

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 **Philippe Peylin¹, Cédric Bacour², Natasha MacBean¹, Sébastien Leonard¹, Peter**
6 **Rayner^{1,3}, Sylvain Kuppel^{1,4}, Ernest Koffi¹, Abdou Kane¹, Fabienne Maignan¹,**
7 **Frédéric Chevallier¹, Philippe Ciais¹, Pascal Prunet²**

8 [1]{Laboratoire des Sciences du Climat et de l'Environnement, UMR 8212 CEA-CNRS-
9 UVSQ, 91191 Gif-sur-Yvette cedex, France}

10 [2]{Noveltis, Parc Technologique du Canal, 2 avenue de l'Europe, 31520 Ramonville-Saint-
11 Agne, France}

12 [3]{University of Melbourne, 3010, Vic, Melbourne, Australia}

13 [4]{Grupo de Estudios Ambientales, IMASL-CONICET/Universidad Nacional de San Luis,
14 San Luis, Argentina}

15 Correspondence to: Philippe Peylin (philippe.peylin@lsce.ipsl.fr)

16

17

1 **Abstract**

2 Large uncertainties in Land surface models (LSMs) simulations still arise from inaccurate
3 forcing, poor description of land surface heterogeneity (soil and vegetation properties),
4 incorrect model parameter values and incomplete representation of biogeochemical processes.
5 The recent increase in the number and type of carbon cycle related observations, including
6 both in situ and remote sensing measurements, has opened a new road to optimize model
7 parameters via robust statistical model-data integration techniques, in order to reduce the
8 uncertainties of simulated carbon fluxes and stocks. In this study we present a carbon cycle
9 data assimilation system that assimilates three major data streams, namely MODIS-NDVI
10 observations of vegetation activity, net ecosystem exchange (NEE) and latent heat (LE) flux
11 measurements at more than 70 sites (FLUXNET), and atmospheric CO₂ concentrations at 53
12 surface stations, in order to optimize the main parameters (around 180 parameters in total) of
13 the ORCHIDEE LSM (version 1.9.5 used for CMIP5 simulations). The system relies on a
14 step-wise approach that assimilates each data stream in turn, propagating the information
15 gained on the parameters from one step to the next.

16 Overall, the ORCHIDEE model is able to achieve a consistent fit to all three data streams,
17 which suggests that current LSMs have reached the level of development to assimilate these
18 observations. The assimilation of MODIS-NDVI (step 1) reduced the growing season length
19 in ORCHIDEE for temperate and boreal ecosystems, thus decreasing the global mean annual
20 gross primary production. Using FLUXNET data (step 2) led to large improvements in the
21 seasonal cycle of the NEE and LE fluxes for all ecosystems (i.e., increased amplitude for
22 temperate ecosystems). The assimilation of atmospheric CO₂, using the general circulation
23 model of the Laboratoire de Météorologie Dynamique (LMDz; step 3), provides an overall
24 constraint (i.e., constraint on large scale net CO₂ fluxes), resulting in an improvement of the
25 fit to the observed atmospheric CO₂ growth rate. Thus the optimized model predicts a land C
26 sink of around 2.2 PgC.yr⁻¹ (for the 2000-2009 period), which is more compatible with
27 current estimates from the Global Carbon Project than the prior value. The consistency of the
28 step-wise approach is evaluated with back-compatibility checks. The final optimized model
29 (after step 3) does not significantly degrade the fit to MODIS-NDVI and FLUXNET data that
30 were assimilated in the first two steps, suggesting that a stepwise approach can be used
31 instead of the more “challenging” implementation of a simultaneous optimization in which all
32 data streams are assimilated together. Most parameters, including the scalar of the initial soil

1 carbon pool size, changed during the optimization with a large error reduction. This work
2 opens new perspectives for better predictions of the land carbon budgets.

3

4 **1 Introduction**

5 Atmospheric CO₂ 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 CO₂, 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 CO₂, 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 carbon cycle components (Le Quéré et al., 2015). In addition to uncertainties in the mean
16 global annual terrestrial C budget and its trend over time (Sitch et al., 2015), there remain
17 strong discrepancies between LSMs in their predictions of regional budgets (Canadell, 2013)
18 at seasonal and inter-annual timescales and in their sensitivity to climate and atmospheric CO₂
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 CO₂ concentration to constrain the Simple
7 Diagnostic Biosphere Model (SDBM). Later, Kaminski et al (2002) constrain the seasonal
8 cycle of SDBM with the same data stream. This effort was continued by the original Carbon
9 Cycle Data Assimilation System (CCDAS) described in Rayner et al. (2005) and Kaminski et
10 al. (2012) which used both atmospheric CO₂ and satellite-derived Fraction of Absorbed
11 Photosynthetic Radiation (FAPAR) data to optimize vegetation productivity by adjusting the
12 C cycle-related parameters of the Biosphere Energy-Transfer Hydrology (BETHY) model
13 (see a review in Kaminski et al., 2013). Note that although Rayner et al. (2005) did use, in
14 addition to atmospheric CO₂ data, soil moisture and radiation fields from an earlier
15 assimilation from a simpler model version, no parameters were passed between the two
16 assimilations and very little comment was made on the consistency between the two
17 assimilations, an important issue that will be central to this paper. Meanwhile substantial
18 efforts have been put into the use of local eddy covariance flux tower measurements of net
19 exchange of CO₂ and latent and sensible heat fluxes to optimize photosynthesis, respiration
20 and energy-related parameters of terrestrial ecosystem models, both at individual sites (e.g.
21 Wang et al., 2001, 2007; Williams et al., 2005; Braswell et al., 2005; Knorr and Kattge, 2005;
22 Moore et al., 2008; Ricciuto et al., 2008), and more recently using multiple sites together
23 (hereafter multiple sites) from the global FLUXNET network (e.g. Groenendijk et al., 2011;
24 Kuppel et al., 2012, 2014; Alton, 2013; Xiao et al., 2014). Increasingly the focus in carbon
25 cycle data assimilation is moving towards using multiple different data streams as
26 independent constraints, with the aim of bringing more information at different spatial and
27 temporal scales and constraining several processes at once in order to reduce the likelihood of
28 model equifinality (where multiple sets of parameters achieve the same reduction in model-
29 data misfit). Recent examples include the combination of in-situ eddy covariance flux
30 observations and ground-based information on vegetation structure and C stocks (Richardson
31 et al., 2010; Ricciuto et al., 2011; Keenan et al., 2012, 2013; Thum et al., 2015), or in-situ flux
32 data and satellite FAPAR (Kato et al., 2013; Zobitz et al., 2014; Bacour et al., 2015) or
33 atmospheric CO₂ and biomass data using a simple biosphere model (Saito et al., 2014). This is

1 a non-trivial task however, especially when optimizing a complex LSM (see MacBean et al,
2 submitted), which has many parameters acting from local to global scales.

3 When assimilating multiple different data streams we have two options: i) to optimize the
4 model with each data stream in turn, and to propagate the information gained on the
5 parameter values from one step to the next (hereafter referred to as “stepwise” assimilation),
6 or ii) to include all data streams together in the same optimization (hereafter referred to as
7 “simultaneous” assimilation). Kaminski et al. (2012) suggested that it is essential to perform a
8 consistent, simultaneous assimilation that includes all data streams in the same optimization.
9 It is important to note that this is an implementation question. Tarantola (2005) recasts the
10 fundamentals of the approach as the conjunction or multiplication of probability densities.
11 This multiplication is associative so it makes no difference whether it is performed in one step
12 or several (and whether the system is linear or not). In complex problems such as these, one
13 cannot carry or even describe the full structure of the relevant probability densities so which
14 approach will work best in each case is unclear. In particular, technical difficulties associated
15 with the different number of observations for each data stream and the characterization of
16 error correlations between them, in addition to computational constraints to run global LSMs,
17 might result in the preference for a step-wise assimilation framework. Additionally, it may be
18 more straightforward, to expose a restricted set of parameters (following a global sensitivity
19 analysis) to each observation type in a stepwise approach to ensure that each data stream
20 constrains only the most relevant parts of the model. This reduces biases from other poorly-
21 represented processes caused by inadequate model structure. Note finally that more complex
22 approaches based on random generation of parameter sets, such as the multi-objective
23 approach using the Pareto ranking of several cost functions (e.g. Yapo et al., 1998), are not
24 yet affordable for global LSMs from a computational point of view. For these reasons we
25 follow the stepwise approach in this paper.

26 We present the first global-scale CCDAS that assimilates three of the main global data
27 streams that have been used to date to understand the terrestrial carbon cycle – atmospheric
28 CO₂ concentration, satellite-derived information of vegetation greenness (from the MODIS
29 instrument) and multisite eddy covariance net CO₂ and latent heat flux measurements (from
30 FLUXNET) – to optimize the parameters of the Organizing Carbon and Hydrology in
31 Dynamics Ecosystems (ORCHIDEE) process-based LSM (Krinner et al., 2005). This study is
32 the first (to our knowledge) to assimilate these three major data streams in a process-based

1 LSM used as the land component of an Earth System Model (ESM), the French Institut Pierre
2 Simon Laplace ESM. Two contemporary studies also optimize the parameters of the land
3 component of an ESM; however Raoult, et al. (2016) only uses FluxNet observations to
4 optimize the parameters of the JULES model, while Schürmann et al. (2016) only assimilate
5 two data streams (fAPAR and CO₂) in the JSBACH model at coarse resolution (10° x 10°).
6 Note finally that the level of complexity of the ecosystem model (and the spatial resolution) is
7 part of the problem: achieving an optimization with a given model does not guarantee that the
8 framework would work with a more complex or different one.

9 In this context, the main questions that we aim to answer in this paper are as follows:

- 10 i) How and to which extend the optimization of the ORCHIDEE model allows to fit the three
11 data streams that are considered?
- 12 ii) Does the step-wise optimization result in a degradation of the fit to other data streams used
13 in the previous steps?
- 14 iii) What are the main changes in the optimized parameters when using sequentially these
15 three data streams in a global CCDAS and which processes are constrained?
- 16 iv) What are the improvements for the land C cycle in terms of net/gross fluxes and stocks as
17 a result of multi-data stream optimization? What preliminary perspectives can we draw that
18 may help us in improving model predictions of trends, variability and the location of
19 terrestrial C sources and sinks?

20 Following these objectives, the paper first describes the new ORCHIDEE-CCDAS including
21 the concept, the observations, the models and the optimization approach. We then present the
22 results, including the fit to the data, consistency checks (question i) above) as well as mean
23 global and regional C cycle budget for the period 2000-2009. The last section discusses issues
24 and perspectives associated with these results.

25

26 **2 Methods**

27 **2.1 ORCHIDEE-CCDAS concept**

28 We have designed a CCDAS around the ORCHIDEE land surface model (ORCHIDEE-
29 CCDAS, later also referred to as ORCHIDAS for simplicity) that combines a state-of-the-art

1 description of the driving biogeochemical processes within the model with multiple
2 observational constraints in a robust statistical framework, in order to improve the simulation
3 of land carbon fluxes and stocks. The system allows us to retrieve the best estimate, given the
4 observations and prior information, of selected parameters (see §2.3.3) as well as to evaluate
5 their uncertainty. It relies on a stepwise assimilation of a comprehensive set of three C cycle-
6 related observations that are representative of small (100 m) to large (continental) scales (see
7 §2.2):

- 8 • Step 1: Satellite measurements of vegetation activity using the Normalized Difference
9 Vegetation Index (NDVI) from the MODIS instrument over the 2000-2008 period for
10 a randomly selected set of sites for boreal and temperate deciduous vegetation types;
- 11 • Step 2: In-situ eddy-covariance net CO₂ and water (latent heat) flux measurements
12 from the FLUXNET database for a large set of sites, spanning 7 different vegetation
13 types;
- 14 • Step 3: In-situ monthly atmospheric surface CO₂ concentration measurements from
15 the GLOBALVIEW-CO₂ database over three years (2002-2004).

16 The system relies on two models:

- 17 • The ORCHIDEE global LSM, whose main C cycle parameters are optimized (see
18 §2.3)
- 19 • The general circulation model of the Laboratoire de Météorologie Dynamique, LMDz
20 (see §2.3), to relate the surface carbon fluxes to atmospheric CO₂ concentrations.

21 The framework combines the different observational data streams within ORCHIDAS in
22 order to optimize selected model parameters using a variational data assimilation system,
23 described in section 2.4. Figure 1 illustrates the structure of the CCDAS and the different
24 components that are involved. Such a framework distinguishes i) the assimilated observations,
25 ii) an ensemble of forcing and input data streams, iii) the models and optimization framework,
26 as well as iv) an evaluation step, where independent datasets are compared to the optimized
27 model stocks and fluxes. As explained in the introduction, a major feature of the current
28 system is the stepwise approach, in which all data streams are assimilated sequentially (i.e.
29 one after the other). The information retrieved at a given step (retrieved optimal parameter
30 values and associated uncertainty) is propagated to the next step (see Fig. 2 and §2.4). Note
31 that for simplicity we did not propagate the error correlations in this first implementation of

1 the system, a simplification that appeared sufficient (see the consistency analysis in section
2 3.2); section 4 also discusses the potential impact of this simplification.

3 At each step, the parameter optimization relies on a Bayesian framework that explicitly
4 minimizes the difference between the simulated and observed quantities in addition to
5 minimizing the difference between the optimized model parameters and “a priori” values (see
6 §2.4.2). The dependence of the simulated quantities on the optimized variables is non-linear,
7 which thus necessitates the use of an iterative algorithm. Note that all components of the
8 surface C budget need also to be included in the ORCHIDAS, particularly when using
9 atmospheric CO₂ measurements which requires the atmospheric transport model to be
10 prescribed with fossil fuel emissions, CO₂ fluxes associated with biomass burning and ocean
11 CO₂ fluxes (see §2.5) in addition to net ecosystem exchange (NEE) from ORCHIDEE.

12 **2.2 Assimilated observations**

13 **2.2.1 MODIS-NDVI**

14 MODIS collection 5 obtained from surface reflectance data (from 2000-2008) in the red (R)
15 and near-infrared (NIR) bands at 5 km resolution (CMG) are used to optimize the phenology-
16 related parameters of ORCHIDEE in the first step. The R and NIR data were processed to
17 correct for directional effects following Vermote et al. (2009) and then used to calculate the
18 NDVI, which is assumed to be linearly related to the model FAPAR. The NDVI are then i)
19 aggregated to the 0.72° spatial resolution of the ERA-Interim meteorological fields that are
20 used to force ORCHIDEE, ii) interpolated to a daily time series (for practical implementation)
21 and iii) checked for quality (see MacBean et al., 2015 for details). If there is a gap in the
22 observations of more than 15 days, no interpolation is done (i.e., no data during the gap are
23 assimilated). Figure 3 displays the location of the sites that were selected (see §2.4.1).

24 **2.2.2 Eddy covariance flux data**

25 Eddy covariance flux measurements of net surface CO₂ flux – hereafter referred to as net
26 ecosystem exchange (NEE) and latent heat flux (LE) – from 78 observation sites of a network
27 of regional networks (FLUXNET; see Fig. 3) are used to constrain ecosystem physiology and
28 fast C-related processes at daily to seasonal timescales in ORCHIDEE in the second step. We
29 use quality-checked and gap-filled data from a global synthesis called the La Thuile dataset
30 (Papale, 2006). In order to avoid dealing with the large error correlations in the half-hourly

1 data (see Lasslop et al., 2008), daily mean values of NEE and LE are used in the ORCHIDAS.
2 Days with less than 80% of the half-hourly data are left out of the assimilation. The selection
3 of the sites and the data processing (gap-filling, correction for energy balance closure) are
4 detailed in Kuppel et al. (2014). Note that uncertainties due to incomplete sampling of the
5 diurnal cycle are likely very small (less than 5%) as the error in the gap-filling procedure is
6 usually less than 20% (Lasslop et al., 2008).

7

8 2.2.3 Atmospheric CO₂ concentrations

9 Atmospheric CO₂ concentration measurements were taken from an ensemble of selected
10 surface stations around the world (Fig. 3). The spatial concentration gradients relate to the
11 integral of the fluxes over large areas and thus allow the optimization of large-scale global
12 patterns of carbon fluxes. These data were taken from the NOAA Earth System Laboratory
13 (ESRL) GLOBALVIEW-CO₂ collaborative product (GLOBALVIEW-CO₂, 2013) and
14 averaged to monthly means. We assimilated the monthly values for 53 sites for the 2002-2004
15 period inclusive in the last step of the assimilation system. Such restricted period (3 years
16 only) was chosen for practical reasons (computing resources) while constructing the
17 ORCHIDAS system. The station locations, indicated in Fig. 3, favor the background
18 conditions i.e. the surrounding air masses are only weakly influenced by local continental
19 sources, such as power plants. The choice of monthly mean is related to the use of pre-
20 calculated transport fields with LMDZ (see §2.3.2). We also used additional sites to evaluate
21 the result of the optimization (locations indicated in Fig. 3): this included 17 continental sites
22 that are more directly influenced by local fluxes potentially not well captured at the
23 considered LMDz spatial resolution and 7 sites from Pacific Ocean cruises that were not
24 included in the optimization in order not to overweight that the data contribution from that
25 particular region. Note that we did not considered free troposphere aircraft data or column
26 integrated measurements (TCCON sites) in this evaluation, although they are less sensitive to
27 biases in the Planetary Boundary Layer representation, given that i) we are using pre-
28 calculated transport fields previously computed at surface stations only and ii) few scarce free
29 tropospheric datasets will not bring much more information to the additional surface stations.

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 5th 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 (2.5° x 3.75°; see § 2.3.2). However,
28 when we assimilate in situ flux data in the second step, we force the model with the gap-filled
29 half-hourly meteorological data measured at each site. The global PFT map was derived from
30 the 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 the third step, we spun-up
2 the model recycling the 1989-1998 meteorology and then used a transient simulation from
3 1990 to 2001 with changing climate (ERA-Interim) and increasing CO₂, 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 CO₂. Note that the use of soil carbon data, such as from the
7 Harmonized World Soil Database (as well as above ground biomass data), to initialize the
8 model is not straightforward and represents a challenge to keep the internal model
9 consistency, given that the three soil carbon reservoirs of the CENTURY model are in
10 balance with all components of the model, in particular the input through the different litter
11 pools. Computational and scientific issues to avoid a spin-up approach are still under
12 investigation with ORCHIDEE (see discussion section).

13 2.3.2 LMDz model

14 The transport model used in this study is version 3 of the General Circulation Model (GCM),
15 LMDz (Hourdin and Armengaud, 1999) with a horizontal resolution of 3.75° (longitude) x
16 2.5° (latitude) and 19 sigma-pressure layers up to 3 hPa. The calculated winds (u and v) are
17 relaxed to the ECMWF ERA-40 meteorological data (Uppala et al. 2005) with a relaxation
18 time of 2.5h (guiding) in order to realistically account for large-scale advection (Hourdin et
19 al., 2000). Deep convection is parameterized according to the scheme of Tiedtke (1989) and
20 the turbulent mixing in the planetary boundary layer is based on a local second-order closure
21 formalism. The LMDz GCM model has been widely used to model climate (IPCC, 2007,
22 2013) and its derived transport model has been used for the simulation of chemistry of gas
23 and particles and greenhouse gases distributions (Hauglustaine et al., 2004; Folberth et al.,
24 2005; Bousquet et al. 2005, 2006; Rivier et al., 2006). For this study, we used pre-calculated
25 transport fields, as described in Peylin et al. (2005), that correspond to the sensitivity of
26 concentration at each atmospheric site and each month to the surface flux of each model grid-
27 cell for each day (often called influence functions). The sensitivities (using inter-annual
28 winds) were calculated with the “retro-transport” formulation implemented in the LMDz
29 transport model (Hourdin et al. 2006). This approach decreases the computing time of the
30 optimization compared to the use of the full forward LMDz model at each iteration, as the
31 transport is replaced by a matrix multiplication with the vector of surface fluxes. Note that the
32 initial 3D state of the atmospheric concentrations was defined from Chevallier et al. (2010)

1 2.3.3 Parameters optimized

2 The optimized parameters are described in Table 1, and their prior values, uncertainty and
3 range are given in Table 2. In the most recent studies using ORCHIDAS at site scales a large
4 set of ORCHIDEE parameters has been optimized (Kuppel et al., 2014; Santaren et al., 2014;
5 Bacour et al., 2015). In this study a smaller set was chosen, based on a Morris sensitivity
6 analysis (Morris, 1991; results not shown) that determines the sensitivity of the NEE and LE
7 to all model parameters at various FLUXNET sites (for each PFT), in order to reduce the
8 computational cost of the global optimization in step 3 (see §2.5). We considered 9 PFT-
9 dependent and 4 “global” (i.e. non PFT-dependent) parameters that control mostly the fast
10 carbon processes (diurnal to seasonal). In addition, we introduced a new parameter, K_{soilC} , to
11 scale the initial values (after spin-up) of the modeled slow and passive soil carbon pools, in
12 order to take account of all the historical effects not accounted for in the model that would
13 result in a disequilibrium of these pools in reality. For the site-specific optimizations with
14 FLUXNET data, we have one $K_{soilC,site}$ parameter per site. For the global scale optimization
15 step, we used 30 $K_{soilC,reg}$ parameters corresponding to 30 regions potentially coherent for land
16 use and land management history as well as ecosystem and edaphic properties (see Fig. A2).
17 The initial soil carbon pools of all pixels within each region were thus scaled by the same
18 value. The prior value for all K_{soilC} parameters was set to one, i.e. the default state of soil
19 carbon pools is assumed to be in equilibrium.

20 Overall (including all PFT-dependent parameters), we optimize 16 parameters related to
21 phenology, 36 to photosynthesis, 3 to respiration, 1 to the energy budget, 78 soil C pool
22 scalars (one for each FLUXNET site), and 30 regional soil C pool scalars for the global
23 simulations – a total of 184 parameters (16, 134 and 86 in step 1, 2 and 3, respectively). Note
24 that the soil C pool multipliers at the FLUXNET sites are independent from the regional C
25 pool multipliers, as the history of soil carbon over large eco-regions of several millions square
26 kilometers is rather heterogeneous (as it is mainly related to previous land use changes), and
27 most likely, the FLUXNET sites are not representative of larger regions in terms of the soil
28 carbon disequilibrium. The prior standard deviation for each parameter is equal to 40% of the
29 parameter range (lower and higher boundaries) prescribed for each parameter following
30 Kuppel et al. (2012). The parameter ranges were specified following expert judgment of their
31 meaning in the ORCHIDEE equations and based on literature reviews or databases (such as
32 TRY, Kattge et al., 2011).

1 **2.4 System description: a step-wise approach**

2 **2.4.1 Stepwise assimilation of three data streams**

3 The ORCHIDAS system relies on a stepwise assimilation of the three data streams described
4 in section 2.2. Figure 2 illustrates the flow of information in this sequential approach:

5 ***Step 1 – Assimilation of MODIS-NDVI:*** Four parameters related to the seasonal cycle of the
6 vegetation (phenology) are optimized for the temperate and boreal deciduous PFTs (TeBD,
7 BoND, BoBD and NC3 – see caption of Table 2). These four deciduous PFTs alone are
8 considered in step 1 in this ORCHIDAS version because the tropical deciduous phenology
9 modules in ORCHIDEE require further modifications to improve the functions that control
10 leaf growth and fall in response to water availability (MacBean et al., 2015). Evergreen PFTs
11 were also not considered, as there are no phenology modules related to these PFTs in the
12 model. The procedure is similar to that described in detail in MacBean et al. (2015) and
13 therefore only briefly recalled here. A simple linear relationship between the modeled
14 Fraction of Absorbed Photosynthetically Active Radiation (FAPAR) and MODIS-NDVI
15 observations is assumed, based on studies such as Knyazikhin et al. (1998). Given that
16 considerable discrepancies exist between so-called “high-level” satellite products such as LAI
17 or fAPAR when considering their magnitude (D’Odorico et al., 2014), we thus only use the
18 temporal information in the NDVI observations and normalized both the model FAPAR
19 output and the NDVI observations to their 5th and 95th percentiles (following Bacour et al.,
20 (2015) and MacBean et al., 2015). Note that assimilating raw fAPAR data with the
21 ORCHIDEE model led to the degradation of the NEE with the estimation of spurious
22 parameter values (Bacour et al., 2015). The model was run for fifteen randomly selected grid
23 cells for each of the four PFTs using the ERA-Interim meteorological forcing. Only grid cells
24 that included vegetation fraction of greater than 60% for the PFT optimized were considered.
25 We selected a set of grid points instead of the whole grid to substantially decrease the
26 computing time; but the remaining points are used for the evaluation of the optimized model.
27 The fifteen sites for each PFT were included in one optimization for each PFT following a
28 multi-site approach, in which all observations are used simultaneously to optimize the model
29 parameters. The optimized parameters are described in Table 1. They correspond to a scalar
30 on the growing degree days (GDD) threshold for the start of the vegetation ($K_{pheno,crit}$), a
31 parameter controlling the use of carbohydrate reserve during the start of leaf growth

1 ($K_{lai,happy}$), a temperature threshold for the onset of leaf senescence (CT_{senes}) and the critical
2 age for leaves ($L_{agecrit}$).

3 **Step 2 – Assimilation of FLUXNET data:** Mean daily NEE and LE flux measurements for 78
4 sites, including up to 10 years worth of data for each site, are used to optimize a set of model
5 parameters controlling the fast carbon and water processes (photosynthesis, respiration,
6 phenology – see Table 1). The site selection and the choice of a daily time step are described
7 in more details in Kuppel et al. (2014). These sites cover 7 of the PFTs in ORCHIDEE (see
8 Table 2). The posterior parameter values of the four phenology parameters derived in step 1,
9 and their associated uncertainties, are input as prior information in step 2. For the additional
10 parameters, the default ORCHIDEE values are used for the prior and the uncertainties are set
11 as described in §2.3.3. A multi-site optimization is performed for each PFT independently as
12 in step 1. Global parameters, i.e. those that are not PFT-dependent, were optimized for each
13 PFT and the mean across all PFTs was then calculated to define the prior parameter vector in
14 step 3 of the assimilation with atmospheric CO₂ data (at global scale). Such an approach was
15 chosen to allow us to optimize all PFTs in parallel and therefore to simplify the assimilation
16 process.

17 **Step 3 – Assimilation of atmospheric CO₂ concentrations:** We use monthly mean CO₂
18 concentrations from 53 surface stations over three years (2002-2004) to provide a large-scale
19 constraint to the land surface fluxes (i.e. to match the global CO₂ growth rate, mean seasonal
20 cycle and its latitudinal variation, as well as the spatial gradients between stations). We use
21 the LMDz atmospheric transport model (see §2.3.2) to assimilate these observations. The set
22 of parameters optimized in step 2 are included in step 3, except for the albedo scaling
23 parameter ($K_{albedo,veg}$), as the net carbon fluxes are only weakly sensitive to that parameter. We
24 used the posterior parameter distributions from step 2 (parameter optimal values and
25 associated uncertainties) as prior information for step 3, and expanded the parameter vector to
26 include the 30 K_{soilC} parameters that scale the initial soil carbon pools for large “spatially-
27 coherent regions” (see §2.1.2 and Fig. A2). The air-sea fluxes and fossil fuel and biomass
28 burning emissions are also accounted for (but not optimized) in this final step, in order to
29 close the global carbon budget within the atmospheric transport model (see §2.5).

1 2.4.2 Optimization procedure (for all steps):

2 In each step the statistically optimal parameter values are derived with an optimization
3 procedure following the principle of the 4-D variational assimilation systems (developed for
4 numerical weather prediction), using a tangent linear operator (and finite differences for a few
5 parameters, Bacour et al. 2015). Assuming that the errors associated with the parameters, the
6 observations and the model outputs follow Gaussian distributions, the optimal parameter set
7 corresponds to the minimum of a cost function, $J(\mathbf{x})$, that measures the mismatch between i)
8 the observations (\mathbf{y}) and the corresponding model outputs, $H(\mathbf{x})$, (where H is the model
9 operator), and ii) the a priori (\mathbf{x}_b) and optimized parameters (\mathbf{x}), weighted by their error
10 covariance matrices (Tarantola, 1987; Chapter 4):

$$11 \quad J(\mathbf{x}) = \frac{1}{2} \left[(H(\mathbf{x}) - \mathbf{y})^T \mathbf{R}^{-1} (H(\mathbf{x}) - \mathbf{y}) + (\mathbf{x} - \mathbf{x}_b)^T \mathbf{B}^{-1} (\mathbf{x} - \mathbf{x}_b) \right] \quad (1)$$

12 \mathbf{R} represents the error variance/covariance matrix associated with the observations and \mathbf{B} the
13 parameter prior error variance/covariance matrix. At each step a different cost function is
14 defined with the observations and parameters related to that step (see Fig. 2). \mathbf{R} includes the
15 errors on the measurements, the model structure and the meteorological forcing. Model errors
16 are rather difficult to assess and may be much larger than the measurement error itself.
17 Therefore we chose to focus on the structural error and defined the variances in \mathbf{R} as the mean
18 squared difference between the prior model and the observations for both step 1 and step 2
19 (see Kuppel et al. 2013). For simplicity we assumed that the observation error covariances
20 were independent between the different observations and therefore we kept \mathbf{R} diagonal (off-
21 diagonal terms set to zero), given the rapid decline of the model error auto-correlation beyond
22 one day (Kuppel et al., 2013). For step 3 we used a different approach, given the large bias in
23 the model a priori concentrations, and therefore followed the methodology of Peylin et al.
24 (2005) based on the observed and modeled temporal concentration variability at each site.
25 Overall, data uncertainties in the optimization procedure are between 0.1 and 0.45 for NDVI
26 (step 1), around 3-6 $\text{gCm}^{-2}\text{d}^{-1}$ for daily NEE, and 15-30 Wm^{-2} for daily LE (step 2) and
27 between 0.1 ppm at remote oceanic stations and 4 ppm at continental sites (step 3).

28 The determination of the optimal parameter vector that minimizes $J(\mathbf{x})$ is performed by
29 successive calls to a “gradient-descent” minimization algorithm L-BFGS-B (Byrd et al.
30 1995), which is specifically dedicated to solving large nonlinear optimization problems that
31 are subject to simple bounds on the parameter values. In order to find the minimum of $J(\mathbf{x})$ the

1 algorithm requires the gradient of $J(\mathbf{x})$ (Jacobian) with respect to the ORCHIDEE parameters.
2 L-BFGS-B explores each parameter space simultaneously along the gradient of the cost
3 function, and uses an approximation of the Hessian (second derivative) of $J(\mathbf{x})$, which is
4 updated at each iteration, to define the size of the step at each iteration.

5 For step 1 and step 2, the model “ H ” simply corresponds to the land surface model: $H = S$,
6 with $S(x)$ representing the surface fluxes from the ORCHIDEE model using the parameter
7 vector, x . The gradients $dJ(x)/dx$ are calculated from the tangent linear model of ORCHIDEE
8 that was automatically generated by the numerical Transformation of Algorithms in Fortran
9 (www.fastopt.de), except for two parameters linked to the model phenology for which the
10 threshold functions prevent the use of a linear approximation. A finite difference approach
11 was used for these parameters in order to define a mean derivative at any point; we also
12 checked that no spurious oscillations occurred for these parameters during the minimization
13 iterations.

14 For step 3, the model “ H ” corresponds to the composition of the land surface model with the
15 transport model: $H = T \circ S$ (see Kaminski et al. (2002) for details), with T representing the
16 LMDz transport model. T is a linear operator for a non-reactive species: $T(S(\mathbf{x})) = \mathbf{T} \cdot S(\mathbf{x})$,
17 with \mathbf{T} a matrix representation of the transport operator. It corresponds to the sensitivity of
18 CO_2 concentrations at each site and for each month to the daily surface flux of each model
19 grid-cell. It is then combined with the ORCHIDEE surface fluxes ($S(\mathbf{x})$) through a matrix
20 multiplication to derive $H(\mathbf{x})$. \mathbf{T} has been pre-calculated for all atmospheric stations in order
21 to save computing time during the iterative optimization process (see §2.3.2). For simplicity
22 we use monthly mean values for both the fluxes $S(\mathbf{x})$ and the transport sensitivities (\mathbf{T}) in the
23 computation of the gradients $dJ(\mathbf{x})/d\mathbf{x}$.

24 For improved minimization efficiency, the inversion is preconditioned (following Chevallier
25 et al., 2005), which means that L-BFGS-B is fed with the control variable $\mathbf{x}' = \mathbf{B}^{-1/2}(\mathbf{x} -$
26 $\mathbf{x}_b)$, rather than with \mathbf{x} , as this homogenizes the range of variation of the optimized
27 parameters.

28 2.4.3 Error estimation

29 The posterior parameter error covariance matrix, \mathbf{A} , can be approximated to the inverse
30 Hessian of the cost function, using the linearity assumption at the minimum of $J(\mathbf{x})$. It can be

1 derived with the Jacobian of the model at the end of the minimization (i.e. the last iteration),
2 \mathbf{H}_∞ , following Tarantola (1987):

$$3 \quad \mathbf{A} = [\mathbf{H}_\infty^T \cdot \mathbf{R}^{-1} \cdot \mathbf{H}_\infty + \mathbf{B}^{-1}]^{-1} \quad (4)$$

4 Note that for step 3, $\mathbf{H}_\infty = \mathbf{T} \cdot \mathbf{S}_\infty$, where \mathbf{S}_∞ is the Jacobian of the ORCHIDEE model at the
5 last iteration. The posterior parameter error covariance, \mathbf{A} , can then be propagated into the
6 model state variable space (e.g. carbon fluxes and stocks), \mathbf{A}_{var} , given the following matrix
7 product (only used for the global fluxes in step 3):

$$8 \quad \mathbf{A}_{\text{var}} = \mathbf{S}_\infty \cdot \mathbf{A} \cdot \mathbf{S}_\infty^T \quad (5)$$

9 The square root of the diagonal elements of \mathbf{A}_{var} corresponds to the standard deviation, σ , of
10 carbon fluxes/stocks for each grid cell. In order to evaluate the knowledge improvement
11 brought by the assimilation, the uncertainty reduction between the prior (σ_{prior}) and posterior
12 (σ_{post}) is determined as $1 - (\sigma_{\text{post}} / \sigma_{\text{prior}})$.

13 2.4.4 Additional processing steps

14 In order to analyze the fit to the atmospheric CO₂ concentrations in terms of the trend and
15 seasonal cycle, we decomposed the observed and modeled time series by fitting the monthly
16 mean values with a function comprising a first order polynomial term and four harmonics,
17 and then filtered the residuals of that function in frequency space using a low pass filter
18 (cutoff frequency of 65 days), following Thoning et al. (1989). The polynomial term defines
19 the trend while the seasonal cycle corresponds to the harmonics plus the filtered residuals.
20 The amplitude of the seasonal cycle is then calculated as the difference between the monthly
21 mean maximum and minimum for year 2003 (middle year of the optimization period).
22 Finally, we define the Carbon Uptake Period (CUP) as the sum of the days when the values of
23 the seasonal cycle extracted from the CO₂ concentration time series are negative (a negative
24 convention being for CO₂ removed from the atmosphere).

25 2.5 Prescribed emissions of carbon fluxes

26 In this section we describe the other components of the carbon cycle (apart from the surface C
27 exchange with terrestrial vegetation) that are imposed in step 3 of the optimization process as
28 fixed fluxes.

1 2.5.1 Ocean fluxes

2 The ocean contributes to an uptake of about a quarter to a third of the anthropogenic
3 emissions with significant year-to-year variations (Sabine et al., 2004). For this version of the
4 ORCHIDAS, we developed a statistical model to estimate the spatial and temporal variations
5 (monthly) of the ocean surface CO₂ partial pressure (pCO₂^{SW}), and from that the air-sea CO₂
6 fluxes, using satellite and in-situ ocean measurements and model outputs. The air-sea CO₂
7 fluxes are primarily controlled by the ocean biogeochemistry, the horizontal transport and the
8 vertical mixing in the ocean and the atmospheric forcing (CO₂ partial pressure at the interface
9 to the water (pCO₂^{ATM}) and wind); they can be defined from the following equation:

$$10 F_{CO_2} = K_{ex} \times (pCO_2^{SW} - pCO_2^{ATM}) \quad (6)$$

11 where K_{ex} stands for the exchange coefficient and F_{CO_2} the CO₂ flux from the sea surface
12 water to the atmosphere.

13 The computation of pCO₂^{SW} is performed using feedforward artificial neural networks, i.e., a
14 MultiLayer Perceptron (MLP; Rosenblatt 1958) that maps a set of spatio-temporal variables
15 (input) onto observed pCO₂^{SW} data. We use a two-step approach: the first step to derive a
16 monthly mean pCO₂^{SW} climatology and the second step to correct for the year to year
17 variations. The pCO₂^{SW} observations come from the Global Surface pCO₂ (Lamont-Doherty
18 Earth Observatory, LDEO) Database (Takahashi et al., 2009). The inputs are a series of
19 variables connected to the spatial and temporal evolution of pCO₂^{SW}: i) sea surface
20 temperature (SST), sea surface salinity (SSS) and mixed layer depth (MLD) as a proxy of the
21 physical processes (these fields come from a re-analysis of the NEMO-OPA ocean model
22 (Madec et al., 1998) with the assimilation of several satellite observations), ii) chlorophyll
23 content from SeaWiFS, as a proxy of the biogeochemistry (CHL), iii) spatial and temporal
24 coordinates (LAT, LON and MONTH) and the pCO₂^{SW} at previous time step (recursive
25 approach), i.e.:

$$26 \{pCO_2^{SW}\}_m = MLP(\{SST, SSS, MLD, CHL\}_{(m-2, m-1, m)}, \{pCO_2^{SW}\}_{(m-2, m-1)}, LAT, LON) \quad (7)$$

27 with m the monthly index. The available data (20685 points) is divided into two parts: 75% is
28 used for the learning phase of the ANN and 25% for the validation phase. The overall
29 performance of the neural network for extrapolating the spatial and seasonal distribution of
30 pCO₂^{SW} is relatively good, with a spatio-temporal correlation coefficient between the
31 estimated pCO₂^{SW} and the independent observations of 0.80.

1 $p\text{CO}_2^{\text{ATM}}$ at the surface are taken from a global simulation of atmospheric CO_2 concentrations
2 with optimized fluxes (Chevallier et al. 2010). K_{ex} is defined as the product of k , the gas
3 transfer velocity, taken from the Wanninkhof (1992) formulation using winds from ERA-
4 Interim, and s , the solubility of CO_2 , taken from the Weiss formulation (Weiss, 1974). The
5 system is further described in Roedenbeck et al. (2015). The global ocean sink is around 1.60
6 PgC.yr^{-1} for the period 2002-2004 used in step 3. It is within the uncertainty range of the
7 Global Carbon Project estimates (Le Quéré et al., 2015) if we account for the pre-industrial
8 ocean outgazing flux included in our “delta $p\text{CO}_2$ ” approach. Its temporal evolution is
9 depicted in Fig. A1

10 2.5.2 Global fossil fuel and cement CO_2 emissions

11 We have used a recently developed CO_2 fossil fuel and cement emission product (see
12 <http://www.carbones.eu/wcmqs/>) that covers the period 1980 to 2009 at the spatial resolution
13 of $1^\circ \times 1^\circ$ and hourly resolution. It is based on EDGAR v4.2 spatially distributed annual
14 emissions (Olivier et al., 2012) and time profiles developed by the University of Stuttgart
15 (<http://carbones.ier.uni-stuttgart.de/wms/impressum.html>). It was assumed that EDGAR
16 delivers the most up-to-date spatially distributed and sector specific emissions, based on
17 national emission statistics. The IER (Institut für Energiewirtschaft und Rationelle
18 Energieanwendung) further applied country and sector specific time profiles, taking into
19 account monthly, daily, and hourly variations depending on the sector. The derivation of the
20 time profiles relies on different data sets (e.g. Eurostat, ENSTO-E
21 (<https://www.entsoe.eu/about-entso-e/Pages/default.aspx>), UN monthly bulletin) as well as
22 correlations between recorded emissions and climate variables. Currently, the temporal
23 profiles are derived mostly from data sets over Europe that were extrapolated using
24 information on climate zone, average monthly temperature for the seasonal cycles and
25 similarity in socio-economic parameters like population and Gross Domestic Product (GDP).
26 The annual mean emission for the period 2002-2004 is 7.14 PgC.yr^{-1} .

27 2.5.3 Fire emissions:

28 Fire emissions data from the Global Fire Data (GFEDv3 –
29 <http://www.globalfiredata.org/Data/index.html>) are prescribed in the ORCHIDAS (except
30 during the model spin-up). The GFEDv3 data are broken-down into 6 sectors (deforestation,
31 peat fires, savanna fires, agriculture, forest fires, and woodland) that are further grouped into

1 3 main types. We generated fluxes of CO₂ relevant for typical "burning - regrowth" processes,
2 as detailed in Appendix A2. The first type corresponds to deforestation and peat fires with
3 carbon permanently lost to the atmosphere, the second to agriculture and savannah fires which
4 are assumed to be compensated by a sink during the regrowth period (i.e. with zero annual net
5 emission for each pixel) and the third to woodland and burnt forests which are assumed to be
6 at steady state for a given region (10 sub-continental scale regions) over the period covered by
7 GFEDv2 (i.e. regrowth of nearby forest compensates for the burned forest derived in GFED).
8 The sum of these three components leads to the global flux, with a gross emission around 2.1
9 PgC.yr⁻¹ and a net emission after regrowth of only 1.1 PgC.yr⁻¹ (Fig. A2 in Appendix) that is
10 prescribed to the ORCHIDAS over the period 2002-2004.

11

12 **3 Results**

13 **3.1 Model fit to the data**

14 **3.1.1 Step 1: assimilation of MODIS NDVI data**

15 The optimization in Step 1 resulted in an improved fit to the MODIS NDVI observations for
16 the four PFTs considered (TeBD, BoND, BoBD, NC3, see §2.4) as seen in Fig. 4, which
17 shows the mean seasonal cycle across the 2000-2008 period for all sites for each PFT. The
18 most prominent change after the optimization was a substantially shorter growing season for
19 all PFTs due to an earlier start of leaf senescence. This was caused by both a lower critical
20 leaf age ($L_{agecrit}$) and a higher temperature threshold for senescence (CT_{senes}) (see Fig. 9). The
21 impact on the start of leaf growth was less dramatic but important nonetheless, with a shift to
22 a later start of leaf growth as a result of an increase in the $K_{pheno,crit}$ parameter which acts as a
23 scalar on the threshold of Growing Degree Days (GDD) used to trigger leaf onset (see
24 Appendix A in MacBean et al., 2015). Overall, a mean reduction in RMSE of 23, 17, 58 and
25 19% was achieved for TeBD, BoBD, BoND trees and NC3 grasses respectively, with the
26 greatest improvement for BoND trees. The mean correlation between the normalized MODIS-
27 NDVI and modeled FAPAR time series over the 2000 – 2008 period increased for TeBD and
28 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,
29 respectively). The prior correlation of 0.55 remained similar after the assimilation for BoBD
30 trees.

1 Following the improvement at the sites selected for the optimization, we evaluated the impact
2 for each PFT at the global scale using the global median correlation between the MODIS-
3 NDVI and the model FAPAR time series (from all pixels where the fraction of a given PFT is
4 above 60%, see Maignan et al. 2011). The global correlation increased for BoND trees and
5 NC3 grasses from 0.36 to 0.91 and 0.53 to 0.59 (prior to posterior), respectively. It remains
6 stable for BoBD (0.54) or slightly increased for TeBD (0.88 to 0.89).

7 3.1.2 Step 2: assimilation of FLUXNET data

8 The optimization in Step 2 brings an improvement to the simulated NEE and LE for all seven
9 PFTs considered, with Fig. 5 showing the corresponding PFT-averaged mean NEE seasonal
10 cycles (mean across all sites/years). NEE is overestimated by the prior model for all PFTs on
11 average. This is partly due to the model spin-up procedure, which brings each simulated site
12 to a near equilibrium state with a mean NEE close to zero (i.e. no net carbon sink, see §2.1.1).
13 This bias is significantly corrected by the optimization to match the observed carbon uptake at
14 most sites, notably via the scaling of the initial soil carbon pool content at each site
15 (parameters $K_{soilC,site}$; Table 1) which thus significantly reduces the ecosystem respiration
16 (Kuppel et al., 2014). Overall, the largest reductions of model-data RMSE are found in
17 temperate forests (TeNE, TeBE and TeBD), where the RMSE decreased by more than 25%
18 compared with the prior model. The improvements are less significant for the other PFTs,
19 with RMSE reductions between 10 and 18%.

20 In addition, the optimization increases the NEE seasonal amplitude in temperate evergreen
21 forests (TeNE and TeBE) and temperate broadleaf deciduous forests (TeBD), and reduces the
22 amplitude for boreal needle leaf forest (BoNE) and natural C3 grasses (NC3), in agreement
23 with the observations (except for BoNE where the amplitude decrease is too large). Despite
24 the better model-data agreement for evergreen broadleaf forests (TrBE and TeBE), the
25 optimized model still fails to catch some seasonal features such as a persistent carbon uptake
26 (i.e. negative NEE) in the dry season for the tropical regions (TrBE) and nearly-null carbon
27 exchange in the first months of the year for temperate regions (TeBE). These results are
28 discussed further in Kuppel et al. (2014), who used a similar optimization set-up with a
29 slightly different parameter set – see §2.3.3. Similar improvements, although of smaller
30 amplitude, occur for the latent heat fluxes (not shown).

1 3.1.3 Step 3: assimilation of atmospheric CO₂ data

2 The final optimization step with the atmospheric CO₂ concentrations provides a large
3 improvement of the fit to the observed concentrations at most stations. The cost function J
4 was reduced through the minimization by a factor of 5.7 within 37 iterations.

5 Figure 6 illustrates the simulated concentrations for four stations (representative of different
6 conditions), over the assimilation period, with the standard prior parameter vector (used in
7 step 1), the posterior vector from step 2 (used as prior in step 3) and the posterior vector from
8 this last step. The improvement in the fit to the observations can be quantified with the
9 reduction in RMSE (from the prior to the posterior of step 3) - the largest reduction is at the
10 South Pole station (73%) and is on average around 25% across all sites. Note that for a few
11 stations the fit is slightly degraded (up to 10%) except for one Pacific site (regular ship
12 measurements around the equator, POCN00) for which there is a 40% degradation, possibly
13 due to small biases in the simulation of the ITCZ position in LMDz. When calculated with
14 respect to the standard prior (used in step 1) the RMSE decrease is slightly larger on average,
15 especially for the northern mid to high latitude stations. For these stations the optimization
16 performed in step 2 with FLUXNET data led to a significant improvement of the mean
17 seasonal cycle amplitude of the atmospheric CO₂ data, as discussed in Kuppel et al. (2014).

18 We then investigated the fit to the observed CO₂ concentrations in terms of the mean seasonal
19 cycle and trend (see section 2.4.4). With only three years of data the mean trend is more
20 difficult to define as it varies between stations; however, the optimization in step 3 increases
21 the net land carbon sink in order to match the observed trend at most stations (as expected). If
22 we take the Mauna Loa and South Pole stations that are representative of an integration of the
23 fluxes at hemispheric scales, the prior CO₂ trend of 2.8 and 2.9 ppm.yr⁻¹ respectively, is
24 reduced to 2.1 and 2.2 ppm.yr⁻¹ close to the observations (2.1 ppm.yr⁻¹ for both). The left
25 panel of Fig. 7 illustrates changes in the amplitude of the simulated seasonal cycle at each
26 station (see definition in §2.4.4). The values correspond to relative changes between the prior
27 (of step 3) and posterior of the absolute difference between observed and modeled amplitude
28 ($(||\Delta A_{poste}| - |\Delta A_{prior}||)/|\Delta A_{prior}|$). They reveal an improvement in the seasonal cycle
29 amplitude at nearly all stations of the southern hemisphere ($\approx 40\%$ improvement) and at the
30 majority of the northern hemisphere stations ($\approx 15\%$). A few stations in north East Asia (3)
31 and northwest America (4) show a small degradation of the amplitude ($\approx 15\%$). The right
32 panel of Fig. 7 displays the changes of the Carbon Uptake Period (CUP, see §2.4.4) expressed

1 in terms of relative changes between prior and posterior of the absolute values of model-data
2 differences, as for the amplitude. Most stations reveal an improvement of the CUP of around
3 20%, which is slightly lower than the improvement for the seasonal cycle amplitude.

4 Finally, we verified that the improvement is not only valid at the optimization sites but also at
5 independent atmospheric CO₂ stations (see section 2.2.3). On average the mean RMSE for the 24
6 additional sites is 10.5 ppm for the prior of step 1 (prior of ORCHIDEE), 2.8 ppm for the prior or step
7 3 (i.e. posterior of step 2) and 2.1 ppm for the posterior of step 3. The corresponding values for the 53
8 sites used for the optimization are 10.5, 2.45 and 1.8 ppm, respectively. The error reduction during
9 step 3 is thus similar for both the assimilated and the validation data sets, further confirming that the
10 optimization provides a global improvement of the simulated carbon fluxes.

11 **3.2 Consistency of the step-wise optimization**

12 The main issue with a step-wise data assimilation system (versus a simultaneous approach)
13 concerns the potential degradation of the model – data fit for the different data streams that
14 are assimilated in previous steps. We noted that CO₂ concentrations were already improved
15 when NDVI and FLUXNET data are assimilated (see §3.1.3), but we need to check if the
16 final parameter set from step 3 leads to a degradation of the fit to MODIS-NDVI (step 1) and
17 to FLUXNET (step 2) data compared to the fit achieved during the respective steps and, in the
18 case of a significant degradation, if we still have an improvement for these data streams
19 compared to the initial *a priori* fit.

20 Figure 8 summarizes the performance of the model data fit for MODIS-NDVI and
21 FLUXNET-NEE data streams for the prior and posterior of each step by evaluating the
22 median RMSE between the model and the observations across all sites. The values are
23 calculated for each PFT separately. In this section, we keep in mind the fact that we do not
24 optimize the same PFTs with FLUXNET data and with MODIS-NDVI.

25 **Consistency for MODIS-NDVI**

26 First, we notice again the significant RMSE reduction between the prior and step 1, as
27 discussed in section 3.1. The fit to MODIS-NDVI (normalized data) for step 2 and step 3
28 shows only a significant degradation (increased RMSE) for temperate broadleaf deciduous
29 forest (TeBD), which decreases the improvement achieved in step 1 (compared to the prior)
30 by a factor of two. A marginal degradation for natural C3 grassland (NC3) is obtained after
31 step 3: the RMSE increases slightly from 0.24 to 0.26, but is still lower than the prior value of

1 0.3. There is no degradation for boreal needleleaf deciduous trees (BoND), but a surprising
2 small decrease of the RMSE (i.e. improvement in the model-data fit) for boreal broadleaf
3 deciduous forests (BoBD; from 0.26 to 0.23). In this latter case, the use of additional
4 parameters in steps 2 and 3 (see §2.4) allows further improvement of the fit between the
5 normalized FAPAR and NDVI time series. On average the degradation of the fit to NDVI is
6 thus very limited in step 2 and step 3, and in no case is the RMSE worse than the prior.

7 Consistency for FLUXNET data

8 Figure 8 again reveals the significant reduction of the RMSEs for NEE in step 2 compared to
9 the standard prior or to the posterior of step 1 for most PFTs, except BoNE. We see only
10 small degradations (increases) in RMSE between step 2 and step 3 for temperate needle leaf
11 evergreen forests (TeNE: from 1.06 to 1.13 $\text{gC.m}^2.\text{d}^{-1}$), temperate broadleaf evergreen forests
12 (TeBE: from 1.06 to 1.09 $\text{gC.m}^2.\text{d}^{-1}$), temperate broadleaf deciduous forests (TeBD: from 1.06
13 to 1.13 $\text{gC.m}^2.\text{d}^{-1}$) and boreal needle leaf evergreen forests (BoNE: from 0.59 to 0.60 $\text{gC.m}^2.\text{d}^{-1}$).
14 An interesting feature to notice is that the NEE RMSE increases from the prior to the
15 posterior of step 1 (i.e. before NEE has been used in the optimization in step 2). Using remote
16 sensing products of vegetation activity or “greenness” (e.g. NDVI) to calibrate the phenology
17 of ORCHIDEE thus does not always improve the simulated NEE, as they only provide a
18 strong constraint on the timing of the leaf phenology (and thus indirectly the GPP) and a weak
19 constraint on the maximum GPP but no constraint on the respiration fluxes. These reasons
20 were discussed in Bacour et al. (2015) who used the same LSM and assimilation system.
21 Overall, the reduction of the improvement of the model data fit to the NEE (step 3 versus step
22 2) is marginal (limited to a few percent), thus further suggesting the consistency of our step-
23 wise approach. Similar results are also obtained for the latent heat flux (LE) (not shown).

24 3.3 Estimated parameter values

25 We now discuss the parameter values, focusing on the changes obtained through the
26 successive steps. Figure 9 presents the prior and posterior values for each parameter together
27 with their associated uncertainties (estimated through Eq. (4)) and the allowed range of
28 variation. Note that nine parameters are PFT-dependent while four are global (non PFT-
29 dependent). For the global non PFT-dependent parameters included in the step 2 optimization,
30 we took the mean value and error-variance (see §2.4) as the prior for step 3. Note finally that

1 the parameters linked to the initial soil carbon pools ($K_{soilC,site}$, $K_{soilC,reg}$) are not shown in Fig.
2 9 as they are too numerous (though see Fig. A2 for the regional values).

3 If we first consider the phenology parameters optimized in step 1 ($K_{lai,happy}$, $K_{pheno,crit}$, L_{age_crit} ,
4 $C_{T,senes}$; see Table 1) we see that for most PFTs they do not change significantly between step
5 1 and step 3, although they differ significantly from the prior. There are few exceptions,
6 including $K_{pheno,crit}$ (the threshold for the start of the growing season) for Boreal Needleleaf
7 deciduous forests and $K_{lai,happy}$ (level of carbohydrate use) for temperate and boreal broadleaf
8 deciduous forests (TeBD, BoBD). Note that a few phenology parameters hit one of the
9 physical bounds, which may indicate model structural errors or model parameter equifinality.
10 For most phenology parameters, the uncertainties are strongly reduced during their first
11 optimization (step 1), except for a few cases like $C_{T,senes}$ for C3 grassland. Note finally that a
12 more in depth spatio-temporal validation demonstrated the generality of the optimized
13 phenology parameters across multiple sites (for further details see MacBean et al., 2015).

14 For the photosynthesis parameters (V_{cmax} , $G_{s,slope}$, C_{Topt} , SLA , f_{stress} ; see Table 1), we find a
15 similar result with little changes between step 2 and step 3, but still a significant departure
16 from the prior values. Most parameters are well constrained by the inversion, with posterior
17 uncertainties that are greatly reduced compared to the prior, except for Tropical broadleaf
18 rain-green forest (TrBR) and Boreal needle-leaf deciduous forest (BoND) for which there is
19 nearly no constraint on $G_{s,slope}$, and f_{stress} (see Table 1).

20 The non-PFT dependent respiration-related parameters ($HR_{H,c}$, Q_{10} , MR_b) mostly change in
21 step 2 and only slightly in step 3 (with an additional reduction of the error) in order to fit the
22 large-scale constraint provided by the atmospheric observations. The values of the scalar of
23 the initial soil carbon pools size for the FLUXNET site optimizations ($K_{soilC,site}$, one parameter
24 per site, not shown) were largely reduced on average, in order to decrease the heterotrophic
25 respiration (see Kuppel et al. (2014) for additional discussion). In step 3 the same scalars that
26 were defined for an ensemble of large regions ($K_{soilC,reg}$) have decreased in the southern
27 hemisphere (less than 10%; see Fig. A2 in Appendix A3) and slightly increased in the
28 northern hemisphere (around 1%), to achieve a better match to the atmospheric CO₂ growth
29 rate and north-south gradient. Importantly, we notice that for step 3, the fit to the atmospheric
30 CO₂ concentrations (especially to the trend) is achieved mainly by small changes in $K_{soilC,reg}$
31 and in few other respiration-related parameters. Note finally that the parameter controlling the
32 albedo ($K_{albedo,veg}$), modified with the FLUXNET observations only (see §2.4), is not well

1 constrained by the optimization (only a small reduction in uncertainty). Overall, most
2 parameters appear to be well constrained when first optimized, with only small changes in the
3 following steps. This suggests that the targeting of different parameter subspaces in the
4 various optimisation steps was well-chosen.

5 **3.4 Estimated carbon fluxes and uncertainties**

6 The main objective of a carbon cycle data assimilation procedure is to improve the simulated
7 land surface net and gross carbon fluxes as well as the simulated carbon stocks for both
8 present and future conditions. Given the focus of the paper, i.e. to describe the potential of a
9 step-wise global carbon cycle data assimilation system, we only discuss a few large-scale
10 features of the optimized annual net and gross carbon fluxes, as well as one of the carbon
11 stock variables (forest above-ground biomass). We thus do not discuss the inter-annual flux
12 variability.

13 **Large-scale annual mean net fluxes**

14 The mean annual carbon fluxes (NEE) for the globe, northern extra tropics, tropics, and
15 southern extra tropics are reported in Fig. 10 for the 2000-2009 decade for the prior and
16 posterior model simulations for all steps. We ran the optimized model over the full 2000's
17 decade in order to compare with one other estimate of the land surface residual from the
18 Global Carbon Project (GCP, Le Quéré et al, 2015) over the same decade. The prior NEE
19 indicates a total sink of 0.5 PgC.yr^{-1} over this period, from both the northern and tropical
20 regions. Such a prior sink is due to the increase of atmospheric CO_2 during the transient
21 simulation following the spin-up (1990-2009, see section 2.3.1) and climate variability.
22 Changes from the prior are rather small in step 1 (assimilation of MODIS NDVI) with an
23 increase of the northern sink by 0.12 PgC.yr^{-1} and a decrease of the tropical sink by 0.05
24 PgC.yr^{-1} (Fig. 10). Step 2 (assimilation of FLUXNET data) does not significantly change the
25 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

1 (natural and anthropogenic), since that we used atmospheric CO₂ 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.m².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.m².y⁻¹ and
19 posterior values between 22 and 66 gC.m².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.m².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 172, 155, 156 and 157 PgC.yr⁻¹ for the prior and posterior of step 1,
25 2 and 3, respectively. The small overall decrease (9%) 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 $^{18}\text{O}/^{16}\text{O}$ 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 ($\sim 180 \text{ PgC}$) 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 ($\sim 120 \text{ PgC}$).

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 related data streams, namely MODIS-NDVI
30 observations of vegetation activity at 60 sites, FLUXNET NEE and LE measurements at more
31 than 70 sites, and atmospheric CO_2 concentrations at 53 surface stations over three years in

1 order to optimize the C cycle parameters of the ORCHIDEE process-based LSM
2 (ORCHIDEE-CCDAS). The study details the concept, the implementation and the main
3 results of a stepwise assimilation approach where the data streams have been assimilated in
4 three successive steps (including a propagation of the retrieved posterior parameter
5 distributions 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 CO₂ 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 CO₂ 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 CO₂
24 growth rate and thus to derive a land C sink compatible with current best estimates from the
25 GCP. The assimilation of CO₂ data also slightly changed the seasonality of the NEE, which
26 improved the fit to the atmospheric CO₂ seasonal cycle. Note that assimilating only CO₂ 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 CO₂ 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 more simple approach), and not the full error variance-
8 covariance matrix, we were able to propagate enough information to maintain an optimal
9 model-data fit after the last step for the three data streams, as confirmed with the back-
10 compatibility checks. MacBean et al. (2016) provide a more specific analysis of this issue.
11 However, not propagating the covariance terms may have a larger impact for the reduction of
12 the inferred parameter uncertainties (see for instance the large parameter / flux error reduction
13 in Fig. 9 / Fig. 11). The order of the different steps was dictated by the number of parameters
14 we choose to expose to each data stream, from only a few phenology parameters for NDVI up
15 to the largest set for atmospheric CO₂. Recall that under the fundamental theory the order of
16 assimilation is unimportant. Testing whether our system meets this criterion is an important
17 check on the robustness of the method but is not technically feasible with the full-blown
18 system; it is currently being tested with some smaller 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 CO₂ 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

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 CO₂ 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 CO₂ 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 Bloom et al., 2016 and Thum et al., in revision for AFM). Additionally, uncertainties on
19 the other components of the carbon cycle, such as fossil fuel and biomass burning emissions
20 and ocean fluxes, can be also critical when using atmospheric CO₂ as a constraint. Finally,
21 one can mention new approaches based on remote sensing data to account for site level
22 differences in productivity potential due to edaphic variability (soil quality and
23 slope/drainage) within the same vegetation type, as illustrated for high latitudes in North
24 America (Ise and Sato, 2008).

25 To conclude, this work is a step forward in terms of multiple data streams assimilation that
26 opens new perspectives for a better understanding of the carbon cycle and better predictions
27 of the fate of the land carbon sink in the 21st century as a consequence of anthropogenic
28 changes. As ORCHIDEE is part of the IPSL earth system model the impact of the
29 optimization on future climate change predictions will be assessed in a future study. However,
30 we first need to run the ORCHIDAS with a longer atmospheric CO₂ record (i.e. several
31 decades) in order to provide stronger constraints on parameters controlling the impact of
32 climate extremes on the net carbon fluxes at continental to global scales, and the sensitivity of

1 photosynthesis to increasing CO₂ concentration. The optimized model will allow a more in-
2 depth investigation of the trend and inter-annual variations of land surface C fluxes at
3 continental to regional scales, as well as their driving mechanisms. It will offer a more
4 reliable and robust process-based diagnostic of the land C cycle that is compatible with
5 current major data streams. Overall, we have illustrated the benefit of combining multiple
6 data streams in a process-based model to optimize different processes of the model, related to
7 different temporal and spatial scales. The optimization will be updated regularly as new
8 processes are integrated into the ORCHIDEE model, such as for instance land management
9 (Naudts et al., 2015).

10

11 **Code availability**

12 The ORCHIDEE model code and the run environment are open source
13 (<http://forge.ipsl.jussieu.fr/orchidee>) and the associated documentation can be found at
14 <https://forge.ipsl.jussieu.fr/orchidee/wiki/Documentation>. Note that the tangent linear version
15 of the ORCHIDEE model has been generated using commercial software (TAF;
16 <http://www.fastopt.com/products/taf/taf.shtml>). For this reason, only the “forward” version of
17 the ORCHIDEE model is available for sharing. The optimization scheme (in Python) is
18 available through a dedicated web site for data assimilation with ORCHIDEE
19 (<http://orchidas.lsce.ipsl.fr/>). Nevertheless readers interested in running ORCHIDEE are
20 encouraged to contact the corresponding author for full details and latest bug fixes. Finally,
21 the source code of the LMDZ atmospheric transport model can be found at
22 <http://web.lmd.jussieu.fr/trac>.

23

24 **Appendix**

25 **A1. Ocean fluxes**

26 Figure A1 displays the air-sea fluxes from the statistical model.

27 **A2. Fire fluxes**

28 In order to account for fundamental differences between six fire flux categories provided by
29 the GFED product, we grouped these emissions into 3 types with specific treatments. The first

1 group includes C emissions from deforestation and peat fires, which are considered to be
2 permanent carbon lost to the atmosphere, at least over the considered time scales. These
3 fluxes are rescaled to an annual emission of 1.1 PgC.yr^{-1} globally following typical values
4 reported in the literature for deforestation (Houghton R., 2003). The second group consists of
5 C emissions from agriculture and savannah fires, which are compensated by a C sink during
6 the regrowth of these biomes (i.e., savannah and some type of plants on the farmland). These
7 effects are not completely accounted for in ORCHIDEE as the model does not simulate
8 savannah and agriculture fire. Hence, the emissions over the whole period and for each pixel
9 become zero, but their seasonal variations are used. The final group includes emissions from
10 woodland and burnt forests. We considered that at steady state and for a given region certain
11 forests burn but that nearby forests are re-growing following older fires. We thus imposed
12 regrowth at the region scale given that the ORCHIDEE model version that we use does not
13 account for such regrowth. The main assumption is that over century time scale the
14 forest/woodland system is at steady state over a given region (few thousand square km), i.e.
15 there is no net deforestation. We selected an ensemble of small regions over which we
16 calculated the regrowth of these biomes. The derived emissions over the whole period and for
17 each region thus become zero; though we include their spatial and temporal variations. The
18 overall biomass burning flux considered in the CCDAS for the optimization process is the
19 sum of the three fluxes as described above.

20 **A3. Multipliers of the soil initial carbon pools**

21 Figure A2 provides the optimized values of the $K_{soilC,reg}$ parameters that optimize the initial
22 soil carbon pool sizes.

23 .

24 **Acknowledgements**

25 This work was mainly funded by the EU FP7 CARBONES project (contracts FP7-SPACE-
26 2009-1-242316), with also a small contribution from GEOCARBON project
27 (ENV.2011.4.1.1-1-283080). This work used eddy covariance data acquired by the
28 FLUXNET community and in particular by the following networks: AmeriFlux (U.S.
29 Department of Energy, Biological and Environmental Research, Terrestrial Carbon Program
30 (DE-FG02-04ER63917 and DE-FG02-04ER63911)), AfriFlux, AsiaFlux, CarboAfrica,
31 CarboEuropeIP, CarboItaly, CarboMont, ChinaFlux, Fluxnet-Canada (supported by CFCAS,
32 NSERC, BIOCAP, Environment Canada, and NRCan), GreenGrass, KoFlux, LBA, NECC,

1 OzFlux, TCOS-Siberia, USCCC. We acknowledge the financial support to the eddy
2 covariance data harmonization provided by CarboEuropeIP, FAO-GTOS-TCO, iLEAPS, Max
3 Planck Institute for Biogeochemistry, National Science Foundation, University of Tuscia,
4 Université Laval and Environment Canada and US Department of Energy and the database
5 development and technical support from Berkeley Water Center, Lawrence Berkeley National
6 Laboratory, Microsoft Research eScience, Oak Ridge National Laboratory, University of
7 California-Berkeley, University of Virginia. P. C. acknowledges support from the European
8 Research Council through Synergy grant ERC-2013-SyG-610028 “IMBALANCE-P ». The
9 MODIS MOD09CMG collection 5 surface reflectance data are freely available to download
10 from the Land Processes Distributed Active Archive Center (LP DAAC) data portal
11 (<https://lpdaac.usgs.gov>). The authors wish to thank M. Jung for providing access to the GPP
12 MTE data, which were downloaded from the GEOCARBON data portal ([https://www.bgc-
13 jena.mpg.de/geodb/projects/Data.php](https://www.bgc-jena.mpg.de/geodb/projects/Data.php)). The authors are also grateful to computing support and
14 resources provided at LSCE and to the overall ORCHIDEE project that coordinate the
15 development of the code (<http://labex.ipsl.fr/orchidee/index.php/about-the-team>).

16

17

18

19

20

21

22

23

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., Lewis, S. L., Phillips, O. L., Asner, G. P., ... &
6 Berry, N. J., An integrated pan-tropical biomass map using multiple reference datasets. *Global
7 change biology*, 2016.
- 8 Bacour, C., Peylin, P., MacBean, N., Rayner, P. J., Delage, F., Chevallier, F., Weiss, M.,
9 Demarty, J., Santaren, D., Baret, F., Berveiller, D., Dufrêne, E. and Prunet, P.: Joint
10 assimilation of eddy covariance flux measurements and FAPAR products over temperate
11 forests within a process-oriented biosphere model, *J. Geophys. Res. Biogeosciences*, 120, 1–
12 19, doi:10.1002/2015JG002966, 2015.
- 13 Bloom, A. A., Exbrayat, J.-F., van der Velde, I. R., Feng, L. and Williams, M.: The decadal
14 state of the terrestrial carbon cycle: Global retrievals of terrestrial carbon allocation, pools,
15 and residence times., *Proc. Natl. Acad. Sci. U. S. A.*, 113(5), 1285– 1290,
16 doi:10.1073/pnas.1515160113, 2016.
- 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, A limited memory algorithm for bound
25 constrained optimization, *SIAM J. Sci. Stat. Comput.*, 16(5), 1190–1208, 1995.
- 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](http://www.biogeosciences-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 Inferring CO₂ sources and sinks from satellite observations: Method and application to TOVS
31 data, *Journal of Geophysical Research*, 110, D24309, doi:20.1029/2005JD006390, 13pp,
32 2005.
- 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, CO₂ surface fluxes at grid point

1 scale estimated from a global 21-year reanalysis of atmospheric measurements, *Journal of*
2 *Geophysical Research*, **115**, D21307, doi:10.1029/2010JD013887, 2010.

3 Collatz GJ, Ribas-Carbo M, Berry JA, (1992), Coupled Photosynthesis-Stomatal Conductance
4 Model for Leaves of C4 Plants. *Aust J Plant Physiol*, **19**: 519-38.

5 de Rosnay P, Polcher J. (1998), Modelling root water uptake in a complex land surface
6 scheme coupled to a GCM. *Hydrol Earth Syst Sc*, **2**: 239-55.

7 Dee, D. P., Uppala, S. M., Simmons, A. J., Berrisford, P., Poli, P., Kobayashi, S., Andrae, U.,
8 Balmaseda, M. A., Balsamo, G., Bauer, P., Bechtold, P., Beljaars, A. C. M., van de Berg, L.,
9 Bidlot, J., Bormann, N., Delsol, C., Dragani, R., Fuentes, M., Geer, A. J., Haimberger, L.,
10 Healy, S. B., Hersbach, H., Hólm, E. V., Isaksen, L., Kållberg, P., Köhler, M., Matricardi, M.,
11 McNally, A. P., Monge-Sanz, B. M., Morcrette, J. J., Park, B. K., Peubey, C., de Rosnay, P.,
12 Tavolato, C., Thépaut, J. N., and Vitart, F.: The ERA-interim reanalysis: configuration and
13 performance of the data assimilation system, *Q. J. Roy. Meteor. Soc.*, **137**, 553–597,
14 doi:10.1002/qj.828, 2011.

15 Ducoudre NI, Laval K, Perrier A. Sechiba (1993), A New Set of Parameterizations of the
16 Hydrologic Exchanges at the Land Atmosphere Interface within the Lmd Atmospheric
17 General-Circulation Model. *J Climate*, **6**: 248-73.

18 Dufresne, J. L., Foujols, M. A., Denvil, S., Caubel, A., Marti, O., Aumont, O., ... & Bony, S.
19 (2013). Climate change projections using the IPSL-CM5 Earth System Model: from CMIP3
20 to CMIP5. *Climate Dynamics*, **40**(9-10), 2123-2165.

21 Farquhar, G. D., von Caemmerer, S. von and Berry, J. A.: A biochemical model of
22 photosynthetic CO₂ assimilation in leaves of C₃ species, *Planta*, **149**(1), 78–90, 1980.

23 Folberth G., Hauglustaine D.A., Ciais P., et al. (2005), On the role of atmospheric chemistry
24 in the global CO₂ budget, *Geophysical Research Letters*, **32**(8): L08801

25 Fung, I. Y., C. J. Tucker, and K. C. Prentice, Application of Advanced Very High Resolution
26 Radiometer vegetation index to study atmosphere – biosphere exchange of CO₂, *J. Geophys.*
27 *Res.*, **92**, 2999– 3015, 1987.

28 GLOBALVIEW : Cooperative Global Atmospheric Data Integration Project. 2013, updated
29 annually. Multi-laboratory compilation of synchronized and gap-filled atmospheric carbon
30 dioxide records for the period 1979-2012 (obspack_co2_1_GLOBALVIEW-
31 CO2_2013_v1.0.4_2013-12-23). Compiled by NOAA Global Monitoring Division: Boulder,
32 Colorado, U.S.A. Data product accessed at <http://dx.doi.org/10.3334/OBSPACK/1002>.

33 Groenendijk, M., Dolman, a. J., van der Molen, M. K., Leuning, R., Arneth, a., Delpierre,
34 N., Gash, J. H. C., Lindroth, a., Richardson, a. D., Verbeeck, H. and Wohlfahrt, G.:
35 Assessing parameter variability in a photosynthesis model within and between plant
36 functional types using global Fluxnet eddy covariance data, *Agric. For. Meteorol.*, **151**(1),
37 22–38, doi:10.1016/j.agrformet.2010.08.013, 2011.

- 1 Hauglustaine D.A., Hourdin F., Jourdain L., et al. (2004), Interactive chemistry in the
2 Laboratoire de Meteorologie Dynamique general circulation model: Description and
3 background tropospheric chemistry evaluation, *Journal of Geophysical Research -*
4 *Atmosphere*, **109**(D4): D04314
- 5 Houghton, R. A. (2003) Revised estimates of the annual net flux of carbon to the atmosphere
6 from changes in land use and land management 1850-2000. *Tellus* **55B**: 378-390.
- 7 Hourdin F. and Armengaud A. (1999), The use of finite-volume methods for atmospheric
8 advection of trace species. Part I: Test of various formulations in a general circulation
9 model, *Monthly Weather Review*, **127**(5): 822-837
- 10 Hourdin, F. et al. (2006), The LMDZ4 general circulation model: climate performance and
11 sensitivity to parametrized physics with emphasis on tropical convection, *Climate*
12 *Dynamics*, **27**(7), 787-813, 2006.
- 13 IPCC, (2007). Climate Change 2007: The Physical Science Basis. Contribution of Working
14 Group I to the Fourth Assessment Report of the Intergovernmental Panel on Climate Change
15 [Solomon, S., D. Qin, M. Manning (eds.)]
- 16 Ise, T., and H. Sato, Representing subgrid-scale edaphic heterogeneity in a largescale
17 ecosystem model: A case study in the circumpolar boreal regions, *Geophys. Res. Lett.*, **35**,
18 L20407, doi:10.1029/2008GL035701, 2008.
- 19 Kaminski, T., Knorr, W., Scholze, M., Gobron, N., Pinty, B., Giering, R. and Mathieu, P-P.
20 (2012), Consistent assimilation of MERIS FAPAR and atmospheric CO₂ into a terrestrial
21 vegetation model and interactive mission benefit analysis, *Biogeosciences*, **9**, 3173-3184.
- 22 Kaminski, T., Knorr, W., Schürmann, G., Scholze, M., Rayner, P. J., Zaehle, S., Blessing, S.,
23 Dorigo, W., Gayler, V., Giering, R., Gobron, N., Grant, J. P., Heimann, M., Hooker-Stroud,
24 a., Houweling, S., Kato, T., Kattge, J., Kelley, D., Kemp, S., Koffi, E. N., Köstler, C.,
25 Mathieu, P. P., Pinty, B., Reick, C. H., Rödenbeck, C., Schnur, R., Scipal, K., Sebald, C.,
26 Stacke, T., Van Scheltinga, a. T., Vossbeck, M., Widmann, H. and Ziehn, T.: The
27 BETHY/JSBACH Carbon Cycle Data Assimilation System: Experiences and challenges, *J.*
28 *Geophys. Res. Biogeosciences*, **118**(4), 1414–1426, doi:10.1002/jgrg.20118, 2013.
- 29 Kato, T., Knorr, W., Scholze, M., Veenendaal, E., Kaminski, T., Kattge, J. and Gobron, N.:
30 Simultaneous assimilation of satellite and eddy covariance data for improving terrestrial water
31 and carbon simulations at a semi-arid woodland site in Botswana, *Biogeosciences*, **10**(2),
32 789–802, doi:10.5194/bg-10-789-2013, 2013.
- 33 Keenan, T. F., Davidson, E. a., Munger, J. W. and Richardson, A. D.: Rate my data:
34 Quantifying the value of ecological data for the development of models of the terrestrial
35 carbon cycle, *Ecol. Appl.*, **23**(1), 273–286, doi:10.1890/12-0747.1, 2013.
- 36 Keenan, T. F., Davidson, E., Moffat, A. M., Munger, W. and Richardson, A. D.: Using
37 model-data fusion to interpret past trends, and quantify uncertainties in future projections, of
38 terrestrial ecosystem carbon cycling, *Glob. Change Biol.*, **18**(8), 2555–2569,
39 doi:10.1111/j.1365-2486.2012.02684.x, 2012.

- 1 Knorr, W. and Kattge, J.: Inversion of terrestrial ecosystem model parameter values against
2 eddy covariance measurements by Monte Carlo sampling, *Glob. Change Biol.*, 11(8), 1333–
3 1351, doi:10.1111/j.1365-2486.2005.00977.x, 2005.
- 4 Knyazikhin, Y., Martonchik, J.V., Myneni, R.B., Diner, D.J., and Running, S.W. (1998),
5 Synergistic algorithm for estimating vegetation canopy leaf area index and fraction of
6 absorbed photosynthetically active radiation from MODIS and MISR, *Journal of Geophysical*
7 *Research*, **103**, D24, 32,257-32,276.
- 8 Koffi, E. N., Rayner, P. J., Scholze, M., & Beer, C. (2012). Atmospheric constraints on gross
9 primary productivity and net ecosystem productivity: Results from a carbon-cycle data
10 assimilation system. *Global biogeochemical cycles*, 26(1).
- 11 Krinner, G., Viovy, N., de Noblet-Ducoudré, N., Ogée, J., Polcher, J., Friedlingstein, P.,
12 Ciais, P., Sitch, S. and Prentice, I. C.: A dynamic global vegetation model for studies of the
13 coupled atmosphere-biosphere system, *Glob. Biogeochem. Cycles*, 19(1) [online] Available
14 from: <http://onlinelibrary.wiley.com/doi/10.1029/2003GB002199/pdf>, 2005.
- 15 Kuppel, S., Peylin, P., Chevallier, F., Bacour, C., Maignan, F. and Richardson, A. D.:
16 Constraining a global ecosystem model with multi-site eddy-covariance data, *Biogeosciences*,
17 9(10), 3757–3776, doi:10.5194/bg-9-3757-2012, 2012.
- 18 Kuppel, S., Chevallier, F. and Peylin, P.: Quantifying the model structural error in carbon
19 cycle data assimilation systems, *Geosci. Model Dev.*, 6(1), 45–55, doi:10.5194/gmd-6-45-
20 2013, 2013.
- 21 Kuppel, S., Peylin, P., Maignan, F., Chevallier, F., Kiely, G., Montagnani, L. and Cescatti, A.:
22 Model–data fusion across ecosystems: from multisite optimizations to global simulations,
23 *Geosci. Model Dev.*, 7(6), 2581–2597, 2014.
- 24 Lasslop, G., M. Reichstein, D. Papale, A. D. Richardson, A. Arneth, A. Barr, P. Stoy, and G.
25 Wohlfahrt. 2010. Separation of net ecosystem exchange into assimilation and respiration
26 using a light response curve approach: critical issues and global evaluation. *Global Change*
27 *Biology*, **16**, 187-208.
- 28 Lasslop, G., M. Reichstein, J. Kattge, and D. Papale. 2008. Influences of observation errors in
29 eddy flux data on inverse model parameter estimation. *Biogeosciences*, **5**, 1311-1324.
- 30 Le Quéré, C., Moriarty, R., Andrew, R. M., Peters, G. P., Ciais, P., Friedlingstein, P., Jones,
31 S. D., Sitch, S., Tans, P., Arneth, a., Boden, T. a., Bopp, L., Bozec, Y., Canadell, J. G., Chini,
32 L. P., Chevallier, F., Cosca, C. E., Harris, I., Hoppema, M., Houghton, R. a., House, J. I., Jain,
33 a. K., Johannessen, T., Kato, E., Keeling, R. F., Kitidis, V., Klein Goldewijk, K., Koven, C.,
34 Landa, C. S., Landschützer, P., Lenton, a., Lima, I. D., Marland, G., Mathis, J. T., Metzl, N.,
35 Nojiri, Y., Olsen, a., Ono, T., Peng, S., Peters, W., Pfeil, B., Poulter, B., Raupach, M. R.,
36 Regnier, P., Rödenbeck, C., Saito, S., Salisbury, J. E., Schuster, U., Schwinger, J., Séférian,
37 R., Segsneider, J., Steinhoff, T., Stocker, B. D., Sutton, a. J., Takahashi, T., Tilbrook, B.,
38 van der Werf, G. R., Viovy, N., Wang, Y.-P., Wanninkhof, R., Wiltshire, a. and Zeng, N.:
39 Global carbon budget 2014, *Earth Syst. Sci. Data*, 7(1), 47–85, doi:10.5194/essd-7-47-2015,
40 2015.

- 1 Liss, P. and Merlivat, L. (1986). The role of sea-air exchange in geochemical cycling, Ed. P.
2 Menard, chapter Air-sea gas exchange rates: Introduction and synthesis, pages 113-127.
3 Reidel, Dordrecht.
- 4 MacBean, N., Maignan, F., Peylin, P., Bacour, C., François-Marie, B. and Ciais, P.: Using
5 satellite data to improve the leaf phenology of a global Terrestrial Biosphere Model,
6 *Biogeosciences*, 12, 7185-7208, 2015.
- 7 MacBean, N., Peylin, P., Chevallier, F., Scholze, M. and G. Schürmann, Consistent
8 assimilation of multiple data streams in a Carbon Cycle Data Assimilation System,
9 *Biogeosciences*, in revision, 2016.
- 10 Madec, G., P. Delecluse, M. Imbard and C. Lévy, 1998 : OPA 8.1 Ocean General Circulation
11 Model reference manual, Note du Pole de Modelisation, Institut Pierre-Simon Laplace, 11,
12 91pp.
- 13 Maignan, F., Bréon, F.-M., Chevallier, F., Viovy, N., Ciais, P., Garrec, C., Trues, J., and
14 Mancip, M.: Evaluation of a Global Vegetation Model using time series of satellite vegetation
15 indices, *Geosci. Model Dev.*, 4, 1103-1114, doi:10.5194/gmd-4-1103-2011, 2011.
- 16 Moore, D. J. P., Hu, J., Sacks, W. J., Schimel, D. S. and Monson, R. K.: Estimating
17 transpiration and the sensitivity of carbon uptake to water availability in a subalpine forest
18 using a simple ecosystem process model informed by measured net CO₂ and H₂O fluxes,
19 *Agric. For. Meteorol.*, 148(10), 1467–1477, doi:10.1016/j.agrformet.2008.04.013, 2008.
- 20 Naudts, K., Ryder, J., J McGrath, M., Otto, J., Chen, Y., Valade, A., ... & Ghattas, J. (2015).
21 A vertically discretised canopy description for ORCHIDEE (SVN r2290) and the
22 modifications to the energy, water and carbon fluxes. *Geoscientific Model Development*, 8,
23 2035-2065.
- 24 Nightingale, P.D., et al. 2000. In situ evaluation of air-sea gas exchange parameterizations
25 using novel conservative and volatile tracers. *Glob. Biogeochem Cycles*, **14**, 373-387.
- 26 Olson, J., Watts, J.A., and Allison, L.J. (1983), Carbon in Live Vegetation of Major World
27 Ecosystems, ORNL-5862, Oak Ridge National Laboratory, Oak Ridge, Tennessee, 164pp.
- 28 Papale, D. (2006), Towards a standardized processing of Net Ecosystem Exchange measured
29 with eddy covariance technique: algorithms and uncertainty estimation, [online] Available
30 from: <http://dspace.unitus.it/handle/2067/1321>
- 31 Parton W, Stewart J, Cole C, (1988), Dynamics of C, N, P and S in grassland soils: a model.
32 *Biogeochemistry*, **5**: 109-31.
- 33 Peylin P, Rayner PJ, Bousquet P, et al. (2005), Daily CO₂ flux estimates over Europe from
34 continuous atmospheric measurements: 1, inverse methodology, *Atmospheric Chemistry and*
35 *Physics*, **5**: 3173-3186.
- 36 Piao, S., Sitch, S., Ciais, P., Friedlingstein, P., Peylin, P., Wang, X., Ahlström, A., Anav, A.,
37 Canadell, J. G., Cong, N., Huntingford, C., Jung, M., Levis, S., Levy, P. E., Li, J., Lin, X.,

- 1 Lomas, M. R., Lu, M., Luo, Y., Ma, Y., Myneni, R. B., Poulter, B., Sun, Z., Wang, T., Viovy,
2 N., Zaehle, S. and Zeng, N.: Evaluation of terrestrial carbon cycle models for their response to
3 climate variability and to CO₂ trends, *Glob. Change Biol.*, 19(7), 2117–2132,
4 doi:10.1111/gcb.12187, 2013.
- 5 Prentice, I. C., Liang, X., Medlyn, B. E. and Wang, Y.-P.: Reliable, robust and realistic: the
6 three R's of next-generation land-surface modelling, *Atmospheric Chem. Phys.*, 15(10),
7 5987–6005, doi:10.5194/acp-15-5987-2015, 2015.
- 8 Raoult, N. M., Jupp, T. E., Cox, P. M., and Luke, C. M.: Land surface parameter optimization
9 through data assimilation: the adJULES system, *Geosci. Model Dev.*, 9, 2833-2852,
10 doi:10.5194/gmd-9-2833-2016, 2016.
- 11 Rayner, P. J., Scholze, M., Knorr, W., Kaminski, T., Giering, R. and Widmann, H.: Two
12 decades of terrestrial carbon fluxes from a carbon cycle data assimilation system (CCDAS),
13 *Global Biogeochemical Cycles*, 19, doi:10.1029/2004GB002254, 2005.
- 14 Rayner, P. J. (2010). The current state of carbon-cycle data assimilation. *Current Opinion in*
15 *Environmental Sustainability*, 2(4), 289-296.
- 16 Regnier, P., Friedlingstein, P., Ciais, P., Mackenzie, F. T., Gruber, N., Janssens, I. A., ... &
17 Arndt, S. (2013). Anthropogenic perturbation of the carbon fluxes from land to ocean. *Nature*
18 *Geoscience*, 6(8), 597-607.
- 19 Ricciuto, D. M., A. W. King, D. Dragoni, and W. M. Post (2011), Parameter and prediction
20 uncertainty in an optimized terrestrial carbon cycle model: Effects of constraining variables
21 and data record length, *J. Geophys. Res.*, 116, G01033, doi:10.1029/2010JG001400.
- 22 Ricciuto, D. M., Butler, M. P., Davis, K. J., Cook, B. D., Bakwin, P. S., Andrews, A. and
23 Teclaw, R. M.: Causes of interannual variability in ecosystem-atmosphere CO₂ exchange in a
24 northern Wisconsin forest using a Bayesian model calibration, *Agric. For. Meteorol.*, 148(2),
25 309–327, doi:10.1016/j.agrformet.2007.08.007, 2008.
- 26 Richardson, A. D., Williams, M., Hollinger, D. Y., Moore, D. J. P., Dail, D. B., Davidson, E.
27 a., Scott, N. a., Evans, R. S., Hughes, H., Lee, J. T., Rodrigues, C. and Savage, K.: Estimating
28 parameters of a forest ecosystem C model with measurements of stocks and fluxes as joint
29 constraints, *Oecologia*, 164(1), 25–40, doi:10.1007/s00442-010-1628-y, 2010.
- 30 Rivier L., Ciais P., Hauglustaine D.A., et al. (2006), Evaluation of SF₆, C₂Cl₄, and CO to
31 approximate fossil fuel CO₂ in the Northern Hemisphere using a chemistry transport
32 model, *Journal Geophysical Research-Atmosphere*, 111(D16) - D16311.
- 33 Rödenbeck, C., T. J. Conway, and R. L. Langenfelds (2006), The effect of systematic
34 measurement errors on atmospheric CO₂ inversions: A quantitative assessment, *Atmos.*
35 *Chem. Phys.*, 6, 149–161, doi:10.5194/acp-6-149-2006.
- 36 Rödenbeck C., D.C.E. Bakker, N. Gruber, Y. Iida, A. Jacobson, S. Jones, P. Landschutzer, N.
37 Metzl, S. Nakaoka, A. Olsen, G.-H. Park, P. Peylin, K.B. Rodgers, T.P. Sasse, U. Schuster,
38 J.D. Shutler, V. Valsala, R. Wanninkhof, and J. Zeng, 2015. Data-based estimates of the

- 1 ocean carbon sink variability – First results of the Surface Ocean pCO₂ Mapping
2 intercomparison (SOCOM). *Biogeosciences*, 12: 7251-7278. doi:10.5194/bg-12-7251-2015.
- 3 Rosenblatt, F. (1958). The perceptron: a probabilistic model for information storage and
4 organization in the brain. *Psychological review*, 65(6), 386.
- 5 Ruimy A, Dedieu G, Saugier B, (1996), TURC: A diagnostic model of continental gross
6 primary productivity and net primary productivity. *Global Biogeochemical Cycles*, 10: 269-
7 85.
- 8 Sabine, C.L., et al. 2004. The oceanic sink for anthropogenic CO₂. *Science*, 305 (5682), 367-
9 371.
- 10 Saito, M., A. Ito and S. Maksyutov, Optimization of a prognostic biosphere model for
11 terrestrial biomass and atmospheric CO₂ variability, *Geosci. Model Dev.*, 7, 1829-1840,
12 doi:10.5194/gmd-7-1829-2014, 2014.
- 13 Santaren, D., Peylin, P., Bacour, C., Ciais, P. and Longdoz, B.: Ecosystem model
14 optimization using in situ flux observations: benefit of Monte Carlo versus variational
15 schemes and analyses of the year-to-year model performances, *Biogeosciences*, 11(24), 7137–
16 7158, 2014.
- 17 Schurmann, G. J., Kaminski, T., Kostler, C., Carvalhais, N., Vofßbeck, M., Kattge, J., Giering,
18 R., Rodenbeck, C., Heimann, M., and Zaehle, S.: Constraining a land surface model with
19 multiple observations by application of the MPI-Carbon Cycle Data Assimilation System,
20 *Geosci. Model Dev. Discuss.*, doi:10.5194/gmd-2015-263, in review, 2016.
- 21 Sitch, S., Friedlingstein, P., Gruber, N., Jones, S. D., Murray-Tortarolo, G., Ahlström, A.,
22 Doney, S. C., Graven, H., Heinze, C., Huntingford, C., Levis, S., Levy, P. E., Lomas, M.,
23 Poulter, B., Viovy, N., Zaehle, S., Zeng, N., Arneth, A., Bonan, G., Bopp, L., Canadell, J. G.,
24 Chevallier, F., Ciais, P., Ellis, R., Gloor, M., Peylin, P., Piao, S., Le Quéré, C., Smith, B.,
25 Zhu, Z. and Myneni, R.: Recent trends and drivers of regional sources and sinks of carbon
26 dioxide, *Biogeosciences*, 12, 653–679, doi:10.5194/bgd-12-653-2015, 2015.
- 27 Takahashi, et al. 2009, Corrigendum to "Climatological mean and decadal change in surface
28 ocean pCO₂, and net sea-air CO₂ flux over the global oceans" *Deep Sea Res. II*, 56, 554-577.
- 29 Tarantola A. (1987), *Inverse problem theory: Methods for data fitting and parameter*
30 *estimation*. Elsevier, Amsterdam.
- 31 Tarantola, A. (2005), *Inverse problem theory and methods for model parameters*
32 *estimation*, *Society for Industrial and Applied Mathematics*, Philadelphia, ISBN 0-89871-572-
33 5.
- 34 Thum, T., MacBean, N., Peylin, P., Bacour, C., Santaren, D., Longdoz, B., Loustau, D. and
35 Ciais, P., The potential of forest biomass data in addition to carbon and water flux
36 measurements to constrain ecosystem model parameters: Case studies at two temperate forest
37 sites, *Agriculture and Forest Meteorology*, in revision, 2016.

- 1 Tiedtke M. (1989), A comprehensive mass flux scheme for cumulus parameterization in
2 large-scale models, *Monthly Weather Review*, **117**(8): 1779-1800
- 3 Tucker, C. J. (1979). Red and photographic infrared linear combinations for monitoring
4 vegetation. *Remote Sensing of Environment*, **8**, 127-150.
- 5 Twine, T. E., W. P. Kustas, J. M. Norman, D. R. Cook, P. R. Houser, T. P. Meyers, J. H.
6 Prueger, P. J. Starks, and M. L. Wesely (2000), Correcting eddy-covariance flux
7 underestimates over a grassland, *Agric. For. Meteorol.*, **103**(3), 279–300, doi:10.1016/S0168-
8 1923(00)00123-4.
- 9 Verant, S., Laval, K., Polcher, J., and De Castro, M. (2004), Sensitivity of the continental
10 hydrological cycle to the spatial resolution over the Iberian Peninsula, *Journal of*
11 *Hydrometeorology*, **5**, 267-285.
- 12 Vermote, E., C.O. Justice and F-M Breon (2009), Towards a generalized approach for
13 correction of the BRDF effect in MODIS directional reflectances, *IEEE Transactions on*
14 *Geoscience and Remote Sensing*, **47**, 3, 898-908.
- 15 Wang, Y. P., Baldocchi, D., Leuning, R., Falge, E. and Vesala, T.: Estimating parameters in a
16 land-surface model by applying nonlinear inversion to eddy covariance flux measurements
17 from eight FLUXNET sites, *Glob. Change Biol.*, **13**(3), 652–670, doi:10.1111/j.1365-
18 2486.2006.01225.x, 2007.
- 19 Wang, Y. P., Leuning, R., Cleugh, H. and Coppin, P.: Parameter estimation in surface
20 exchange models using nonlinear inversion : how many parameters can we estimate and
21 which measurements are most useful ?, *Glob. Change Biol.*, **7**, 495–510, doi:10.1046/j.1365-
22 2486.2001.00434.x, 2001.
- 23 Wanninkhof, R., 1992. Relationship between wind speed and gas exchange. *J. Geophys.*
24 *Res.* **97**, 7373-7382.
- 25 Weiss, R.F., 1974. Carbon dioxide in water and seawater: the solubility of a non-ideal gas.
26 *Mar. Chem.*, **2**, 203-215.
- 27 Welp, L. R., Keeling, R. F., Meijer, H. A., Bollenbacher, A. F., Piper, S. C., Yoshimura, K.,
28 ... & Wahlen, M. (2011). Interannual variability in the oxygen isotopes of atmospheric CO₂
29 driven by El Nino. *Nature*, **477**(7366), 579-582.
- 30 Williams, M., Schwarz, P. a, Law, B. E., Irvine, J. and Kurpius, M. R.: An improved analysis
31 of forest carbon dynamics using data assimilation, *Glob. Change Biol.*, **11**(1), 89–105,
32 doi:10.1111/j.1365-2486.2004.00891.x, 2005.
- 33 Xiao, J., Davis, K. J., Urban, N. M. and Keller, K.: Uncertainty in model parameters and
34 regional carbon fluxes: A model-data fusion approach, *Agric. For. Meteorol.*, 189-190, 175–
35 186, doi:10.1016/j.agrformet.2014.01.022, 2014.
- 36 Yapo, P.O., Gupta, H.V. and Sorooshian, S., Multi-objective global optimization for
37 hydrologic models, *J. Hydrol.*, **204**, 83-97, 1998.

1 Zobitz, J. M., D. J. P. Moore, T. Quaife, B. H. Braswell, A. Bergeson, J. A. Anthony, and R.
2 K. Monson (2014), Joint data assimilation of satellite reflectance and net ecosystem exchange
3 data constrains ecosystem carbon fluxes at a high-elevation subalpine forest, *Agric. For.*
4 *Meteorol.*, 195–196, 73–88.

5 Zobler, L (1986), A world soil file for global climate modeling, NASA Technical
6 Memorandum 87802. NASA Goddard Institute for Space Studies, New York, U.S.A.

7

8

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 ($\mu\text{mol}\cdot\text{m}^{-2}\cdot\text{s}^{-1}$)	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 ($\text{m}^2\cdot\text{g}^{-1}$)	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
$K_{soilC,reg}$		30 regions	CO ₂
K_{albedo}	Multiplicative factor of the vegetation albedo	Global	Flux, CO ₂

4

5

6

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.

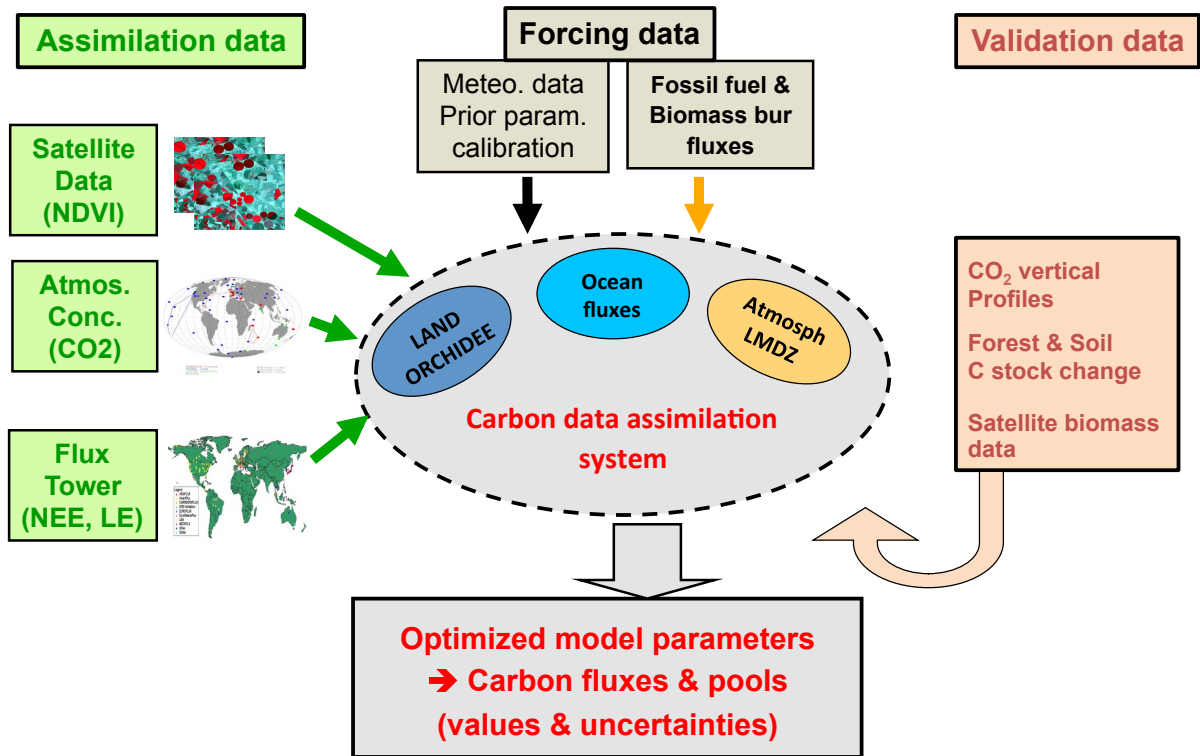
Parameter	Plant functional type								
	TrBE	TrBR	TeNE	TeBE	TeBD	BoNE	BoBD	BoND	NC3
V_{cmax}	65 ± 24 [35; 95]	65 ± 24 [35; 95]	35 ± 12.8 [19; 51]	45 ± 16 [25; 65]	55 ± 20 [30; 80]	35 ± 12.8 [19; 51]	45 ± 16 [25; 65]	35 ± 12.8 [19; 51]	70 ± 25.6 [38; 102]
$G_{s,slope}$	6.0 ± 2.4 [6; 12]	6.0 ± 2.4 [6; 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 **Main Figures**

2

3

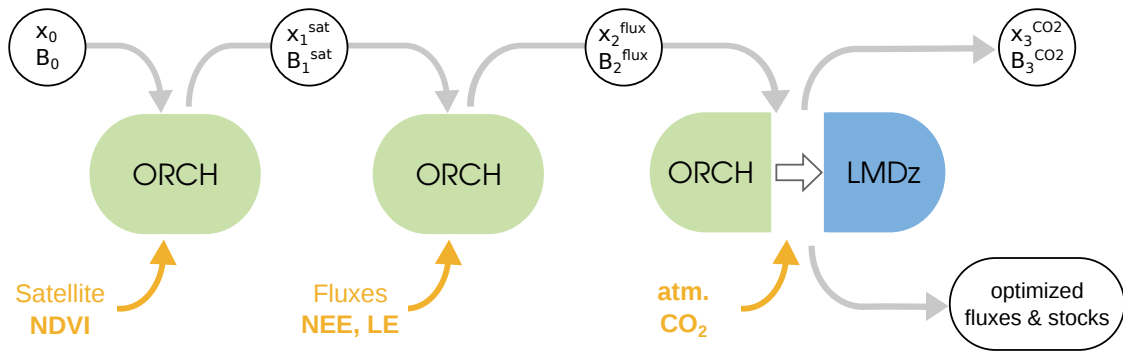
4



5

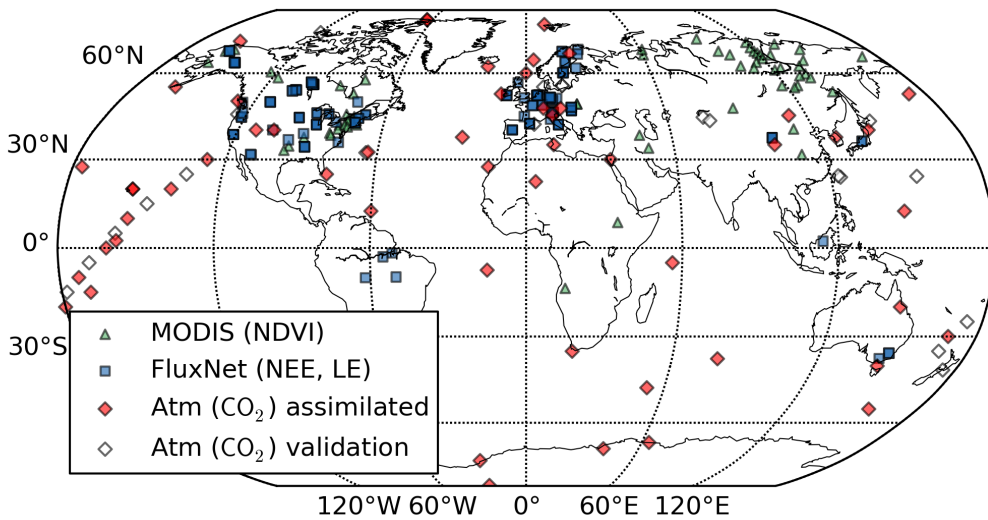
6 Figure 1. Schematic of the ORCHIDEE Carbon Cycle Data Assimilation System
7 (ORCHIDAS).

8

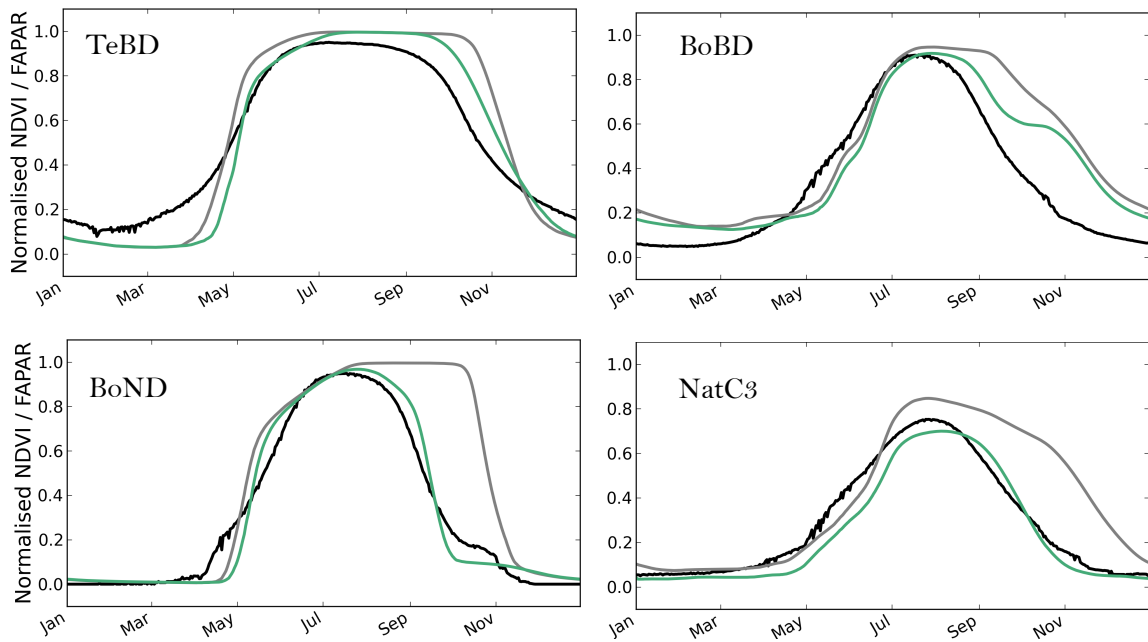


2 Figure 2. Illustration of the step-wise data assimilation approach used for the assimilation of
 3 multiple data streams in the ORCHIDEE-CCDAS. The list of parameters for each step is
 4 summarized in Table 1.

5
 6
 7

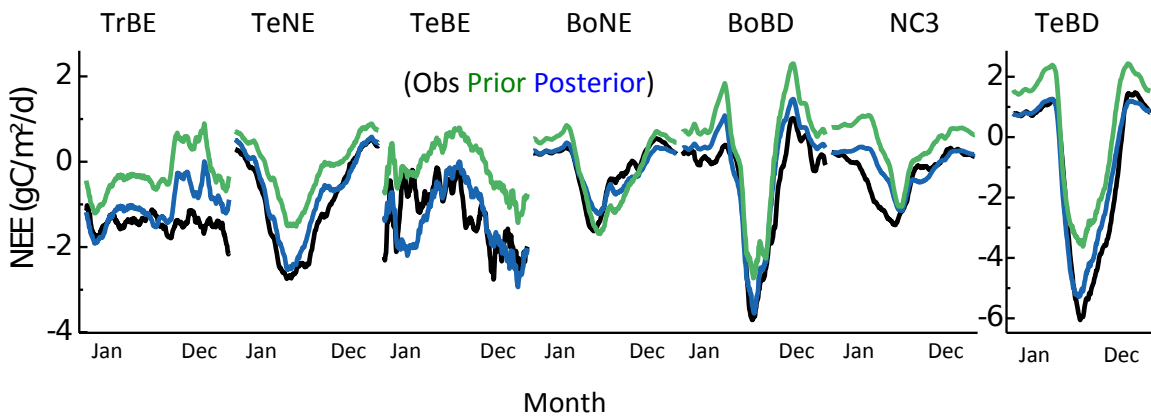


9 Figure 3: Location of the different observations used for each data stream assimilated in the
 10 system: MODIS-NDVI measurements, FLUXNET sites with NEE and LE measurements and
 11 atmospheric CO₂ stations (both the sites that aer assimilated and the sites used for the
 12 validation).

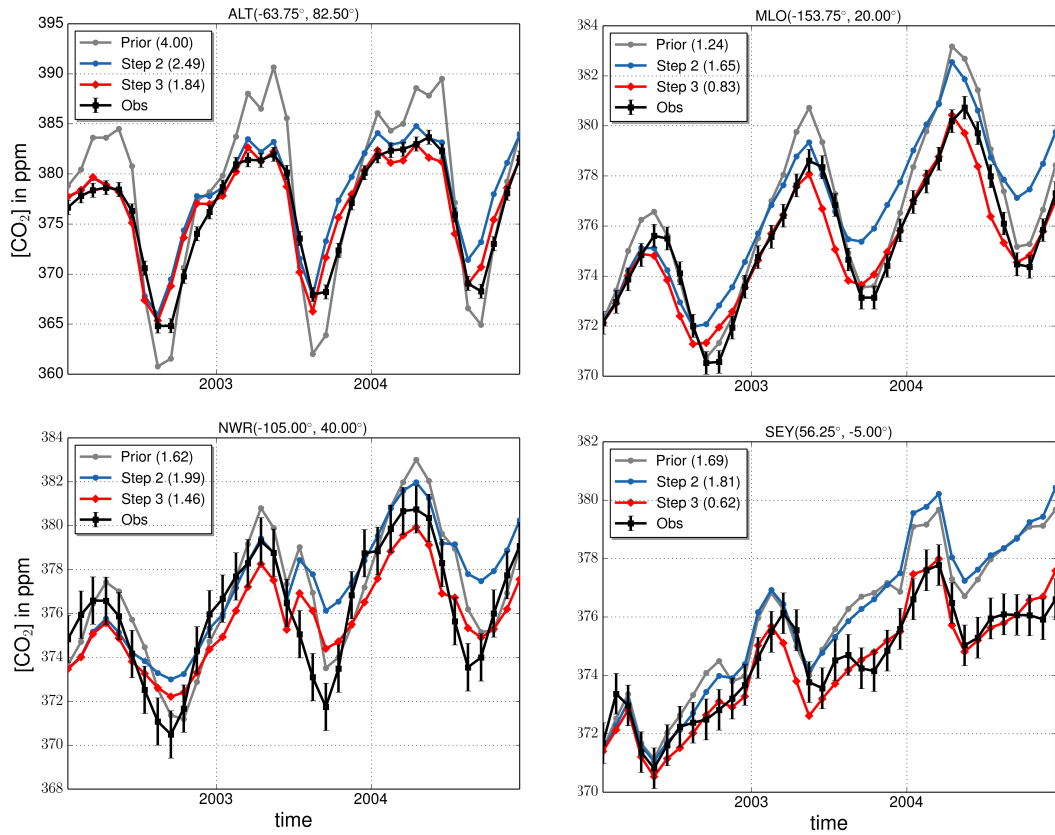


1
 2 Figure 4. Mean seasonal cycle (2000-2008) of the normalised modelled FAPAR before and
 3 after optimisation, compared to that of the MODIS NDVI data, for the temperate and boreal
 4 deciduous PFTs (TeBD, BoBD, BoND and NatC3). Black = MODIS NDVI data; Grey =
 5 prior simulation (default ORCHIDEE 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 optimized in Step 2 of the assimilation. The mean across all
 11 sites for a given PFT is provided for the observations (black), the posterior of step 1 (green)
 12 and the posterior of step 2 (blue).



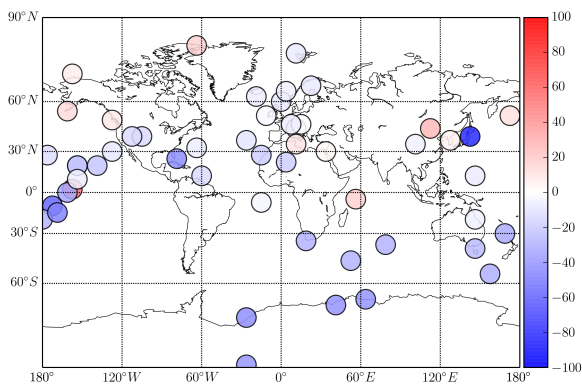
1

2 Figure 6: Monthly mean atmospheric CO₂ concentrations after step 3 of the optimization, for
 3 several stations over the period 2002-2004 of the optimization. The observations (black), the
 4 prior model (grey) and the posterior model after step 2 (blue) and step 3 (red) are displayed.
 5 Numbers in parenthesis correspond to RMSEs.

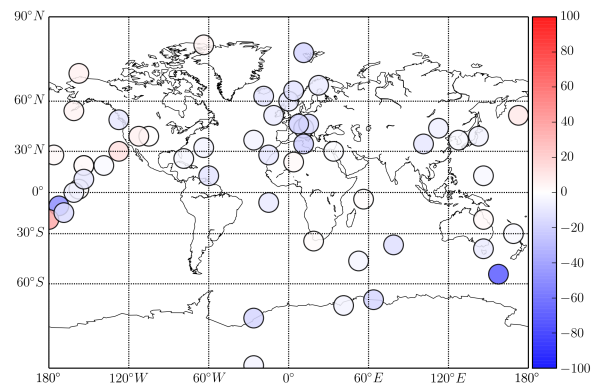
6

7

a) Seasonal amplitude: relative changes



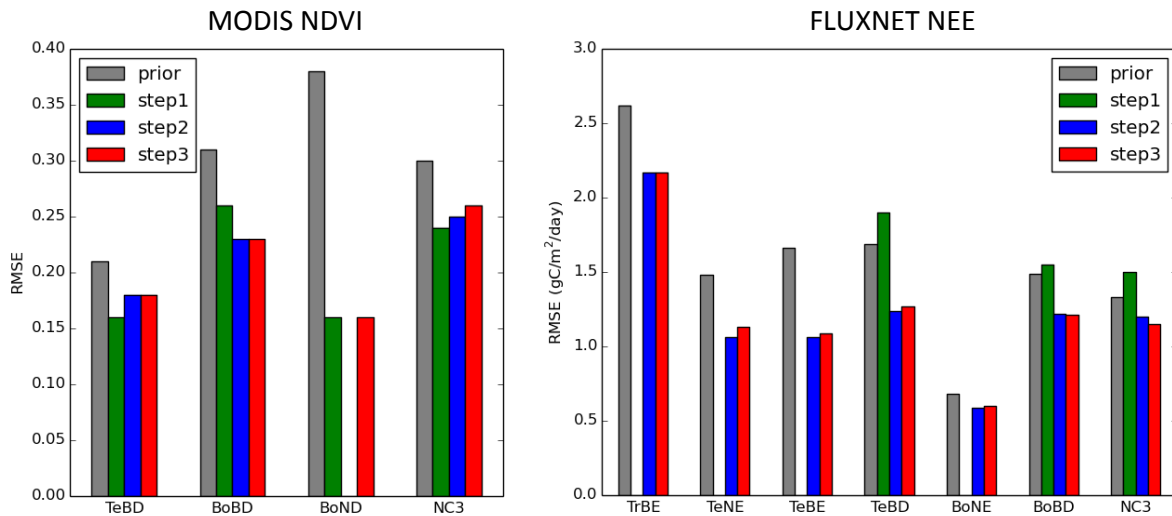
b) Length of CUP: relative changes



8

1 Figure 7: Changes in the mean seasonal cycle of the atmospheric CO₂ concentrations after
 2 step 3 of the optimization at all atmospheric stations. Left: Relative changes (in percentage)
 3 between the prior of step 3 and posterior absolute model-data differences for the amplitude of
 4 the 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).

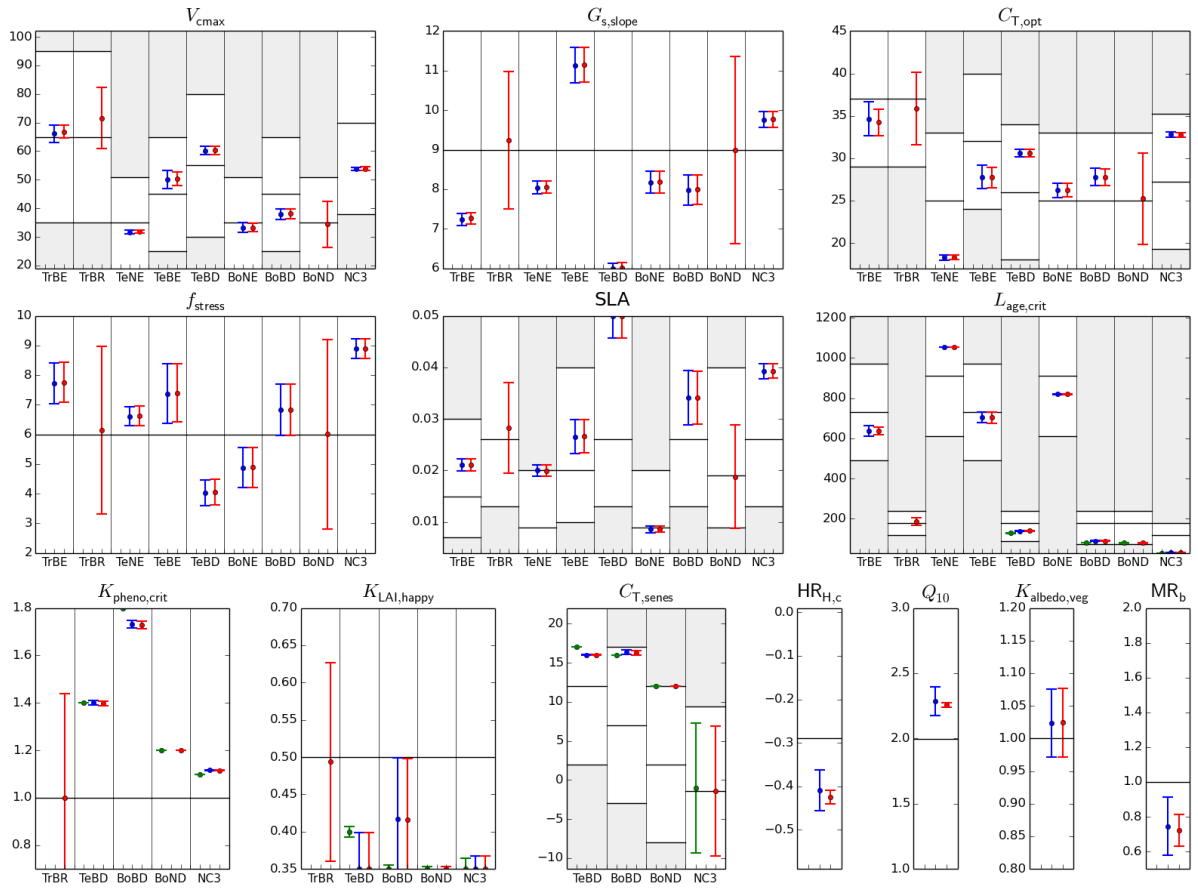
6
 7
 8



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 CO₂
 9 data (red symbol). The parameter values are depicted with the symbols and the estimated
 10 uncertainties with the vertical line ($\pm \sigma$).

11

12

13

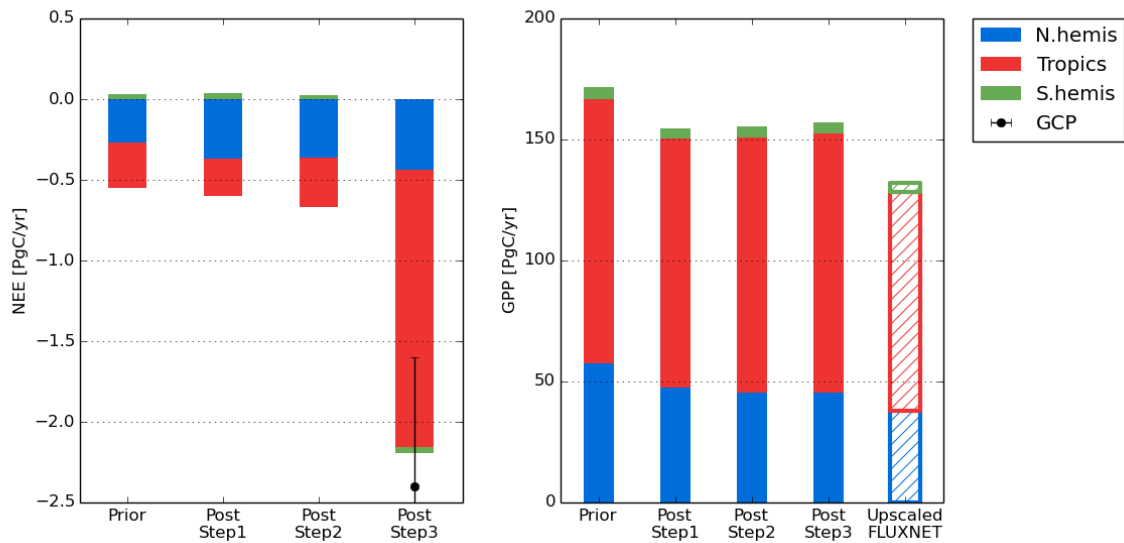
14

15

16

17

18



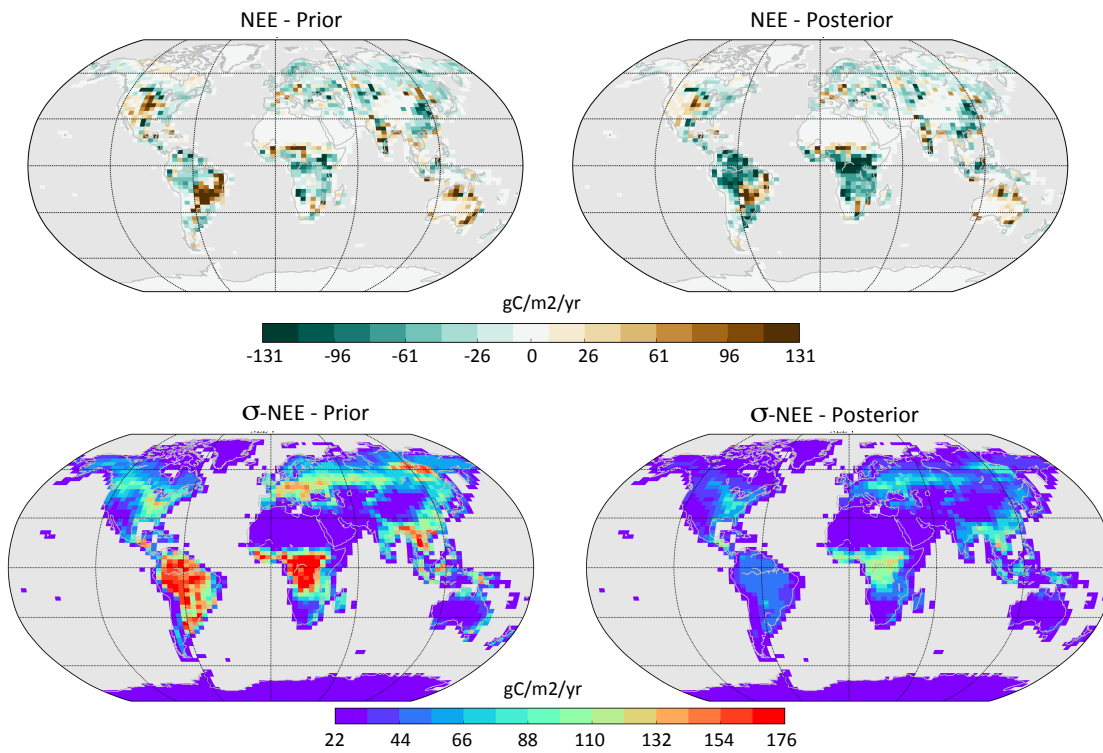
1

2 Figure 10: Left: Net Ecosystem Exchange (NEE) for three regions (North of 35°N, Tropics,
 3 South of 35°S) for the prior model, and after each step of the optimizations (mean over 2002-
 4 2004). The NEE estimate from the Global Carbon Project (GCP) for the same period (Le
 5 Quéré et al. 2015) is provided for step 3 with its error bar. Right: same but for Gross Primary
 6 Production where the data driven estimate (MTE product using FluxNet data; Jung et al.,
 7 2009) is provided for comparison.

8

9

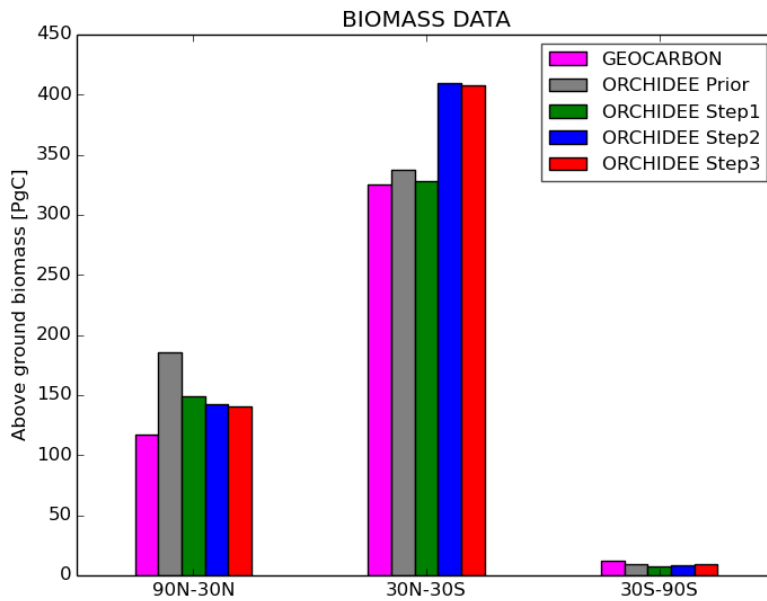
1



2

3 Figure 11: Simulated annual net carbon exchange (NEE) for the land ecosystems prior to any
4 optimization (left column) and after step 3 of the optimization process (right column). Upper
5 figures correspond to the mean NEE (in $\text{gC}\cdot\text{m}^{-2}\cdot\text{y}^{-1}$) over the assimilation period (2001-2003)
6 and lower figures to the associated monthly flux uncertainties (averaged over the whole
7 period and expressed in $\text{gC}\cdot\text{m}^{-2}\cdot\text{y}^{-1}$) due to the parameter uncertainties (see text).

8



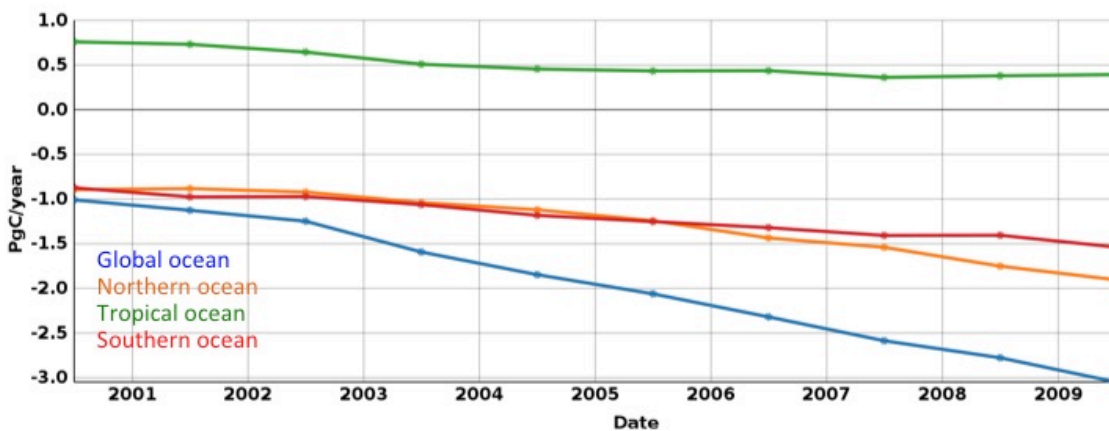
1

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 Appendix figures

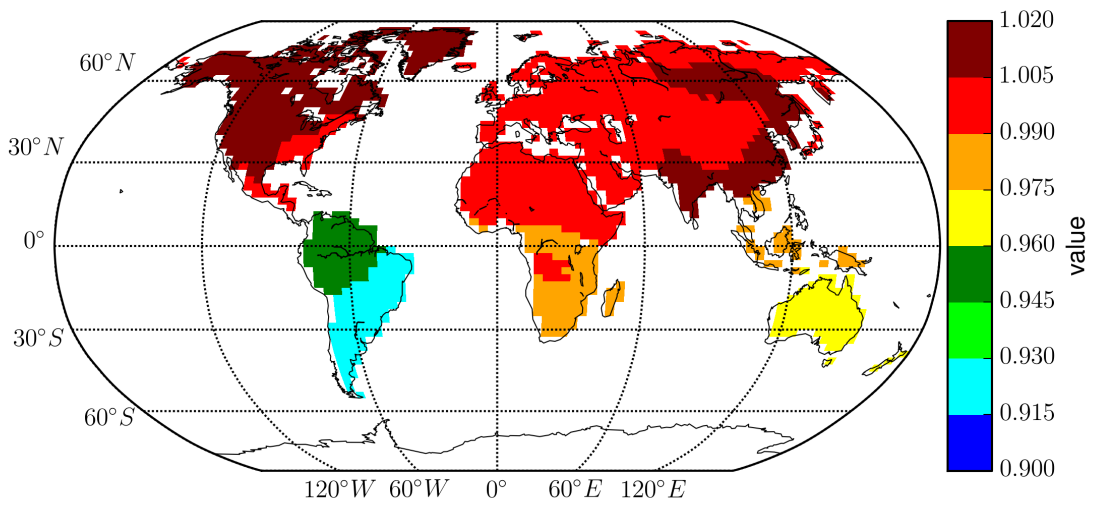
8



9

10 Figure A1: CO₂ air-sea fluxes including the natural ocean out-gazing, used as input to the
 11 ORCHIDEE-CCDAS and estimated from a neural network approach using observed pCO₂
 12 data (see main text, section 2.5.1). The Northern, Tropical and Southern ocean contributions
 13 to the global ocean flux (blue curve) are also provided.

1



2

3 Figure A2: Map of the posterior values of the coefficient scaling the initial carbon pool sizes
4 per regions.

5

6