, can then be propagated into the

17  model state variable space (e.g. carbon fluxes and stocks), $\mathbf{A_{var}}$, given the following matrix

18  product (only used for the global fluxes in step 3):

19
$$\mathbf{A_{var}} = \mathbf{S_\infty}.\mathbf{A}.\mathbf{S_\infty^T}$$
(5)

20  The square root of the diagonal elements of $\mathbf{A_{var}}$ corresponds to the standard deviation, $\sigma$, of

21  carbon fluxes/stocks for each grid cell. In order to evaluate the knowledge improvement

22  brought by the assimilation, the uncertainty reduction between the prior ($\sigma_{prior}$) and posterior

23  ($\sigma_{post}$) is determined as $1 - (\sigma_{post} / \sigma_{prior})$.

### 2.4.4 Additional processing steps

25  In order to analyze the fit to the atmospheric $CO_2$ concentrations in terms of the trend and

26  seasonal cycle, we decomposed the observed and modeled time series by fitting the monthly

27  mean values with a function comprising a first order polynomial term and four harmonics,

28  and then filtered the residuals of that function in frequency space using a low pass filter

29  (cutoff frequency of 65 days), following Thoning et al. (1989). The polynomial term defines

1    the trend while the seasonal cycle corresponds to the harmonics plus the filtered residuals.

2    The amplitude of the seasonal cycle is then calculated as the difference between the monthly

3    mean maximum and minimum for year 2003 (middle year of the optimization period).

4    Finally, we define the Carbon Uptake Period (CUP) as the sum of the days when the values of

5    the seasonal cycle extracted from the $CO_2$ concentration time series are negative (a negative

6    convention being for $CO_2$ removed from the atmosphere).

7    **2.5 Prescribed emissions of carbon fluxes**

8    In this section we describe the other components of the carbon cycle (apart from the surface C

9    exchange with terrestrial vegetation) that are imposed in step 3 of the optimization process as

10   fixed fluxes.

11   2.5.1 Ocean fluxes

12   The ocean contributes to an uptake of about a quarter to a third of the anthropogenic

13   emissions with significant year-to-year variations (Sabine et al., 2004). For this version of the

14   ORCHIDAS, we developed a statistical model to estimate the spatial and temporal variations

15   (monthly) of the ocean surface $CO_2$ partial pressure ($pCO_2^{SW}$), and from that the air-sea $CO_2$

16   fluxes, using satellite and in-situ ocean measurements and model outputs. The air-sea $CO_2$

17   fluxes are primarily controlled by the ocean biogeochemistry, the horizontal transport and the

18   vertical mixing in the ocean and the atmospheric forcing ($CO_2$ partial pressure at the interface

19   to the water ($pCO_2^{ATM}$) and wind); they can be defined from the following equation:

20   $$F_{CO2} = K_{ex} \times (pCO_2^{SW} - pCO_2^{ATM}) \qquad (6)$$

21   where $K_{ex}$ stands for the exchange coefficient and $F_{CO2}$ the $CO_2$ flux from the sea surface

22   water to the atmosphere.

23   The computation of $pCO_2^{SW}$ is performed using feedforward artificial neural networks, i.e., a

24   MultiLayer Perceptron (MLP; Rosenblatt 1958) that maps a set of spatio-temporal variables

25   (input) onto observed $pCO_2^{SW}$ data. We use a two-step approach: the first step to derive a

26   monthly mean $pCO2^{SW}$ climatology and the second step to correct for the year to year

27   variations. The $pCO_2^{SW}$ observations come from the Global Surface pCO2 (Lamont-Doherty

28   Earth Observatory, LDEO) Database (Takahashi et al., 2009). The inputs are a series of

29   variables connected to the spatial and temporal evolution of $pCO_2^{SW}$: i) sea surface

30   temperature (SST), sea surface salinity (SSS) and mixed layer depth (MLD) as a proxy of the

1   physical processes (these fields come from a re-analysis of the NEMO-OPA ocean model

2   (Madec et al., 1998) with the assimilation of several satellite observations), ii) chlorophyll

3   content from SeaWiFS, as a proxy of the biogeochemistry (CHL), iii) spatial and temporal

4   coordinates (LAT, LON and MONTH) and the $pCO_2^{SW}$ at previous time step (recursive

5   approach), i.e.:

6   $$\{PCO_2^{SW}\}_m = MLP\left(\{SST, SSS, MLD, CHL\}_{(m-2,m-1,m)}, \{PCO_2^{SW}\}_{(m-2,m-1)} \, LAT, LON\right) \quad (7)$$

7   with $m$ the monthly index. The available data (20685 points) is divided into two parts: 75% is

8   used for the learning phase of the ANN and 25% for the validation phase. The overall

9   performance of the neural network for extrapolating the spatial and seasonal distribution of

10  $pCO_2^{SW}$ is relatively good, with a spatio-temporal correlation coefficient between the

11  estimated $pCO_2^{SW}$ and the independent observations of 0.80.

12  $pCO_2^{ATM}$ at the surface are taken from a global simulation of atmospheric $CO_2$ concentrations

13  with optimized fluxes (Chevallier et al. 2010). $K_{ex}$ is defined as the product of $k$, the gas

14  transfer velocity, taken from the Wanninkhof (1992) formulation using winds from ERA-

15  Interim, and $s$, the solubility of $CO_2$, taken from the Weiss formulation (Weiss, 1974). The

16  system is further described in Roedenbeck et al. (2015). The global ocean sink is around 1.60

17  $PgC.yr^{-1}$ for the period 2002-2004 used in step 3. It is within the uncertainty range of the

18  Global Carbon Project estimates (Le Quéré et al., 2015) if we account for the pre-industrial

19  ocean out-gazing flux included in our "delta $pCO_2$" approach. Its temporal evolution is

20  depicted in Fig. A1

21  ## 2.5.2 Global fossil fuel and cement $CO_2$ emissions

22  We have used a recently developed $CO_2$ fossil fuel and cement emission product (see

23  http://www.carbones.eu/wcmqs/) that covers the period 1980 to 2009 at the spatial resolution

24  of 1° x 1° and hourly resolution. It is based on EDGAR v4.2 spatially distributed annual

25  emissions (Olivier et al., 2012) and time profiles developed by the University of Stuttgart. It

26  was assumed that EDGAR delivers the most up-to-date spatially distributed and sector

27  specific emissions, based on national emission statistics. IER further applied country and

28  sector specific time profiles, taking into account monthly, daily, and hourly variations

29  depending on the sector. The derivation of the time profiles relies on different data sets (e.g.

30  Eurostat, ENSTO-E, UN monthly bulletin) as well as correlations between recorded

31  emissions and climate variables. Currently, the temporal profiles are derived mostly from data

Geoscientific
Model Development
Discussions

1    sets over Europe that were extrapolated using information on climate zone, average monthly

2    temperature for the seasonal cycles and similarity in socio-economic parameters like

3    population and Gross Domestic Product (GDP). The annual mean emission for the period

4    2002-2004 is 7.14 PgC.yr$^{-1}$.

5    ## 2.5.3 Fire emissions:

6    Fire    emissions    data    from    the    Global    Fire    Data    (GFEDv3    –

7    http://www.globalfiredata.org/Data/index.html) are prescribed in the ORCHIDAS. The

8    GFEDv3 data are broken-down into 6 sectors (deforestation, peat fires, savanna fires,

9    agriculture, forest fires, and woodland) that are further grouped into 3 main types. We

10   generated fluxes of $CO_2$ relevant for typical "burning - regrowth" processes, as detailed in

11   Appendix A2. The first type corresponds to deforestation and peat fires with carbon

12   permanently lost to the atmosphere, the second to agriculture and savannah fires which are

13   assumed to be compensated by a sink during the regrowth period (i.e. with zero annual net

14   emission for each pixel) and the third to woodland and burnt forests which are assumed to be

15   at steady state for a given region (10 sub-continental scale regions) over the period covered by

16   GFEDv2 (i.e. regrowth of nearby forest compensates for the burned forest derived in GFED).

17   The sum of these three components leads to the global flux, with a gross emission around 2.1

18   PgC.yr$^{-1}$ and a net emission after regrowth of only 1.1 PgC.yr$^{-1}$ (Fig. A2 in Appendix) that is

19   prescribed to the ORCHIDAS over the period 2002-2004.

20

21   **3   Results**

22   **3.1 Model fit to the data**

23   3.1.1 Step 1: assimilation of MODIS NDVI data

24   The optimization in Step 1 resulted in an improved fit to the MODIS NDVI observations for

25   the four PFTs considered (TeBD, BoND, BoBD, NC3, see §2.4) as seen in Fig. 4, which

26   shows the mean seasonal cycle across the 2000-2008 period for all sites for each PFT. The

27   most prominent change after the optimization was a substantially shorter growing season for

28   all PFTs due to an earlier start of leaf senescence. This was caused by both a lower critical

29   leaf age ($L_{agecrit}$) and a higher temperature threshold for senescence ($CT_{senes}$) (Fig. 8). The

30   impact on the start of leaf growth was less dramatic but important nonetheless, with a shift to

18

1  a later start of leaf growth as a result of an increase in the $K_{pheno,crit}$ parameter which acts as a

2  scalar on the threshold of Growing Degree Days (GDD) used to trigger leaf onset (see

3  Appendix A in MacBean et al., 2015). Overall, a mean reduction in RMSE of 23, 17, 58 and

4  19% was achieved for TeBD, BoBD, BoND trees and NC3 grasses respectively, with the

5  greatest improvement for BoND trees. The mean correlation between the normalized MODIS-

6  NDVI and modeled FAPAR time series over the 2000 – 2008 period increased for TeBD and

7  BoND trees and NC3 grasses (prior and posterior of 0.9 to 0.93, 0.42 to 0.91 and 0.6 to 0.66,

8  respectively). The prior correlation of 0.55 remained similar after the assimilation for BoBD

9  trees.

10  Following the improvement at the sites selected for the optimization, we evaluated the impact

11  for each PFT at the global scale using the global median correlation between the MODIS-

12  NDVI and the model FAPAR time series (from all pixels where the fraction of a given PFT is

13  above 60%, see Maignan et al. 2011). The global correlation increased for BoND trees and

14  NC3 grasses from 0.36 to 0.91 and 0.53 to 0.59 (prior to posterior), respectively. It remains

15  stable for BoBD (0.54) or slightly increased for TeBD (0.88 to 0.89).

16  ### 3.1.2 Step 2: assimilation of FLUXNET data

17  The optimization in Step 2 brings an improvement to the simulated NEE and LE for all seven

18  PFTs considered, with Fig. 5 showing the corresponding PFT-averaged mean NEE seasonal

19  cycles (mean across all sites/years). NEE is overestimated by the prior model for all PFTs on

20  average. This is partly due to the model spin-up procedure, which brings each simulated site

21  to a near equilibrium state with a mean NEE close to zero (i.e. no net carbon sink, see §2.1.1).

22  This bias is significantly corrected by the optimization to match the observed carbon uptake at

23  most sites, notably via the scaling of the initial soil carbon pool content at each site

24  (parameters $K_{soilC,site}$; Table 1) which thus significantly reduces the ecosystem respiration

25  (Kuppel et al., 2014). Overall, the largest reductions of model-data RMSE are found in

26  temperate forests (TeNE, TeBE and TeBD), where the RMSE decreased by more than 25%

27  compared with the prior model. The improvements are less significant for the other PFTs,

28  with RMSE reductions between 10 and 18%.

29  In addition, the optimization increases the NEE seasonal amplitude in temperate evergreen

30  forests (TeNE and TeBE) and temperate broadleaf deciduous forests (TeBD), and reduces the

31  amplitude for boreal needle leaf forest (BoNE) and natural C3 grasses (NC3), in agreement

Geoscientific
Model Development
Discussions

1    with the observations (except for BoNE where the amplitude decrease is too large). Despite

2    the better model-data agreement for evergreen broadleaf forests (TrBE and TeBE), the

3    optimized model still fails to catch some seasonal features such as a persistent carbon uptake

4    (i.e. negative NEE) in the dry season for the tropical regions (TrBE) and nearly-null carbon

5    exchange in the first months of the year for temperate regions (TeBE). These results are

6    discussed further in Kuppel et al. (2014), who used a similar optimization set-up with a

7    slightly different parameter set – see §2.3.3. Similar improvements, although of smaller

8    amplitude, occur for the latent heat fluxes (not shown).

9    ### 3.1.3 Step 3: assimilation of atmospheric $CO_2$ data

10    The final optimization step with the atmospheric $CO_2$ concentrations provides a large

11    improvement of the fit to the observed concentrations at most stations. The cost function $J$

12    was reduced through the minimization by a factor of 5.7 within 37 iterations.

13    Figure 6 illustrates the simulated concentrations for four stations (representative of different

14    conditions) with the standard prior parameter vector (used in step 1), the posterior vector from

15    step 2 (used as prior in step 3) and the posterior vector from this last step. The improvement

16    in the fit to the observations can be quantified with the reduction in RMSE (from the prior to

17    the posterior of step 3) - the largest reduction is at the South Pole station (73%) and is on

18    average around 25% across all sites. Note that for a few stations the fit is slightly degraded

19    (up to 10%) except for one Pacific site (regular ship measurements around the equator,

20    POCN00) for which there is a 40% degradation, possibly due to small biases in the simulation

21    of the ITCZ position in LMDz. When calculated with respect to the standard prior (used in

22    step 1) the RMSE decrease is slightly larger on average, especially for the northern mid to

23    high latitude stations. For these stations the optimization performed in step 2 with FLUXNET

24    data led to a significant improvement of the mean seasonal cycle amplitude of the

25    atmospheric $CO_2$ data, as discussed in Kuppel et al. (2014).

26    We then investigated the fit to the observed $CO_2$ concentrations in terms of the mean seasonal

27    cycle and trend (see section 2.4.4). With only three years of data the mean trend is more

28    difficult to define as it varies between stations; however, the optimization in step 3 increases

29    the net land carbon sink in order to match the observed trend at most stations. If we take the

30    Mauna Loa and South Pole stations that are representative of an integration of the fluxes at

31    hemispheric scales, the prior $CO_2$ trend of 2.8 and 2.9 ppm.yr$^{-1}$ respectively, is reduced to 2.1

1    and 2.2 ppm.yr$^{-1}$ close to the observations (2.1 ppm.yr$^{-1}$ for both). The left panel of Fig. 7

2    illustrates changes in the amplitude of the simulated seasonal cycle at each station (see

3    definition in §2.4.4). The values correspond to relative changes between the prior and

4    posterior of the absolute difference between observed and modeled amplitude ($[\left|\Delta A_{poste}\right| -$

5    $\left|\Delta A_{prior}\right|]/\left|\Delta A_{prior}\right|$). They reveal an improvement in the seasonal cycle amplitude at nearly

6    all stations of the southern hemisphere ($\approx$ 40% improvement) and at the majority of the

7    northern hemisphere stations ($\approx$ 15%). A few stations in north East Asia (3) and northwest

8    America (4) show a small degradation of the amplitude ($\approx$ 15%). The right panel of Fig. 7

9    displays the changes of the Carbon Uptake Period (CUP, see §2.4.4) expressed in terms of

10    relative changes between prior and posterior of the absolute values of model-data differences,

11    as for the amplitude. Most stations reveal an improvement of the CUP of around 20%, which

12    is slightly lower than the improvement for the seasonal cycle amplitude.

### 13    3.2 Consistency of the step-wise optimization

14    The main issue with a step-wise data assimilation system (versus a simultaneous approach)

15    concerns the potential degradation of the model – data fit for the different data streams that

16    are assimilated in previous steps. We noted that $CO_2$ concentrations were already improved

17    when NDVI and FLUXNET data are assimilated (see §3.1.3), but we need to check if the

18    final parameter set from step 3 leads to a degradation of the fit to MODIS-NDVI (step 1) and

19    to FLUXNET (step 2) data compared to the fit achieved during the respective steps and, in the

20    case of a significant degradation, if we still have an improvement for these data streams

21    compared to the initial *a priori* fit.

22    Figure 8 summarizes the performance of the model data fit for MODIS-NDVI and

23    FLUXNET-NEE data streams for the prior and posterior of each step by evaluating the

24    median RMSE between the model and the observations across all sites. The values are

25    calculated for each PFT separately. In this section, we keep in mind the fact that we do not

26    optimize the same PFTs with FLUXNET data and with MODIS-NDVI.

### 27    Consistency for MODIS-NDVI

28    First, we notice again the significant RMSE reduction between the prior and step 1, as

29    discussed in section 3.1. The fit to MODIS-NDVI (normalized data) for step 2 and step 3

30    shows only a significant degradation (increased RMSE) for temperate broadleaf deciduous

Geoscientific
Model Development
Discussions

1   forest (TeBD), which decreases the improvement achieved in step 1 (compared to the prior)

2   by a factor of two. A marginal degradation for natural C3 grassland (NC3) is obtained after

3   step 3: the RMSE increases slightly from 0.24 to 0.26, but is still lower than the prior value of

4   0.3. There is no degradation for boreal needleleaf deciduous trees (BoND), but a surprising

5   small decrease of the RMSE (i.e. improvement in the model-data fit) for boreal broadleaf

6   deciduous forests (BoBD; from 0.26 to 0.23). In this latter case, the use of additional

7   parameters in steps 2 and 3 (see §2.4) allows further improvement of the fit between the

8   normalized FAPAR and NDVI time series. On average the degradation of the fit to NDVI is

9   thus very limited in step 2 and step 3, and in no case is the RMSE worse than the prior.

10  Consistency for FLUXNET data

11  Figure 8 again reveals the significant reduction of the RMSEs for NEE in step 2 compared to

12  the standard prior or to the posterior of step 1 for most PFTs, except BoNE. We see only

13  small degradations (increases) in RMSE between step 2 and step 3 for temperate needle leaf

14  evergreen forests (TeNE: from 1.06 to 1.13 $gC.m^2.d^{-1}$), temperate broadleaf evergreen forests

15  (TeBE: from 1.06 to 1.09 $gC.m^2.d^{-1}$), temperate broadleaf deciduous forests (TeBD: from 1.06

16  to 1.13 $gC.m^2.d^{-1}$) and boreal needle leaf evergreen forests (BoNE: from 0.59 to 0.60 $gC.m^2.d^{-}$

17  $^1$). An interesting feature to notice is that the NEE RMSE increases between the prior to the

18  posterior of step 1 (i.e. before NEE has been used in the optimization in step 2). Using remote

19  sensing products of vegetation activity or "greenness" (e.g. NDVI) to calibrate the phenology

20  of ORCHIDEE thus does not always improve the simulated NEE, the possible reasons for

21  which were discussed in Bacour et al. (2015) who used the same LSM and assimilation

22  system. Overall, the reduction of the improvement of the model data fit to the NEE (step 3

23  versus step 2) is marginal (limited to a few percent), thus further suggesting the consistency of

24  our step-wise approach. Similar results are also obtained for the latent heat flux (LE) (not

25  shown).

26  **3.3 Estimated parameter values**

27  We now discuss the parameter values, focusing on the changes obtained though the

28  successive steps. Figure 9 presents the prior and posterior values for each parameter together

29  with their associated uncertainties (estimated through Eq. (4)) and the allowed range of

30  variation. Note that nine parameters are PFT-dependent while four are global (non PFT-

31  dependent). For the global non PFT-dependent parameters included in the step 2 optimization,

1  we took the mean value (see §2.4) as the prior for step 3. Note finally that the parameters

2  linked to the initial soil carbon pools ($K_{soilC,site}$, $K_{soilC,reg}$) are not shown in Fig. 9 as they are

3  too numerous (though see Fig. A2 for the regional values).

4  If we first consider the phenology parameters optimized in step 1 ($K_{lai,happy}$, $K_{pheno,crit}$, $L_{age\_crit}$,

5  $C_{T,senes}$; see Table 1) we see that for most PFTs they do not change significantly between step

6  1 and step 3, although they differ significantly from the prior. There are few exceptions,

7  including $K_{pheno,crit}$ (the threshold for the start of the growing season) for Boreal Needleleaf

8  deciduous forests and $K_{lai,happy}$ (level of carbohydrate use) for temperate and boreal broadleaf

9  deciduous forests (TeBD, BoBD). Note that a few phenology parameters hit one of the

10  physical bounds, which may indicate model structural errors or model parameter equifinality.

11  For most phenology parameters, the uncertainties are strongly reduced during their first

12  optimization (step 1), except for a few cases like $C_{T,senes}$ for C3 grassland. Note finally that a

13  more in depth spatio-temporal validation demonstrated the generality of the optimized

14  phenology parameters across multiple sites (for further details see MacBean et al., 2015).

15  For the photosynthesis parameters ($V_{cmax}$, $G_{s,slope}$, $C_{Topt}$, $SLA$, $f_{stress}$; see Table 1), we find a

16  similar result with little changes between step 2 and step 3, but still a significant departure

17  from the prior values. Most parameters are well constrained by the inversion, with posterior

18  uncertainties that are greatly reduced compared to the prior, except for Tropical broadleaf

19  rain-green forest (TrBR) and Boreal needle-leaf deciduous forest (BoND) for which there is

20  nearly no constraint on $G_{s,slope}$, and $f_{stress}$ (see Table 1).

21  The non-PFT dependent respiration-related parameters ($HR_{H,c}$, $Q_{10}$, $MR_b$) mostly change in

22  step 2 and only slightly in step 3 (with an additional reduction of the error) in order to fit the

23  large-scale constraint provided by the atmospheric observations. The values of the scalar of

24  the initial soil carbon pools size for the FLUXNET site optimizations ($K_{soilC,site}$, one parameter

25  per site, not shown) were largely reduced on average, in order to decrease the heterotrophic

26  respiration (see Kuppel et al. (2014) for additional discussion). In step 3 the same scalars that

27  were defined for an ensemble of large regions ($K_{soilC,reg}$) have decreased in the southern

28  hemisphere (less than 10%; see Fig. A2 in Appendix A3) and slightly increased in the

29  northern hemisphere (around 1%), to achieve a better match to the atmospheric $CO_2$ growth

30  rate and north-south gradient. Importantly, we notice that for step 3, the fit to the atmospheric

31  $CO_2$ concentrations (especially to the trend) is achieved mainly by small changes in $K_{soilC,reg}$

32  and in few other respiration-related parameters. Note finally that the parameter controlling the

1   albedo ($K_{albedo,veg}$), modified with the FLUXNET observations only (see §2.4), is not well

2   constrained by the optimization (only a small reduction in uncertainty). Overall, most

3   parameters appear to be well constrained when first optimized, with only small changes in the

4   following steps. This suggests that the targeting of different parameter subspaces in the

5   various optimisation steps was well-chosen.

6   **3.4 Estimated carbon fluxes and uncertainties**

7   The main objective of a carbon cycle data assimilation procedure is to improve the simulated

8   land surface net and gross carbon fluxes as well as the simulated carbon stocks for both

9   present and future conditions. Given the focus of the paper, i.e. to describe the potential of a

10  step-wise global carbon cycle data assimilation system, we only discuss a few large-scale

11  features of the optimized annual net and gross carbon fluxes, as well as one of the carbon

12  stock variables (forest above-ground biomass). We thus do not discuss the inter-annual flux

13  variability.

14  Large-scale annual mean net fluxes

15  The mean annual carbon fluxes (NEE) for the globe, northern extra tropics, tropics, and

16  southern extra tropics are reported in Fig. 10 for the 2000-2009 decade for the prior and

17  posterior model simulations for all steps together with one other estimate of the land surface

18  residual from the Global Carbon Project (GCP, Le Quéré et al, 2015) over the same decade.

19  The prior NEE indicates a total sink of 0.5 PgC.yr$^{-1}$ over this period, from both the northern

20  and tropical regions. Such a prior sink is due to the increase of atmospheric $CO_2$ during the

21  transient simulation following the spin-up (1990-2009, see section 2.3.1) and climate

22  variability. Changes from the prior are rather small in step 1 (assimilation of MODIS NDVI))

23  with an increase of the northern sink by 0.12 PgC.yr$^{-1}$ and a decrease of the tropical sink by

24  0.05 PgC.yr$^{-1}$ (Fig. 10). Step 2 (assimilation of FLUXNET data) does not significantly change

25  the net C sink from step 1, with only a small increase in the tropical sink by 0.1 PgC.yr$^{-1}$. The

26  assimilation of atmospheric $CO_2$ data in step 3 provides a large-scale constraint, as already

27  discussed, and increases the total land sink to 2.2 PgC.yr$^{-1}$, a value in much closer agreement

28  with the estimate by the GCP. A larger tropical NEE uptake is responsible for the large

29  increase of the terrestrial biosphere C sink (from 0.3 PgC.yr$^{-1}$ in step 2 to 1.7 PgC.yr$^{-1}$) while

30  the sink in the north increases by less than 0.1 PgC.yr$^{-1}$. The comparison with the GCP

31  number should be taken with caution. The ORCHIDAS estimated sink include all effects

Geoscientific
Model Development
Discussions

1 (natural and anthropogenic), since that we used atmospheric $CO_2$ as a global constraint. Thus

2 the optimized parameters must account for any missing processes like nitrogen limitation or a

3 proper description of agricultural processes and management. However, the GCP number is

4 only for the anthropogenic uptake, excluding the pre-industrial sink due for instance to river

5 export of carbon (around 0.45 PgC.yr$^{-1}$; Regnier et al. 2013).

6 Spatial distribution of the annual mean net flux

7 Figure 11 shows the spatial distribution of NEE averaged over 2002-2004 for the standard

8 prior and posterior after step 3. The large tropical net land carbon sink that is inferred in step

9 3 is mainly explained by an increase of the carbon uptake for the tropical forests of the

10 Amazon basin and equatorial Africa, as well as a decrease of the carbon release on the

11 southern edge of the Amazon basin (tropical rain-green forests and grasses). In the northern

12 mid-high latitudes only smaller regional changes from the prior occur. For Europe, most of

13 north Asia and Canada, the strength of the C sink slightly decreased from the prior (up to 30

14 gC.m$^2$.yr$^{-1}$), while for central USA the strength of C source slightly decreased. If we now

15 consider the uncertainties on the net annual carbon flux that arise from the parameter

16 uncertainty (second row of Fig. 10; Eq. (5)) we observe a very large reduction (compared to

17 the prior) in the monthly flux uncertainty (averaged over the three years used in step 3) over

18 tropical forests. It is reduced by a factor four with initial values around 150 gC.m$^2$.y$^{-1}$ and

19 posterior values between 22 and 66 gC.m$^2$.y$^{-1}$. For mid-to-high latitude boreal ecosystems, the

20 uncertainty reduction is smaller, but the posterior errors are slightly lower than over the

21 tropics, between 18 and 55 gC.m$^2$.y$^{-1}$.

22 Large-scale annual mean Gross Primary Production (GPP)

23 For the GPP the relative changes from the prior are smaller than for the NEE (Fig. 10b). The

24 mean annual total GPP is 169, 160, 154 and 156 PgC.yr$^{-1}$ for the prior and posterior of step 1,

25 2 and 3, respectively. The small overall decrease (8%) brings the GPP slightly closer to the

26 estimate by Jung et al. (2011), around 120 PgC.yr$^{-1}$, based on a statistical Model Tree

27 Ensemble (MTE) that upscaled the in-situ flux measurements (resulting from the partition of

28 measured NEE into GPP and total ecosystem respiration). The decrease in GPP occurs mainly

29 in the northern hemisphere after step 1 (-10 PgC.yr$^{-1}$) following the decrease in $V_{cmax}$ (Fig. 9)

30 while it remains relatively stable over the tropics across all steps. Note that i) the study of

31 Welp et al. (2011) suggests a GPP around 150 PgC.yr$^{-1}$, similar to our estimate, based on

1    measurements of $^{18}O/^{16}O$ ratio in atmospheric $CO_2$ and ii) Koffi et al. (2012) found optimized

2    GPP of 146 PgC.yr$^{-1}$ from a CCDAS using the BETHY model.

3    Above-ground forest biomass

4    We analyze the impact of the optimization on the forest above-ground biomass at equilibrium

5    (i.e. after spin-up; see Fig. 12) as an example of the impact on model C stocks, and compare

6    the simulated values, for the same three latitude bands than above, to the estimate based on

7    field observations and remote sensing data. This product, which was produced in the

8    GEOCARBON project (and thus is referred to by the same name), integrates a pan-tropical

9    biomass map (Avitabile et al., 2016) with a boreal forest biomass product (Santoro et al.,

10   2015).

11   For the northern extra tropics, the prior above-ground C stock (~180 PgC) is reduced by the

12   optimization to 140 PgC, mainly through the decrease of the growing season length in step 1

13   with the assimilation of MODIS-NDVI. The significant decrease in GPP during step 1 (18 %)

14   led indeed to a similar decrease of the forest biomass (16%). Parameter changes through the

15   assimilation of FLUXNET and $CO_2$ data have a smaller impact (a change of less than 5 PgC).

16   These changes in the northern extra tropics bring the estimates by the ORCHIDEE model

17   closer to the satellite-based GEOCARBON product (~ 120 PgC).

18   For the tropics, while there is nearly no change with the assimilation of MODIS-NDVI in step

19   1, the use of FLUXNET data leads to a significant increase of the forest above ground

20   biomass (close to 25%). Such an increase does not correspond to an increase of the GPP (Fig.

21   10) but to changes in the autotrophic respiration parameter ($MR_b$) that lead to a decrease of

22   autotrophic respiration and an increase of NPP. The value does not change through step 3 and

23   remains significantly higher than the data-driven estimate. Note however that the lower value

24   in the GEOCARBON product could be partly due to the fact that we did not yet account for

25   land use effects in the CCDAS, such as deforestation in the Amazon.

26

27   **4   Discussion and conclusions**

28   In this paper we have described a first global Carbon Cycle Data Assimilation System that

29   assimilates three major carbon-cycle data streams, namely MODIS-NDVI observations of

30   vegetation activity at 60 sites, FLUXNET NEE and LE measurements at more than 70 sites,

31   and atmospheric $CO_2$ concentrations at 53 surface stations over three years in order to

Geoscientific
Model Development
Discussions

1    optimize the C cycle parameters of the ORCHIDEE process-based LSM (ORCHIDEE-

2    CCDAS). The study details the concept, the implementation and the main results of a

3    stepwise assimilation approach where the data streams have been assimilated in three

4    successive steps (including a propagation of the retrieved posterior parameter distributions

5    from one step to the next).

6    The assimilation of MODIS-NDVI (60 grid cell points, step 1) improved the phenology of

7    ORCHIDEE with a significant reduction of the growing season length and thus a direct

8    impact on the GPP. The results are similar to those presented in MacBean et al. (2015) who

9    describe the impact of such optimization on the global FAPAR simulations and the

10   improvement in the bias of the calculated leaf onset and senescence dates in more detail. The

11   optimization with FLUXNET data (78 sites, step 2) led to large improvements in the seasonal

12   cycle of the NEE and LE fluxes, constraining primarily the photosynthetic processes. Some

13   discrepancies remain due to site heterogeneity (i.e. different species and edaphic conditions)

14   that the model does not fully capture, and due to missing processes in the model (see Kuppel

15   et al. (2014) for a more thorough discussion). However, without the assimilation of

16   atmospheric $CO_2$ concentrations, the global (and continental) net carbon balance after step 2

17   was still clearly outside the admitted range (as reported by the GCP in Le Quéré et al. (2015),

18   which highlights the importance of assimilating a data stream such as this that provides

19   information at larger scales (constraining large scale respiration fluxes). The use of

20   atmospheric $CO_2$ concentration as an overall constraint in step 3 was technically challenging

21   as it required the coupling of ORCHIDEE with an atmospheric transport model in forward

22   and reverse mode (i.e. to compute the cost function and its gradients at each step of the

23   minimization process). As a result of the final step, we were able to fit the atmospheric $CO_2$

24   growth rate and thus to derive a land C sink compatible with current best estimates from the

25   GCP. The assimilation of $CO_2$ data also slightly changed the seasonality of the NEE, which

26   improved the fit to the atmospheric $CO_2$ seasonal cycle. Note that assimilating only $CO_2$ data

27   would lead to a similar global land C sink but with a different model parameter set not

28   compatible with the information provided by MODIS-NDVI and FLUXNET data.

29   The consistency of the stepwise approach has been evaluated with back-compatibility checks

30   after the final step (step 3: assimilation of atmospheric $CO_2$ concentration). The optimized

31   model with the final set of parameters does not degrade the fit to MODIS-NDVI and

32   FLUXNET data that were assimilated in the first two steps (only minor changes of the

1    RMSEs occur; see Fig. 8). This result has two important consequences. Most importantly it

2    suggests that current state of the art LSMs (at least ORCHIDEE) have reached a level of

3    development where consistent assimilation of multiple data streams is finally possible. This

4    overcomes the most important limitation noted by Rayner (2010) to the widespread use of

5    CCDAS systems. At a more technical level it suggests that stepwise assimilation is a valid

6    and feasible approach. Although we only carried the estimated parameter uncertainties from

7    one step to the next (as a first simple approach), and not the full error variance-covariance

8    matrix, we were able to propagate enough information to maintain an optimal model-data fit

9    after the last step for the three data streams (see MacBean et al. (2016) for a more specific

10   analysis of this issue). However, not propagating the covariance terms may have a larger

11   impact for the reduction of the inferred parameter uncertainties (see for instance the large

12   parameter / flux error reduction in Fig. 9 / Fig. 11). The order of the different steps was

13   dictated by the number of parameters we choose to expose to each data stream, from only a

14   few phenology parameters for NDVI up to the largest set for atmospheric $CO_2$. Recall that

15   under the fundamental theory the order of assimilation is unimportant. Testing whether our

16   system meets this criterion is an important check on the robustness of the method but is not

17   technically feasible with the full-blown system; it is currently being tested with some smaller

18   models.

19   Most of the optimized parameter values have significantly changed compared to their prior

20   values, with a large error reduction for most (Fig. 9) that results in a strong constraint on the

21   simulated fluxes (Fig. 11). In the last step, the assimilation of atmospheric $CO_2$ data mainly

22   led to the optimization of respiration-related parameters, especially the regional soil carbon

23   multipliers ($K_{soilC,reg}$). Note that this was also the case for the BETHY-CCDAS, as described

24   in Rayner et al. (2005) (see their Table 2). This is linked to the difficult issue of representing

25   the effects of historical changes in land cover and land management as well as soil texture

26   effects on soil carbon dynamics, and the necessary choice of a standard spin-up procedure to

27   account for these effects. Ideally one would need to perform the optimization of the model

28   over a long historical period with LULCC and land management practices included and the

29   optimization of related parameters. However, this is not currently feasible at global scale and

30   uncertainties in the forcing would introduce as much difficulty as uncertainties in the initial

31   condition. The adjustment of the initial C pool contents is thus a logical compromise and

32   further investigations into the impact of the selected set-up (number of regions for $K_{soilC,reg}$,

33   their associated uncertainties) on the C fluxes simulated in the future are needed. Note that a

Geoscientific
Model Development
Discussions

1 first improvement would be to include LULCC in the transient simulation (to define the initial

2 state) before the assimilation period.

3 Nonetheless, several limitations, inherent to the optimization of model parameters in a

4 CCDAS, need to be called to mind when evaluating these results (see also Rayner et al.,

5 2010). First, the structure of the land surface model (i.e. how biogeochemical processes are

6 represented) is critical. Any missing/misrepresented processes may have a direct impact and

7 thus lead to biases in the selected parameters. Note that this limitation could be even more

8 severe when using atmospheric $CO_2$ measurements, as these data provide a direct constraint

9 on the overall net C exchange between the atmosphere and the vegetation, thus including all

10 processes. As an example, the model sensitivity to atmospheric $CO_2$ increase (e.g. through the

11 parameters $V_{cmax}$ and $G_{s,slope}$) could be non optimal as the current model version does not

12 include explicit nitrogen and phosphorus limitations on photosynthesis. Second, the chosen

13 set of observations does not provide specific constraints on long term C processes such as tree

14 mortality, disturbance effects, or C allocation within a plant. For instance Fig. 12 illustrates

15 that the optimized model may still significantly overestimate tropical forest biomass. The

16 assimilation of above-ground biomass or soil carbon stock observations (i.e. site-level

17 measurements or regional estimates) should thus provide critical complementary information

18 (see Thum et al., in revision for AFM).

19 To conclude, this work is a step forward in terms of multiple data streams assimilation that

20 opens new perspectives for a better understanding of the carbon cycle and better predictions

21 of the fate of the land carbon sink in the 21st century as a consequence of anthropogenic

22 changes. As ORCHIDEE is part of the IPSL earth system model the impact of the

23 optimization on future climate change predictions will be assessed in a future study. However,

24 we first need to run the ORCHIDAS with a longer atmospheric $CO_2$ record (i.e. several

25 decades) in order to provide stronger constraints on parameters controlling the impact of

26 climate extremes on the net carbon fluxes at continental to global scales, and the sensitivity of

27 photosynthesis to increasing $CO_2$ concentration. The optimized model will allow a more in-

28 depth investigation of the trend and inter-annual variations of land surface C fluxes at

29 continental to regional scales, as well as their driving mechanisms. It will offer a more

30 reliable and robust process-based diagnostic of the land C cycle that is compatible with

31 current major data streams. Overall, we have illustrated the benefit of combining multiple

32 data streams in a process-based model to optimize different processes of the model, related to

Geoscientific
Model Development
Discussions

1 different temporal and spatial scales. The optimization will be updated regularly as new

2 processes are integrated into the ORCHIDEE model, such as for instance land management

3 (Naudts et al., 2015).

4

5 **Code availability**

6 The ORCHIDEE model code and the run environment are open source

7 (http://forge.ipsl.jussieu.fr/orchidee) and the associated documentation can be found at

8 https://forge.ipsl.jussieu.fr/orchidee/wiki/Documentation. Note that the tangent linear version

9 of the ORCHIDEE model has been generated using commercial software (TAF;

10 http://www.fastopt.com/products/taf/taf.shtml). For this reason, only the "forward" version of

11 the ORCHIDEE model is available for sharing. The optimization scheme (in Python) is

12 available through a dedicated web site for data assimilation with ORCHIDEE

13 (http://orchidas.lsce.ipsl.fr/). Nevertheless readers interested in running ORCHIDEE are

14 encouraged to contact the corresponding author for full details and latest bug fixes. Finally,

15 the source code of the LMDZ atmospheric transport model can be found at

16 http://web.lmd.jussieu.fr/trac.

17

18 # Appendix

19 **A1. Ocean fluxes**

20 Figure A1 displays the air-sea fluxes from the statistical model.

21 **A2. Fire fluxes**

22 In order to account for fundamental differences between six fire flux categories provided by

23 the GFED product, we grouped these emissions into 3 types with specific treatments. The first

24 group includes C emissions from deforestation and peat fires, which are considered to be

25 permanent carbon lost to the atmosphere, at least over the considered time scales. These

26 fluxes are rescaled to an annual emission of 1.1 $PgC.yr^{-1}$ globally following typical values

27 reported in the literature for deforestation (Houghton R., 2003). The second group consists of

28 C emissions from agriculture and savannah fires, which are compensated by a C sink during

29 the regrowth of these biomes (i.e., savannah and some type of plants on the farmland). These

1   effects are not completely accounted for in ORCHIDEE as the model does not simulate

2   savannah and agriculture fire. Hence, the emissions over the whole period and for each pixel

3   become zero, but their seasonal variations are used. The final group includes emissions from

4   woodland and burnt forests. We considered that at steady state and for a given region certain

5   forests burn but that nearby forests are re-growing following older fires. We thus imposed

6   regrowth at the region scale given that the ORCHIDEE model version that we use does not

7   account for such regrowth. The main assumption is that over century time scale the

8   forest/woodland system is at steady state over a given region (few thousand square km), i.e.

9   there is no net deforestation. We selected an ensemble of small regions over which we

10   calculated the regrowth of these biomes. The derived emissions over the whole period and for

11   each region thus become zero; though we include their spatial and temporal variations. The

12   overall biomass burning flux considered in the CCDAS for the optimization process is the

13   sum of the three fluxes as described above.

14   **A3. Multipliers of the soil initial carbon pools**

15   Figure A2 provides the optimized values of the $K_{soilC,reg}$ parameters that optimize the initial

16   soil carbon pool sizes.

17   .

18   **Acknowledgements**

9

10

11

12

13

14

15

16

# 1 References

2 Alton, P. B.: From site-level to global simulation: Reconciling carbon, water and energy
3 fluxes over different spatial scales using a process-based ecophysiological land-surface
4 model, Agric. For. Meteorol., 176, 111–124, doi:10.1016/j.agrformet.2013.03.010, 2013.

5 Avitabile, V., Herold, M., Heuvelink, G. B. M., Lewis, S. L., Phillips, O. L., Asner, G. P.,
6 Armston, J., Ashton, P. S., Banin, L., Bayol, N., Berry, N. J., Boeckx, P., de Jong, B. H. J.,
7 DeVries, B., Girardin, C. A. J., Kearsley, E., Lindsell, J. A., Lopez-Gonzalez, G., Lucas, R.,
8 Malhi, Y., Morel, A., Mitchard, E. T. A., Nagy, L., Qie, L., Quinones, M. J., Ryan, C. M.,
9 Ferry, S. J. W., Sunderland, T., Laurin, G. V., Gatti, R. C., Valentini, R., Verbeeck, H.,
10 Wijaya, A. and Willcock, S. (2016), An integrated pan-tropical biomass map using multiple
11 reference datasets. Glob Change Biol. doi:10.1111/gcb.13139.

12 Bacour, C., Peylin, P., MacBean, N., Rayner, P. J., Delage, F., Chevallier, F., Weiss, M.,
13 Demarty, J., Santaren, D., Baret, F., Berveiller, D., Dufrêne, E. and Prunet, P.: Joint
14 assimilation of eddy covariance flux measurements and FAPAR products over temperate
15 forests within a process-oriented biosphere model, J. Geophys. Res. Biogeosciences, 120, 1–
16 19, doi:10.1002/2015JG002966, 2015.

17 Bousquet P., D. Hauglustaine, P. Peylin, C. Carouge, and P. Ciais, Two decades of OH
18 variability as inferred by an inversion of atmospheric transport and chemistry of methyl
19 chloroform, *Atmos. Chem. and Phys.*, 5, 263-2656, ISI:000232370800002, 2005.

20 Braswell, B. H., Sacks, W. J., Linder, E. and Schimel, D. S.: Estimating diurnal to annual
21 ecosystem parameters by synthesis of a carbon flux model with eddy covariance net
22 ecosystem exchange observations, Glob. Change Biol., 11(2), 335–355, doi:10.1111/j.1365-
23 2486.2005.00897.x, 2005.

24 Byrd, R. H., P. Lu, J. Nocedal, and C. Zhu (1995), A limited memory algorithm for bound
25 constrained optimization, SIAM J. Sci. Stat. Comput., 16(5), 1190–1208.

26 Canadell, J. G., Ciais, P., Sabine, C., and Joos, F. (Eds.): REgional Carbon Cycle Assessment
27 and Processes (RECCAP), Special issue, Biogeosciences, http://www.biogeosciences-
28 discuss.net/ special_issue83.html, 2013.

29 Chevallier F., Fisher M., Peylin P., Serrar S., Bousquet P., Bréon F-M., Chédin A., Ciais P.
30 (2005), Inferring CO2 sources and sinks from satellite observations: Method and application
31 to TOVS data, *Journal of Geophysical Research*, **110**, D24309, doi:20.1029/2005JD006390,
32 13pp.

33 Chevallier, F., P. Ciais, T. J. Conway, T. Aalto, B. E. Anderson, P. Bousquet, E. G. Brunke,
34 L. Ciattaglia, Y. Esaki, M. Fröhlich, A.J. Gomez, A.J. Gomez-Pelaez, L. Haszpra, P.
35 Krummel, R. Langenfelds, M. Leuenberger, T. Machida, F. Maignan, H. Matsueda, J. A.
36 Morguí, H. Mukai, T. Nakazawa, P. Peylin, M. Ramonet, L. Rivier, Y. Sawa, M. Schmidt, P.
37 Steele, S. A. Vay, A. T. Vermeulen, S. Wofsy, D. Worthy, (2010), CO2 surface fluxes at grid
38 point scale estimated from a global 21-year reanalysis of atmospheric measurements, *Journal
39 of Geophysical Research*, **115**, D21307, doi:10.1029/2010JD013887.

Geoscientific
Model Development
Discussions

1   Collatz GJ, Ribas-Carbo M, Berry JA, (1992), Coupled Photosynthesis-Stomatal Conductance
2   Model for Leaves of C4 Plants. *Aust J Plant Physiol*, **19**: 519-38.

3   de Rosnay P, Polcher J. (1998), Modelling root water uptake in a complex land surface
4   scheme coupled to a GCM. *Hydrol Earth Syst Sc*, **2**: 239-55.

5   Dee, D. P., Uppala, S. M., Simmons, A. J., Berrisford, P., Poli, P., Kobayashi, S., Andrae, U.,
6   Balmaseda, M. A., Balsamo, G., Bauer, P., Bechtold, P., Beljaars, A. C. M., van de Berg, L.,
7   Bidlot, J., Bormann, N., Delsol, C., Dragani, R., Fuentes, M., Geer, A. J., Haimberger, L.,
8   Healy, S. B., Hersbach, H., Hólm, E. V., Isaksen, L., Kållberg, P., Köhler, M., Matricardi, M.,
9   Mcnally, A. P., Monge-Sanz, B. M., Morcrette, J. J., Park, B. K., Peubey, C., de Rosnay, P.,
10  Tavolato, C., Thépaut, J. N., and Vitart, F.: The ERA-interim reanalysis: configuration and
11  performance of the data assimilation system, Q. J. Roy. Meteor. Soc., 137, 553–597,
12  doi:10.1002/qj.828, 2011.

13  Ducoudre NI, Laval K, Perrier A. Sechiba (1993), A New Set of Parameterizations of the
14  Hydrologic Exchanges at the Land Atmosphere Interface within the Lmd Atmospheric
15  General-Circulation Model. *J Climate*, **6**: 248-73.

16  Dufresne, J.-L., et al. (2011), Climate change projections using the IPSL-CM5 Earth System
17  Model: from CMIP3 to CMIP5, submitted to Clim. Dynam.

18  Farquhar, G. D., von Caemmerer, S. von and Berry, J. A.: A biochemical model of
19  photosynthetic CO2 assimilation in leaves of C3 species, Planta, 149(1), 78–90, 1980.

20  Folberth G., Hauglustaine D.A., Ciais P., et al. (2005), On the role of atmospheric chemistry
21  in the global $CO_2$ budget, *Geophysical Research Letters*, **32**(8): L08801

22  Fung, I. Y., C. J. Tucker, and K. C. Prentice, Application of Advanced Very High Resolution
23  Radiometer vegetation index to study atmosphere – biosphere exchange of $CO_2$, J. Geophys.
24  Res., 92, 2999– 3015, 1987.

25  GLOBALVIEW : Cooperative Global Atmospheric Data Integration Project. 2013, updated
26  annually. Multi-laboratory compilation of synchronized and gap-filled atmospheric carbon
27  dioxide records for the period 1979-2012 (obspack_co2_1_GLOBALVIEW-
28  CO2_2013_v1.0.4_2013-12-23). Compiled by NOAA Global Monitoring Division: Boulder,
29  Colorado, U.S.A. Data product accessed at http://dx.doi.org/10.3334/OBSPACK/1002.

30  Groenendijk, M., Dolman, a. J., van der Molen, M. K., Leuning, R., Arneth, a., Delpierre,
31  N., Gash, J. H. C., Lindroth, a., Richardson, a. D., Verbeeck, H. and Wohlfahrt, G.:
32  Assessing parameter variability in a photosynthesis model within and between plant
33  functional types using global Fluxnet eddy covariance data, Agric. For. Meteorol., 151(1),
34  22–38, doi:10.1016/j.agrformet.2010.08.013, 2011.

35  Hauglustaine D.A., Hourdin F., Jourdain L., et al. (2004), Interactive chemistry in the
36  Laboratoire de Meteorologie Dynamique general circulation model: Description and
37  background tropospheric chemistry evaluation, *Journal of Geophysical Research -*
38  *Atmosphere*, **109**(D4): D04314

1   Houghton, R. A. (2003) Revised estimates of the annual net flux of carbon to the atmosphere
2   from changes in land use and land management 1850-2000. *Tellus* **55B**: 378-390.

3   Hourdin F. and Armengaud A. (1999), The use of finite-volume methods for atmospheric
4   advection of trace species. Part I: Test of various formulations in a general circulation
5   model, *Monthly Weather Review*, **127**(5): 822-837

6   Hourdin, F. et al. (2006), The LMDZ4 general circulation model: climate performance and
7   sensitivity to parametrized physics with emphasis on tropical convection, *Climate
8   Dynamics*, **27**(7), 787-813, 2006.

9   IPCC, (2007). Climate Change 2007: The Physical Science Basis. Contribution of Working
10  Group I to the Fourth Assessment Report of the Intergovernmental Panel on Climate Change
11  [Solomon, S., D. Qin, M. Manning (eds.)]

12  Kaminski, T., Knorr, W., Scholze, M., Gobron, N., Pinty, B., Giering, R. and Mathieu, P-P.
13  (2012), Consistent assimilation of MERIS FAPAR and atmospheric $CO_2$ into a terrestrial
14  vegetation model and interactive mission benefit analysis, *Biogeosciences*, **9**, 3173-3184.

15  Kaminski, T., Knorr, W., Schürmann, G., Scholze, M., Rayner, P. J., Zaehle, S., Blessing, S.,
16  Dorigo, W., Gayler, V., Giering, R., Gobron, N., Grant, J. P., Heimann, M., Hooker-Stroud,
17  a., Houweling, S., Kato, T., Kattge, J., Kelley, D., Kemp, S., Koffi, E. N., Köstler, C.,
18  Mathieu, P. P., Pinty, B., Reick, C. H., Rödenbeck, C., Schnur, R., Scipal, K., Sebald, C.,
19  Stacke, T., Van Scheltinga, a. T., Vossbeck, M., Widmann, H. and Ziehn, T.: The
20  BETHY/JSBACH Carbon Cycle Data Assimilation System: Experiences and challenges, J.
21  Geophys. Res. Biogeosciences, 118(4), 1414–1426, doi:10.1002/jgrg.20118, 2013.

22  Kato, T., Knorr, W., Scholze, M., Veenendaal, E., Kaminski, T., Kattge, J. and Gobron, N.:
23  Simultaneous assimilation of satellite and eddy covariance data for improving terrestrial water
24  and carbon simulations at a semi-arid woodland site in Botswana, Biogeosciences, 10(2),
25  789–802, doi:10.5194/bg-10-789-2013, 2013.

26  Keenan, T. F., Davidson, E. a., Munger, J. W. and Richardson, A. D.: Rate my data:
27  Quantifying the value of ecological data for the development of models of the terrestrial
28  carbon cycle, Ecol. Appl., 23(1), 273–286, doi:10.1890/12-0747.1, 2013.

29  Keenan, T. F., Davidson, E., Moffat, A. M., Munger, W. and Richardson, A. D.: Using
30  model-data fusion to interpret past trends, and quantify uncertainties in future projections, of
31  terrestrial ecosystem carbon cycling, Glob. Change Biol., 18(8), 2555–2569,
32  doi:10.1111/j.1365-2486.2012.02684.x, 2012.

33  Knorr, W. and Kattge, J.: Inversion of terrestrial ecosystem model parameter values against
34  eddy covariance measurements by Monte Carlo sampling, Glob. Change Biol., 11(8), 1333–
35  1351, doi:10.1111/j.1365-2486.2005.00977.x, 2005.

36  Knyazikhin, Y., Martonchik, J.V., Myneni, R.B., Diner, D.J., and Running, S.W. (1998),
37  Synergistic algorithm for estimating vegetation canopy leaf area index and fraction of
38  absorbed photosynthetically active radiation from MODIS and MISR, *Journal of Geophysical
39  Research*, **103**, D24, 32,257-32,276.

Geoscientific
Model Development
Discussions

1  Koffi, E. N., Rayner, P. J., Scholze, M., & Beer, C. (2012). Atmospheric constraints on gross
2  primary productivity and net ecosystem productivity: Results from a carbon  cycle data
3  assimilation system. *Global biogeochemical cycles*, *26*(1).

4  Krinner, G., Viovy, N., de Noblet-Ducoudré, N., Ogée, J., Polcher, J., Friedlingstein, P.,
5  Ciais, P., Sitch, S. and Prentice, I. C.: A dynamic global vegetation model for studies of the
6  coupled atmosphere-biosphere system, Glob. Biogeochem. Cycles, 19(1) [online] Available
7  from: http://onlinelibrary.wiley.com/doi/10.1029/2003GB002199/pdf (Accessed 23
8  November 2015), 2005.

9  Kuppel, S., Chevallier, F. and Peylin, P.: Quantifying the model structural error in carbon
10  cycle data assimilation systems, Geosci. Model Dev., 6(1), 45–55, doi:10.5194/gmd-6-45-
11  2013, 2013.

12  Kuppel, S., Peylin, P., Chevallier, F., Bacour, C., Maignan, F. and Richardson, A. D.:
13  Constraining a global ecosystem model with multi-site eddy-covariance data, Biogeosciences,
14  9(10), 3757–3776, doi:10.5194/bg-9-3757-2012, 2012.

15  Kuppel, S., Peylin, P., Maignan, F., Chevallier, F., Kiely, G., Montagnani, L. and Cescatti, A.:
16  Model–data fusion across ecosystems: from multisite optimizations to global simulations,
17  Geosci. Model Dev., 7(6), 2581–2597, 2014.

18  Lasslop, G., M. Reichstein, D. Papale, A. D. Richardson, A. Arneth, A. Barr, P. Stoy, and G.
19  Wohlfahrt. 2010. Separation of net ecosystem exchange into assimilation and respiration
20  using a light response curve approach: critical issues and global evaluation. *Global Change*
21  *Biology*, **16**, 187-208.

22  Lasslop, G., M. Reichstein, J. Kattge, and D. Papale. 2008. Influences of observation errors in
23  eddy flux data on inverse model parameter estimation. *Biogeosciences*, **5**, 1311-1324.

24  Le Quéré, C., Moriarty, R., Andrew, R. M., Peters, G. P., Ciais, P., Friedlingstein, P., Jones,
25  S. D., Sitch, S., Tans, P., Arneth,  a., Boden, T. a., Bopp, L., Bozec, Y., Canadell, J. G., Chini,
26  L. P., Chevallier, F., Cosca, C. E., Harris, I., Hoppema, M., Houghton, R. a., House, J. I., Jain,
27  a. K., Johannessen, T., Kato, E., Keeling, R. F., Kitidis, V., Klein Goldewijk, K., Koven, C.,
28  Landa, C. S., Landschützer, P., Lenton,  a., Lima, I. D., Marland, G., Mathis, J. T., Metzl, N.,
29  Nojiri, Y., Olsen,  a., Ono, T., Peng, S., Peters, W., Pfeil, B., Poulter, B., Raupach, M. R.,
30  Regnier, P., Rödenbeck, C., Saito, S., Salisbury, J. E., Schuster, U., Schwinger, J., Séférian,
31  R., Segschneider, J., Steinhoff, T., Stocker, B. D., Sutton,  a. J., Takahashi, T., Tilbrook, B.,
32  van der Werf, G. R., Viovy, N., Wang, Y.-P., Wanninkhof, R., Wiltshire,  a. and Zeng, N.:
33  Global carbon budget 2014, Earth Syst. Sci. Data, 7(1), 47–85, doi:10.5194/essd-7-47-2015,
34  2015.

35  Liss, P. and Merlivat, L. (1986). The role of sea-air exchange in geochemical cycling, Ed. P.
36  Menard, chapter Air-sea gas exchange rates: Introduction and synthesis, pages 113-127.
37  Reidel, Dordrecht.

38  MacBean, N., Maignan, F., Peylin, P., Bacour, C., François-Marie, B. and Ciais, P.: Using
39  satellite data to improve the leaf phenology of a global Terrestrial Biosphere Model,
40  Biogeosciences, 12, 7185-7208, 2015.

Geoscientific
Model Development
Discussions

1  MacBean, N., Peylin, P., Chevallier, F.,  Scholze, M. and G. Schürmann, Consistent
2  assimilation of multiple data streams in a Carbon Cycle Data Assimilation System,
3  Biogeosciences, submitted, 2016.

4  Madec, G., P. Delecluse, M. Imbard and C. Lévy, 1998 : OPA 8.1 Ocean General Circulation
5  Model reference manual, Note du Pole de Modelisation, Institut Pierre-Simon Laplace, 11,
6  91pp.

7  Maignan, F., Bréon, F.-M., Chevallier, F., Viovy, N., Ciais, P., Garrec, C., Trules, J., and
8  Mancip, M.: Evaluation of a Global Vegetation Model using time series of satellite vegetation
9  indices, Geosci. Model Dev., 4, 1103-1114, doi:10.5194/gmd-4-1103-2011, 2011.

10  Moore, D. J. P., Hu, J., Sacks, W. J., Schimel, D. S. and Monson, R. K.: Estimating
11  transpiration and the sensitivity of carbon uptake to water availability in a subalpine forest
12  using a simple ecosystem process model informed by measured net CO2 and H2O fluxes,
13  Agric. For. Meteorol., 148(10), 1467–1477, doi:10.1016/j.agrformet.2008.04.013, 2008.

14  Naudts, K., Ryder, J., J McGrath, M., Otto, J., Chen, Y., Valade, A., ... & Ghattas, J. (2015).
15  A vertically discretised canopy description for ORCHIDEE (SVN r2290) and the
16  modifications to the energy, water and carbon fluxes. *Geoscientific Model Development,* 8,
17  2035-2065.

18  Nightingale, P.D., et al. 2000. In situ evaluation of air-sea gas exchange parameterizations
19  using novel conservative and volatile tracers. *Glob. Biogeochem Cycles*, **14**, 373-387.

20  Olson, J., Watts, J.A., and Allison, L.J. (1983), Carbon in Live Vegetation of Major World
21  Ecosystems, ORNL-5862, Oak Ridge National Laboratory, Oak Ridge, Tennessee, 164pp.

22  Papale, D. (2006), Towards a standardized processing of Net Ecosystem Exchange measured
23  with eddy covariance technique: algorithms and uncertainty estimation, [online] Available
24  from: http://dspace.unitus.it/handle/2067/1321 (Accessed 27 August 2013)

25  Parton W, Stewart J, Cole C, (1988), Dynamics of C, N, P and S in grassland soils: a model.
26  *Biogeochemistry*, **5**: 109-31.

27  Peylin P, Rayner PJ, Bousquet P, et al. (2005), Daily $CO_2$ flux estimates over Europe from
28  continuous atmospheric measurements: 1, inverse methodology, *Atmospheric Chemistry and*
29  *Physics*, **5**: 3173-3186.

30  Piao, S., Sitch, S., Ciais, P., Friedlingstein, P., Peylin, P., Wang, X., Ahlström, A., Anav, A.,
31  Canadell, J. G., Cong, N., Huntingford, C., Jung, M., Levis, S., Levy, P. E., Li, J., Lin, X.,
32  Lomas, M. R., Lu, M., Luo, Y., Ma, Y., Myneni, R. B., Poulter, B., Sun, Z., Wang, T., Viovy,
33  N., Zaehle, S. and Zeng, N.: Evaluation of terrestrial carbon cycle models for their response to
34  climate variability and to CO2 trends, Glob. Change Biol., 19(7), 2117–2132,
35  doi:10.1111/gcb.12187, 2013.

36  Prentice, I. C., Liang, X., Medlyn, B. E. and Wang, Y.-P.: Reliable, robust and realistic: the
37  three R's of next-generation land-surface modelling, Atmospheric Chem. Phys., 15(10),
38  5987–6005, doi:10.5194/acp-15-5987-2015, 2015.

Geoscientific
Model Development
Discussions

1   Rayner, P. J., Scholze, M., Knorr, W., Kaminski, T., Giering, R. and Widmann, H.: Two
2   decades of terrestrial carbon fluxes from a carbon cycle data assimilation system ( CCDAS ), ,
3   19, doi:10.1029/2004GB002254, 2005.

4   Rayner, P. J. (2010). The current state of carbon-cycle data assimilation. *Current Opinion in*
5   *Environmental Sustainability*, *2*(4), 289-296.

6   Regnier, P., Friedlingstein, P., Ciais, P., Mackenzie, F. T., Gruber, N., Janssens, I. A., ... &
7   Arndt, S. (2013). Anthropogenic perturbation of the carbon fluxes from land to ocean. *Nature*
8   *Geoscience*, *6*(8), 597-607.

9   Ricciuto, D. M., A. W. King, D. Dragoni, and W. M. Post (2011), Parameter and prediction
10  uncertainty in an optimized terrestrial carbon cycle model: Effects of constraining variables
11  and data record length, J. Geophys. Res., 116, G01033, doi:10.1029/2010JG001400.

12  Ricciuto, D. M., Butler, M. P., Davis, K. J., Cook, B. D., Bakwin, P. S., Andrews, A. and
13  Teclaw, R. M.: Causes of interannual variability in ecosystem-atmosphere CO2 exchange in a
14  northern Wisconsin forest using a Bayesian model calibration, Agric. For. Meteorol., 148(2),
15  309–327, doi:10.1016/j.agrformet.2007.08.007, 2008.

16  Richardson, A. D., Williams, M., Hollinger, D. Y., Moore, D. J. P., Dail, D. B., Davidson, E.
17  a., Scott, N. a., Evans, R. S., Hughes, H., Lee, J. T., Rodrigues, C. and Savage, K.: Estimating
18  parameters of a forest ecosystem C model with measurements of stocks and fluxes as joint
19  constraints, Oecologia, 164(1), 25–40, doi:10.1007/s00442-010-1628-y, 2010.

20  Rivier L., Ciais P., Hauglustaine D.A., et al. (2006), Evaluation of $SF_6$, $C_2Cl_4$, and CO to
21  approximate fossil fuel $CO_2$ in the Northern Hemisphere using a chemistry transport
22  model, *Journal Geophysical Research-Atmosphere*, **111**(D16) - D16311.

23  Rödenbeck, C., T. J. Conway, and R. L. Langenfelds (2006), The effect of systematic
24  measurement errors on atmospheric CO2 inversions: A quantitative assessment, Atmos.
25  *Chem. Phys.*, **6**, 149–161, doi:10.5194/ acp-6-149-2006.

26  Rödenbeck C., D.C.E. Bakker, N. Gruber, Y. Iida, A. Jacobson, S. Jones, P. Landschutzer, N.
27  Metzl, S. Nakaoka, A. Olsen, G.-H. Park, P. Peylin, K.B. Rodgers, T.P. Sasse, U. Schuster,
28  J.D. Shutler, V. Valsala, R. Wanninkhof, and J. Zeng, 2015. Data-based estimates of the
29  ocean carbon sink variability – First results of the Surface Ocean pCO2 Mapping
30  intercomparison (SOCOM). Biogeosciences, 12: 7251-7278. doi:10.5194/bg-12-7251-2015.

31  Rosenblatt, F. (1958). The perceptron: a probabilistic model for information storage and
32  organization in the brain. *Psychological review*, *65*(6), 386.

33  Ruimy A, Dedieu G, Saugier B, (1996), TURC: A diagnostic model of continental gross
34  primary productivity and net primary productivity. *Global Biogeochemical Cycles,* **10**: 269-
35  85.

36  Sabine, C.L., et al. 2004. The oceanic sink for anthropogenic $CO_2$. *Science*, **305** (5682), 367-
37  371.

1   Santaren, D., Peylin, P., Bacour, C., Ciais, P. and Longdoz, B.: Ecosystem model
2   optimization using in situ flux observations: benefit of Monte Carlo versus variational
3   schemes and analyses of the year-to-year model performances, Biogeosciences, 11(24), 7137–
4   7158, 2014.

5   Santoro M, Beaudoin A, Beer C, Cartus O, Fransson JE, Hall RJ et al. (2015). Forest growing
6   stock volume of the northern hemisphere : Spatially explicit estimates for 2010 derived from
7   Envisat ASAR, *Remote Sensing of Environment*, **168**, 316-334.

8   Sitch, S., Friedlingstein, P., Gruber, N., Jones, S. D., Murray-Tortarolo, G., Ahlström, A.,
9   Doney, S. C., Graven, H., Heinze, C., Huntingford, C., Levis, S., Levy, P. E., Lomas, M.,
10  Poulter, B., Viovy, N., Zaehle, S., Zeng, N., Arneth, A., Bonan, G., Bopp, L., Canadell, J. G.,
11  Chevallier, F., Ciais, P., Ellis, R., Gloor, M., Peylin, P., Piao, S., Le Quéré, C., Smith, B.,
12  Zhu, Z. and Myneni, R.: Recent trends and drivers of regional sources and sinks of carbon
13  dioxide, Biogeosciences, 12, 653–679, doi:10.5194/bgd-12-653-2015, 2015.

14  Takahashi, et al. 2009, Corrigendum to "Climatological mean and decadal change in surface
15  ocean $pCO_2$, and net sea-air $CO_2$flux over the global oceans" *Deep Sea Res. II*, **56**, 554-577.

16  Tarantola A. (1987), Inverse problem theory: Methods for data fitting and parameter
17  estimation. Elsevier, Amsterdam.

18  Tarantola, A. (2005), Inverse problem theory and methods for model parameters
19  estimation, *Society for Industrial and Applied Mathematics*, Philadelphia, ISBN 0-89871-572-
20  5.

21  Thum, T., MacBean, N., Peylin, P., Bacour, C., Santaren, D., Longdoz, B., Loustau, D. and
22  Ciais, P., The potential of forest biomass data in addition to carbon and water flux
23  measurements to constrain ecosystem model parameters: Case studies at two temperate forest
24  sites, Agriculture and Forest Meteorology, in revision, 2015.

25  Tiedtke M. (1989), A comprehensive mass flux scheme for cumulus parameterization in
26  large-scale models, *Monthly Weather Review*, **117**(8): 1779-1800

27  Tucker, C. J. (1979). Red and photographic infrared linear combinations for monitoring
28  vegetation. *Remote Sensing of Environment*, **8**, 127-150.

29  Twine, T. E., W. P. Kustas, J. M. Norman, D. R. Cook, P. R. Houser, T. P. Meyers, J. H.
30  Prueger, P. J. Starks, and M. L. Wesely (2000), Correcting eddy-covariance flux
31  underestimates over a grassland, *Agric. For. Meteorol.*, **103**(3), 279–300, doi:10.1016/S0168-
32  1923(00)00123-4.

33  Verant, S., Laval, K., Polcher, J., and De Castro, M. (2004), Sensitivity of the continental
34  hydrological cycle to the spatial resolution over the Iberian Peninsula, Journal of
35  *Hydrometeorology*, **5**, 267-285.

36  Vermote, E., C.O. Justice and F-M Breon (2009), Towards a generalized approach for
37  correction of the BRDF effect in MODIS directional reflectances, *IEEE Transactions on
38  Geoscience and Remote Sensing*, **47**, 3, 898-908.

1   Wang, Y. P., Baldocchi, D., Leuning, R., Falge, E. and Vesala, T.: Estimating parameters in a
2   land-surface model by applying nonlinear inversion to eddy covariance flux measurements
3   from eight FLUXNET sites, Glob. Change Biol., 13(3), 652–670, doi:10.1111/j.1365-
4   2486.2006.01225.x, 2007.

5   Wang, Y. P., Leuning, R., Cleugh, H. and Coppin, P.: Parameter estimation in surface
6   exchange models using nonlinear inversion : how many parameters can we estimate and
7   which measurements are most useful ?, Glob. Change Biol., 7, 495–510, doi:10.1046/j.1365-
8   2486.2001.00434.x, 2001.

9   Wanninkhof, R., 1992. Relationship between wind speed and gas exchange. *J. Geophys.*
10  *Res.* **97**, 7373-7382.

11  Weiss, R.F., 1974. Carbon dioxide in water and seawater: the solubility of a non-ideal gas.
12  Mar. Chem., **2**, 203-215.

13  Welp, L. R., Keeling, R. F., Meijer, H. A., Bollenbacher, A. F., Piper, S. C., Yoshimura, K.,
14  ... & Wahlen, M. (2011). Interannual variability in the oxygen isotopes of atmospheric CO2
15  driven by El Nino. *Nature*, *477*(7366), 579-582.

16  Williams, M., Schwarz, P. a, Law, B. E., Irvine, J. and Kurpius, M. R.: An improved analysis
17  of forest carbon dynamics using data assimilation, Glob. Change Biol., 11(1), 89–105,
18  doi:10.1111/j.1365-2486.2004.00891.x, 2005.

19  Xiao, J., Davis, K. J., Urban, N. M. and Keller, K.: Uncertainty in model parameters and
20  regional carbon fluxes: A model-data fusion approach, Agric. For. Meteorol., 189-190, 175–
21  186, doi:10.1016/j.agrformet.2014.01.022, 2014.

22  Zobitz, J. M., D. J. P. Moore, T. Quaife, B. H. Braswell, A. Bergeson, J. A. Anthony, and R.
23  K. Monson (2014), Joint data assimilation of satellite reflectance and net ecosystem exchange
24  data constrains ecosystem carbon fluxes at a high-elevation subalpine forest, Agric. For.
25  Meteorol., 195–196, 73–88.

26  Zobler, L (1986), A world soil file for global climate modeling, NASA Technical
27  Memorandum 87802. NASA Goddard Institute for Space Studies, New York, U.S.A.

28

29

Geosci. Model Dev. Discuss., doi:10.5194/gmd-2016-13, 2016
Manuscript under review for journal Geosci. Model Dev.
Published: 28 January 2016
© Author(s) 2016. CC-BY 3.0 License.

1 **Tables**

2  Table 1. Parameters description, generality (PFT dependent, global, specific to FLUXNET

3  sites or for a set of regions) and data stream(s) that were used to constrain them.

| Parameter | Description | Dependent | Constraint |
|---|---|---|---|
| $V_{cmax}$ | Maximum carboxylation rate ($\mu mol \cdot m^{-2} \cdot s^{-1}$) | PFT | Flux, $CO_2$ |
| $G_{s,slope}$ | Ball-Berry slope | PFT | Flux, $CO_2$ |
| $c_{T,opt}$ | Optimal photosynthesis temperature (°C) | PFT | Flux, $CO_2$ |
| $SLA$ | Specific leaf area ($m^2 \cdot g^{-1}$) | PFT | Flux, $CO_2$ |
| $K_{LAI,happy}$ | LAI threshold to stop using carbohydrate reserves | PFT | Sat, Flux, $CO_2$ |
| $K_{pheno,crit}$ | Multiplicative parameter of the threshold that determines the start of the growing season | PFT | Sat, Flux, $CO_2$ |
| $L_{age,crit}$ | Average critical age of leaves (days) | PFT | Sat, Flux, $CO_2$ |
| $C_{T,sen}$ | Temperature threshold for senescence (°C) | PFT | Sat, Flux, $CO_2$ |
| $F_{stress,h}$ | Parameter reducing the hydric limitation of photosynthesis | PFT | Flux, $CO_2$ |
| $MR_{offset}$ | Offset of the temperature dependence of maintenance respiration | Global | Flux, $CO_2$ |
| $Q10$ | Temperature dependency of heterotrophic respiration | Global | Flux, $CO_2$ |
| $HR_{Hc}$ | Offset of the soil/litter moisture control function | Global | Flux, $CO_2$ |
| $K_{soilC,site}$ | Multiplicative factor of the initial soil carbon pools | per Site | Flux |
| $K_{soilC,reg}$ | | 36 regions | $CO_2$ |
| $K_{albedo}$ | Multiplicative factor of the vegetation albedo | Global | Flux, $CO_2$ |

4

5

6

Geoscientific
Model Development
Discussions

1    Table 2. Prior information for all parameters except initial soil C pool multipliers: prior value,

2    uncertainty and range of variation for the different plant functional types (Tropical Broadleaf

3    Evergreen/Raingreen forests (TrBE / TrBR), Temperate Needle leaf / Broadleaf Evergreen

4    forests (TeNE, TeBE), Temperate Broadleaf Deciduous forest (TeBD), Boreal Needle leaf

5    Evergreen forests (BoNE), Boreal Broadleaf / Needle leaf Deciduous forests (BoBD / BoND)

6    and C3 grassland.

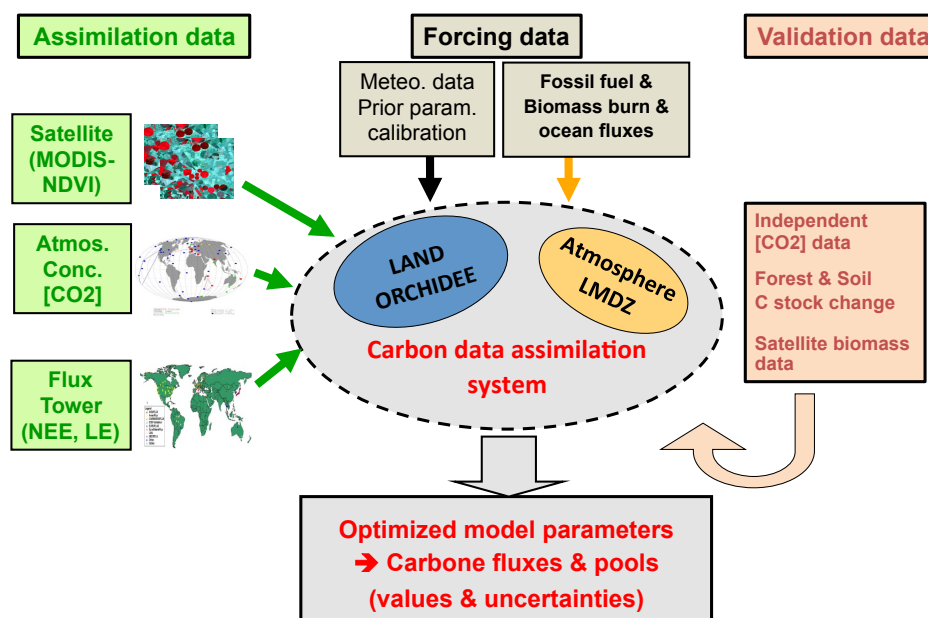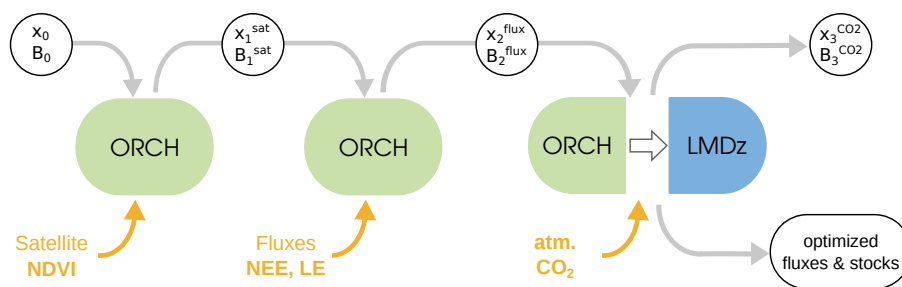| Parameter | Plant functional type | | | | | | | | |
|---|---|---|---|---|---|---|---|---|---|
| | TrBE | TrBR | TeNE | TeBE | TeBD | BoNE | BoBD | BoND | NC3 |
| $V_{cmax}$ | 65 ± 24 [35; 95] | 65 ± 24 [35; 95] | 35 ± 12.8 [19; 51] | 45 ± 16 [25; 65] | 55 ± 20 [30; 80] | 35 ± 12.8 [19; 51] | 45 ± 16 [25; 65] | 35 ± 12.8 [19; 51] | 70 ± 25.6 [38; 102] |
| $G_{s,slope}$ | 6.0 ± 2.4 [6; 12] | 6.0 ± 2.4 [6; 12] | 6.0 ± 2.4 [6; 12] | 6.0 ± 2.4 [6; 12] | 6.0 ± 2.4 [6; 12] | 6.0 ± 2.4 [6; 12] | 6.0 ± 2.4 [6; 12] | 6.0 ± 2.4 [6; 12] | 6.0 ± 2.4 [6; 12] |
| $c_{T,opt}$ | 37 ± 6.4 [29; 45] | 37 ± 6.4 [29; 45] | 25 ± 6.4 [17; 33] | 32 ± 6.4 [24; 40] | 26 ± 6.4 [18; 34] | 25 ± 6.4 [17; 33] | 25 ± 6.4 [17; 33] | 25 ± 6.4 [17; 33] | 27.25 ± 6.4 [19.25; 35.25] |
| $SLA$ | 0.015 ± 0.0092 [0.007; 0.03] | 0.026 ± 0.0148 [0.013; 0.05] | 0.009 ± 0.0064 [0.004; 0.02] | 0.02 ± 0.012 [0.01; 0.04] | 0.026 ± 0.0148 [0.013; 0.05] | 0.009 ± 0.0064 [0.004; 0.02] | 0.026 ± 0.0148 [0.013; 0.05] | 0.009 ± 0.0064 [0.004; 0.02] | 0.026 ± 0.0148 [0.013; 0.05] |
| $K_{LAI,happy}$ | 0.50 ± 0.14 [0.35; 0.70] | 0.50 ± 0.14 [0.35; 0.70] | 0.50 ± 0.14 [0.35; 0.70] | 0.50 ± 0.14 [0.35; 0.70] | 0.50 ± 0.14 [0.35; 0.70] | 0.50 ± 0.14 [0.35; 0.70] | 0.50 ± 0.14 [0.35; 0.70] | 0.50 ± 0.14 [0.35; 0.70] | 0.50 ± 0.14 [0.35; 0.70] |
| $K_{pheno,crit}$ | — | 1.0 ± 0.44 [0.7; 1.8] | — | — | 1.0 ± 0.44 [0.7; 1.8] | — | 1.0 ± 0.44 [0.7; 1.8] | 1.0 ± 0.44 [0.7; 1.8] | 1.0 ± 0.44 [0.7; 1.8] |
| $L_{age,crit}$ | 730 ± 192 [490; 970] | 180 ± 48 [120; 240] | 910 ± 240 [610; 1210] | 730 ± 192 [490; 970] | 180 ± 48 [120; 240] | 910 ± 240 [610; 1210] | 180 ± 48 [120; 240] | 180 ± 48 [120; 240] | 120 ± 60 [30; 180] |
| $C_{T,sen}$ | — | — | — | — | 12 ± 8 [2; 22] | — | 7 ± 8 [−3; 17] | 2 ± 8 [−8; 12] | −1.375 ± 8 [−11.375; 9.375] |
| $F_{stress,h}$ | 6.0 ± 3.2 [2; 10] | 6.0 ± 3.2 [2; 10] | 6.0 ± 3.2 [2; 10] | 6.0 ± 3.2 [2; 10] | 6.0 ± 3.2 [2; 10] | 6.0 ± 3.2 [2; 10] | 6.0 ± 3.2 [2; 10] | 6.0 ± 3.2 [2; 10] | 6.0 ± 3.2 [2; 10] |
| $MR_{offset}$ | 1.0 ± 0.6 [0.5; 2.0] | | | | | | | | |
| $Q10$ | 1.99372 ± 0.8 [1.0; 3.0] | | | | | | | | |
| $HR_{Hc}$ | −0.29 ± 0.24 [−0.59; 0.01] | | | | | | | | |
| $K_{albedo}$ | 1.0 ± 0.16 [0.8; 1.2] | | | | | | | | |

1 **Figures**

2

3

4

5



6

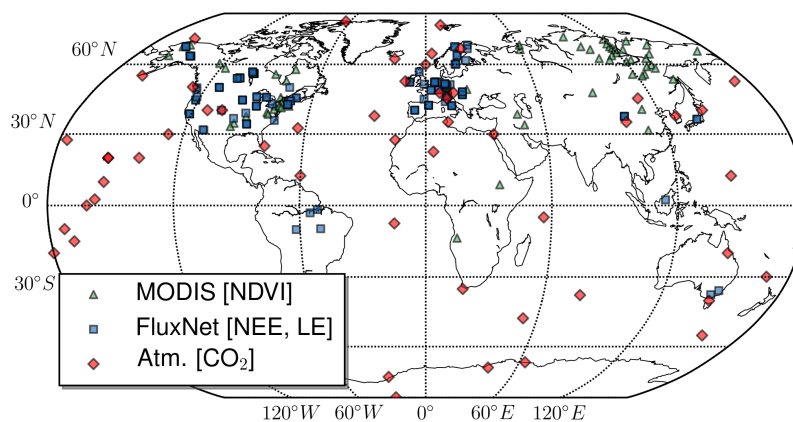7 Figure 1. Schematic of the ORCHIDEE Carbon Cycle Data Assimilation System
8 (ORCHIDAS).

9

Figure 2. Illustration of the step-wise data assimilation approach used for the assimilation of multiple data streams in the ORCHIDEE-CCDAS. The list of parameters for each step is summarized in Table 1.
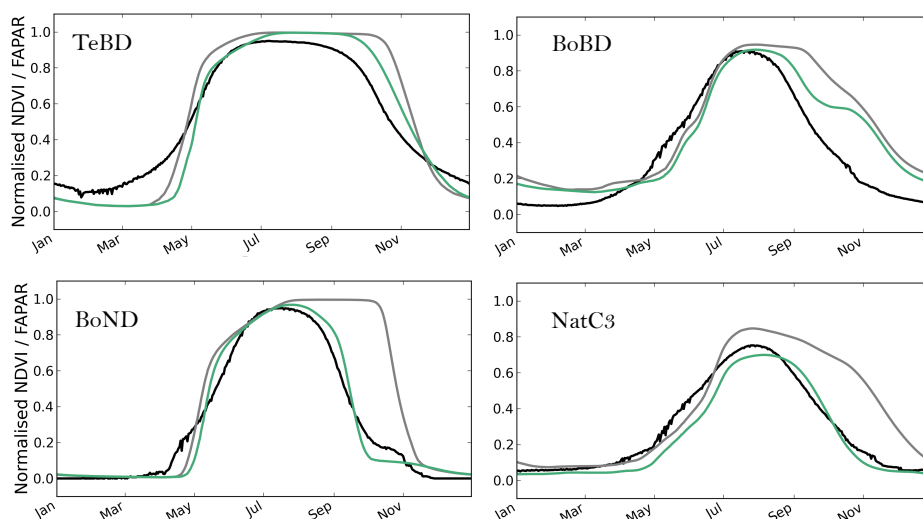


Figure 3: Location of the different observations used for each data stream assimilated in the system: MODIS-NDVI measurements, FLUXNET sites with NEE and LE measurements and atmospheric $CO_2$ stations.

1



2   Figure 4. Mean seasonal cycle (2000-2008) of the normalised modelled FAPAR before and
3   after optimisation, compared to that of the MODIS NDVI data, for the temperate and boreal
4   deciduous PFTs (TeBD, BoBD, BoND and NatC3). Black = MODIS NDVI data; Grey =
5   prior simulation (default ORCHIDEE parameters); Green = posterior multi-site optimisation.

6

7

8



9    Figure 5: Mean seasonal cycle of the Net Carbon Ecosystem Exchange (NEE) for the
10   different plant functional type optimize in Step 2 of the assimilation. The mean across all sites
11   for a given PFT is provided for the observations (black), the prior ORCHIDEE (grey), the
12   posterior of step 1 (green) and the posterior of step 2 (blue).

1



Figure 6: Monthly mean atmospheric $CO_2$ concentrations after step 3 of the optimization, for several stations over the period 2002-2004 of the optimization. The observations (black), the prior model (grey) and the posterior model after step 2 (blue) and step 3 (red) are displayed. Numbers in parenthesis correspond to RMSEs.

6

7



8

Geoscientific
Model Development
Discussions

1    Figure 7: Changes in the mean seasonal cycle of the atmospheric $CO_2$ concentrations after
2    step 3 of the optimization at all atmospheric stations. Left: Relative changes (in percentage)
3    between the prior and posterior absolute model-data differences for the amplitude of the
4    seasonal cycle. Right: Same metric but for the length of the Carbon Uptake Period (CUP),
5    measured as the sum of the days when the de-trended concentrations are negative (see text).

6

7

8



9

10   Figure 8: RMSE between model outputs and observations for two types of observations:
11   MODIS-NDVI on the left and FluxNet-NEE on the right, for different Plant Functional Types
12   (PFT: TrBE, TeNE, TeBE, TeBD, BoBD, BoND, NC3) and for the prior model simulation
13   and the posterior of each step of the assimilation framework. Missing bars correspond to the
14   fact that no data were available to constrain a given PFT.
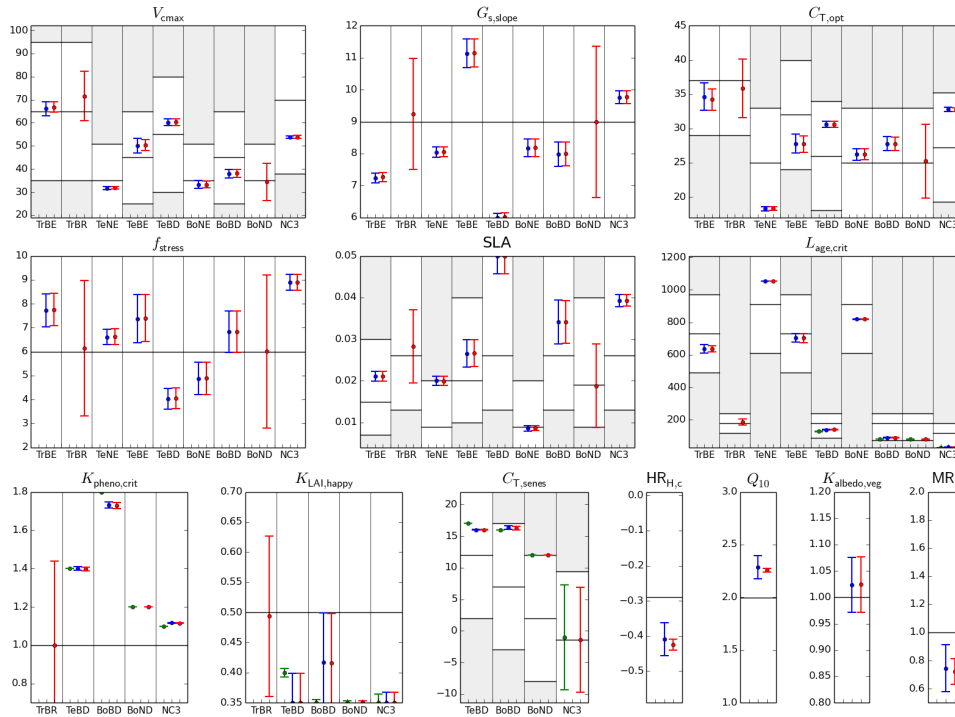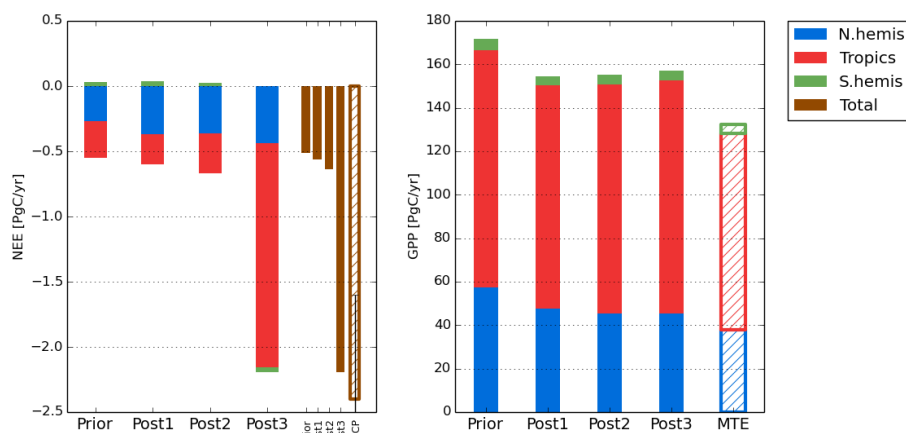
15

1
2 Figure 9: Prior and posterior parameter values and uncertainties for a set of optimized
3 parameters (9 PFT dependent and 4 non-PFT dependent). The prior value corresponds to the
4 horizontal black line and the physical allowed range of variation to the "y" range (i.e. the
5 white zone). For PFT-dependent parameters, there are 9 sub-plots corresponding to PFTs that
6 were optimized (except for $K_{pheno\_crit}$ with only 5 PFTs). For each parameter, there are 3
7 estimated values for the three successive steps: step1: assimilation of MODIS-NDVI data
8 (green symbol); step2: adding FLUXNET data (blue symbol); step3: adding atmospheric $CO_2$
9 data (red symbol). The parameter values are depicted with the symbols and the estimated
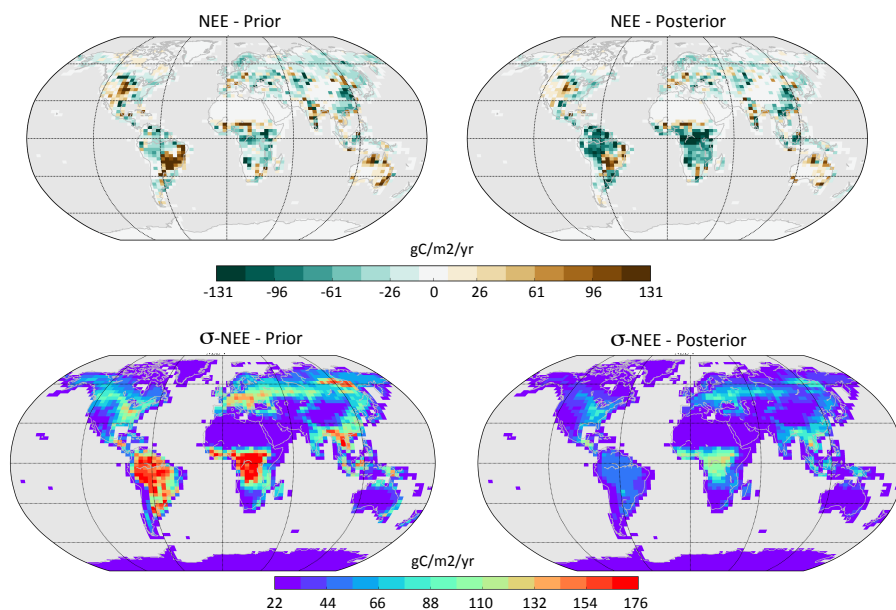10 uncertainties with the vertical line (± sigma).

11

12

13

14

15

16

17

18

Figure 10: Left: Net Ecosystem Exchange (NEE) for three regions (North of 35°N, Tropics, South of 35°S) for the prior model, and after each step of the optimizations (mean over 2002-2004). The total NEE is indicated with the vertical brown bar and compared to the Global Carbon Project (GCP) estimate for the same period (Le Quéré et al. 2015). Right: same but for Gross Primary Production where the data driven estimate (MTE product using FluxNet data; Jung et al., 2009) is provided for comparison.
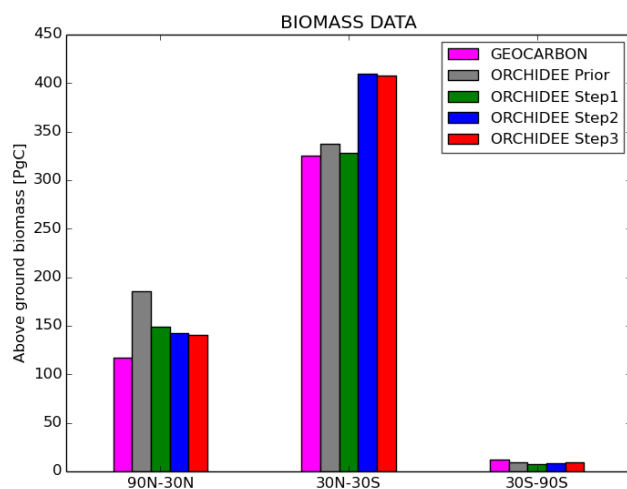
1



2

Figure 11: Simulated annual net carbon exchange (NEE) for the land ecosystems prior to any optimization (left column) and after step 3 of the optimization process (right column). Upper figures correspond to the mean NEE (in $gC.m^{-2}.y^{-1}$) over the assimilation period (2001-2003) and lower figures to the associated monthly flux uncertainties (averaged over the whole period and expressed in $gC.m^{-2}.y^{-1}$) due to the parameter uncertainties (see text).

8

1



2     Figure 12: Above ground forest biomass data for the prior ORCHIDEE model and after step
3     1, step 2 and step 3 of the optimization process. Estimates from satellite observations (Santoro
4     et al., 2015) and referred as "GEOCARBON" (following the EU-GEOCARBON project) are
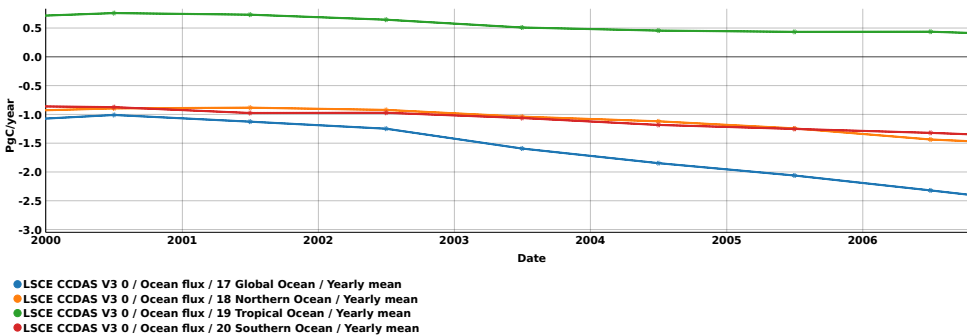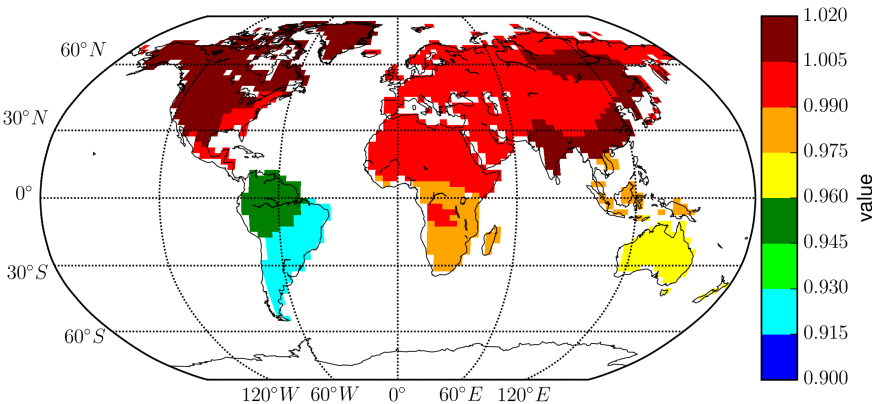5     provided for comparison.

6

7

# 1   Appendix figures

2



3

Figure A1: $CO_2$ air-sea fluxes including the natural ocean out-gazing, used as input to the
ORCHIDEE-CCDAS and estimated from a neural network approach using observed $pCO_2$
data (see main text, section 2.5.1). The Northern, Tropical and Southern ocean contributions
to the global ocean flux (blue curve) are also provided.

8



9

Figure A2: Map of the posterior values of the coefficient scaling the initial carbon pool sizes
per regions.

12

13

14