



1 **Consistent assimilation of multiple data streams in a** 2 **carbon cycle data assimilation system**

3

4 **Natasha MacBean¹, Philippe Peylin¹, Frédéric Chevallier¹, Marko Scholze²,**
5 **Gregor Schürmann³**

6 [1]{Laboratoire des Sciences du Climat et de l'Environnement, LSCE/IPSL, CEA-CNRS-
7 UVSQ, Université Paris-Saclay, F-91191 Gif-sur-Yvette, France}

8 [2]{Department of Physical Geography and Ecosystem Science, Lund University, Lund,
9 Sweden}

10 [3]{Max Planck Institute for Biogeochemistry, Jena, Germany}

11 Correspondence to: N. MacBean (nlmacbean@gmail.com)

12

13 **Abstract**

14 Data assimilation methods provide a rigorous statistical framework for constraining
15 the parametric uncertainty of land surface models (LSMs), with the aim of improving our
16 predictive capability as well as identifying areas in which the models need improvement. The
17 increase in the number of available datasets in recent years allows us to address different
18 aspects of the model at a variety of spatial and temporal scales. However, combining data
19 streams in a DA system is not a trivial task. In this study we highlight some of the challenges
20 surrounding multiple data stream assimilation, with a particular focus on the carbon cycle
21 component of LSMs. We examine the impact of biases and inconsistencies between the
22 observations and the model (resulting in non Gaussian error distributions) and the impact of
23 non-linearity in model dynamics. In addition we explore the differences between performing a
24 simultaneous assimilation (in which all data streams are included in one optimisation) and a
25 step-wise approach (in which each data stream is assimilated sequentially), given the
26 assumptions inherent to the inversion algorithm chosen for this study. We demonstrate some
27 of these issues by assimilating synthetic observations into two simple models: the first a
28 simplified version of the carbon cycle processes represented in many LSMs, and the second a
29 non-linear toy model. We further discuss these experimental results in the context of recent



1 studies in the carbon cycle data assimilation literature, and finally we provide some
2 perspectives and advice to other land surface modellers wishing to use multiple data streams
3 to constrain their models.

4 Keywords: data assimilation, carbon cycle, biogeochemical cycles, land surface model

5

6 **1 Introduction**

7 The carbon cycle is an important component of the Earth system, especially when
8 considering the climatic impact of rising greenhouse gases concentrations from fossil fuel
9 emissions and land use change. It is estimated that the oceans and land surface absorb
10 approximately half of the CO₂ emissions due to anthropogenic activity, but uncertainties
11 remain in the strength and location of sources and sinks, as well as in predictions of future
12 trends (Ciais et al., 2013). Observations allow us to understand the system up until the present
13 day, but they cannot tell us about the future, and can be limited in their spatial coverage. They
14 also cannot distinguish between the complex interactions that may occur between different
15 processes. Incorporating our current knowledge of physical mechanisms of biogeochemical
16 cycles, including carbon, C, dynamics, into Land Surface Models (LSMs) represents a
17 promising approach to analyse these interacting effects, to upscale observations to larger
18 regions, and to make future predictions. However, the models can be limited by the lack of
19 process representation, either due to gaps in our knowledge or in our technical and computing
20 capability. As a result, model evaluations reveal that not all variables are well-captured by the
21 model under current conditions (Anav et al., 2013), and the spread between model projections
22 is still very large (Sitch et al., 2015).

23 Aside from model structural and forcing errors, one source of uncertainty is related to the
24 parameter (i.e. fixed) values of a model. Model-data fusion, or data assimilation (DA), allows
25 the calibration, or optimisation, of these values by reducing the model-data misfit while
26 accounting for the uncertainties inherent in both the model and data in a statistically rigorous
27 framework. The C cycle component of most LSMs is complex and contains a large number of
28 parameters. Luckily however, there are an increasing number of in-situ and remote sensing-
29 based data streams that can be used for parameter optimisation. These data bring information
30 on different spatial and temporal scales, such as:



- 1 • Atmospheric CO₂ concentration data measured at surface stations at continental to
2 global scales, which provide information from synoptic timescales to inter-annual
3 variability (IAV) and long-term trends.
- 4 • Eddy covariance net CO₂ (net ecosystem exchange – NEE) and latent (LE) and
5 sensible heat fluxes measured at half-hourly intervals at many sites across different
6 ecosystems/regions, providing information at seasonal to inter-annual timescales.
- 7 • Satellite-derived measures of vegetation dynamics, including “greenness” indices (i.e.
8 the Normalised Difference Vegetation Index – NDVI), fraction of absorbed
9 photosynthetically active radiation (FAPAR) and leaf area index (LAI) at global scales
10 and at daily time step spanning more than a decade, thus capturing IAV and long-term
11 trends (though usually with a trade-off between spatial and temporal resolution).
- 12 • Satellite-derived measurements of soil moisture and land surface temperature at the
13 same temporal and spatial scales as the satellite-derived observations of vegetation
14 productivity.
- 15 • Aboveground biomass measurements are currently taken at only one or a few points in
16 time at plot scale up to regional scale from aircraft and satellite data, or are estimated
17 from allometric relationships at each site.
- 18 • Soil C stock estimates usually are only taken at one point in time at plot scale.
- 19 • Ancillary data on vegetation characteristics such as tree height or budburst – one
20 measured at certain well-instrumented sites.

21

22 Increasingly, researchers are attempting to bring these sources of information together to
23 constrain different parts of a model at different spatio-temporal scales within a multiple data
24 stream assimilation framework (e.g. Richardson et al., 2010; Keenan et al., 2012; Kaminski et
25 al., 2012; Forkel et al., 2014; Bacour et al., 2015). However, whilst the potential benefit of
26 adding in extra data streams to constrain the C cycle of LSMs is clear, multiple data stream
27 assimilation is not as simple as it may seem. When using more than one data stream there is
28 the option to include all data streams together in the same optimisation (simultaneous
29 approach), or to take a sequential (step-wise) approach. Mathematically, the optimal approach
30 is the simultaneous, but computational constraints related to the inversion of large matrixes or



1 the requirement of numerous simulations (especially for global datasets), and/or the weight of
2 different data streams in the optimisation, may complicate a simultaneous optimisation. On
3 the other hand, in a step-wise assimilation the parameter error covariance matrix has to be
4 propagated at each step, which implies that it can be computed. If the parameter error
5 covariance matrix can be properly estimated and is propagated between each step, the step-
6 wise approach can be mathematically equal to simultaneous. However, many inversion
7 algorithms (e.g. derivative based methods that use the gradient of the cost function to find its
8 minimum) require assumptions of model (quasi-) linearity and Gaussian parameter and
9 observation error distributions. If these assumptions are violated, or the error distributions are
10 poorly defined, it is likely that the step-wise will not be equal to the simultaneous, and that
11 information will be lost at each step. An incorrect description of the observation (– model)
12 error distribution could result from the wrong assumption about the distribution of the
13 residuals between the observation and the model, a poor characterisation of the error
14 correlations, an incompatibility between the model and the data (possibly due to a model
15 structural issue or differences in how a variable is characterised), or a bias in the observations
16 that is not unaccounted for (i.e. is treated as a random error). Whilst a simultaneous
17 optimisation is mathematically more rigorous in the sense that the error correlations are
18 treated within the same inversion, if the prior distributions are not properly characterised any
19 bias may be aliased to the wrong parameters (Wutzler and Carvalhais, 2014), more so than in
20 a step-wise approach.

21 This tutorial-style paper demonstrates some of the challenges of multiple data stream
22 assimilation discussed above with two simple models: one a simplified version of the carbon
23 dynamics included in many LSMs, and the other a “toy” model designed to demonstrate the
24 issues that arise with complex, non-linear models. Section 2 provides a description of these
25 models, the inversion algorithm used to optimise the model parameters and the experiments
26 performed, followed by the results for each test case. Section 3 further discusses the
27 challenges outlined in Section 2 with reference to recent carbon cycle multiple data stream
28 assimilation studies in the literature. Finally Section 4 provides some advice to land surface
29 modellers wishing to carry out multiple data stream assimilation.

30



1 **2 Demonstration with two simple models and synthetic data**

2 **2.1 Methods**

3 **2.1.1 Simple carbon model**

4 To demonstrate the challenges of multiple data stream assimilation in a carbon cycle
 5 context, we have chosen a test model that represents a simplified version of the carbon cycle
 6 dynamics typically implemented in most LSMs. The model has been well-documented in
 7 Raupach (2007) and has been used previously in the OptIC DA inter-comparison project
 8 (Trudinger et al., 2007). It is based on two equations that describe the temporal evolution of
 9 two carbon pools, s_1 and s_2 :

$$10 \quad \frac{ds_1}{dt} = F(t) \left(\frac{s_1}{p_1 + s_1} \right) \left(\frac{s_2}{p_2 + s_2} \right) - k_1 s_1 + s_0 \quad (1)$$

$$11 \quad \frac{ds_2}{dt} = k_1 s_1 - k_2 s_2 \quad (2)$$

12 In this model formulation, s_1 and s_2 are approximately equivalent to above- and belowground
 13 biomass stocks. The unknown parameters p_1 , p_2 , k_1 and k_2 will be optimised in the inversions.
 14 The first term on the right-hand side of Eq. (1) corresponds to the Net Primary Production
 15 (NPP) i.e. the carbon assimilated into the system as a function of time, $F(t)$, weighted by
 16 factors that account for the size of both pools in order to introduce a limitation on NPP (the
 17 two fractions in parentheses). The litterfall is an output of s_1 and an input to s_2 and is a
 18 constant fraction of the aboveground carbon reserve as represented by $k_1 s_1$. Heterotrophic
 19 respiration (Rh) is an output of s_2 and is represented $k_2 s_2$. The constant s_0 is a “seed
 20 production” term set to 0.01 (i.e. not optimised) to ensure the model does not verge towards
 21 zero. A more detailed description of the properties of the model is given in Trudinger et al.
 22 (2007) and an in-depth analysis of the model behaviour is provided in Raupach (2007).
 23 Synthetic observations of both s_1 and s_2 variables were used to optimise all the unknown
 24 parameters in the model (see Section 2.1.5).

25

26 **2.1.2 Non-linear toy model**

27 Although the simple carbon model contains a non-linear term it is essentially still a
 28 quasi-linear model. In order to illustrate the challenges associated with multiple data stream



1 data assimilation for more complex non-linear models, we defined a simple non-linear toy
2 model based on two equations with two unknown parameters:

$$3 \quad s_1 = a \exp^b + at^2 \quad (3)$$

$$4 \quad s_2 = \sin(10a + 10b) + 10t^2 \quad (4)$$

5 where s_1 and s_2 also correspond to two model state variables (as for the simple C model), a
6 and b are the unknown parameters included in the optimisation, and t is the independent
7 variable, which could represent time in a real-world scenario. Note that this model is not
8 based on any particular physical process associated with land surface biogeochemical cycles,
9 but it does contain typical mathematical functions that are observed in reality and
10 implemented in LSMs. For example, the sinusoidal function (Eq. (4)) could represent diurnal
11 variations of various processes such as photosynthesis and respiration. Exponential response
12 functions (such as in Eq. (3)) are also observed for certain processes, including the
13 temperature sensitivity of soil microbial decomposition. As for the simple carbon model,
14 synthetic observations corresponding to the s_1 and s_2 variables were used to optimise both
15 parameters (see Section 2.1.5).

16

17 2.1.3 Bayesian inversion algorithm

18 Most data assimilation approaches follow a Bayesian formalism which, simply put,
19 allows prior knowledge of a system (in this case the model parameters) to be updated, or
20 optimised, based on new information (from the observations). In order to achieve this we
21 define a “cost function” that describes the misfit between the data and the model, taking into
22 account their respective uncertainties, as well as the uncertainty on the prior information. If
23 we follow a Bayesian formalism and least-squares minimisation approach, and assume
24 Gaussian probability distributions for the model parameter and observation error
25 variance/covariance, we derive the following cost-function (Tarantola, 1987):

$$26 \quad J(\mathbf{x}) = \frac{1}{2} [(H(\mathbf{x}) - \mathbf{y})^T \cdot \mathbf{R}^{-1} \cdot (H(\mathbf{x}) - \mathbf{y}) + (\mathbf{x} - \mathbf{x}^b)^T \cdot \mathbf{B}^{-1} \cdot (\mathbf{x} - \mathbf{x}^b)] \quad (5)$$

27 where \mathbf{y} is the observation vector, $H(\mathbf{x})$ the model outputs given parameter vector \mathbf{x} , \mathbf{R} the
28 observation error covariance matrix (including measurement and model errors), \mathbf{x}^b the a priori



1 parameter values, and \mathbf{B} the prior parameter error covariance matrix. This framework leads to
2 a Gaussian posterior parameter probability distribution function.

3 The aim of the inversion algorithm is to find the minimum of this cost function,
4 thereby achieving the best possible fit between the model simulations and the measurements,
5 conditioned on their respective uncertainties and prior information. For cases where there is a
6 strong linear dependence of the model to the parameters (at least for variations in \mathbf{x} of the size
7 of those expected in the data assimilation system), and where the dimensions of the problem
8 are not too large, the solution can be derived analytically. If not, as is usually the case with
9 LSMs, there are different numerical methods to find the most optimal parameter values.
10 These include global search methods that randomly search the parameter space and test the
11 likelihood of a particular parameter set at each iteration, and derivative methods, which
12 calculate the gradient of the cost function at each iteration to find its minimum. In this study
13 we use the latter class of methods. More specifically we use a quasi-Newton algorithm that
14 uses both the gradient of the cost function and its derivative (Hessian) to evaluate if the
15 minimum has been reached (i.e. where the gradient is zero). Thus we obtain the following
16 algorithm for iteratively finding the minimum (Tarantola, 1987, p195):

$$17 \quad \mathbf{x}_{i+1} = \mathbf{x}_i + \varepsilon_i [\mathbf{H}^T \mathbf{R}^{-1} \mathbf{H} + \mathbf{B}^{-1}]^{-1} \mathbf{H}^T \mathbf{R}^{-1} (\mathbf{y} - H(\mathbf{x}) - \mathbf{B}^{-1}(\mathbf{x}_i - \mathbf{x}^b)) \quad (6)$$

18 where i is the iteration number and \mathbf{H} is the Jacobian, or first-order derivatives, of H , which in
19 this study is determined using a finite difference method. Note that as we are potentially
20 dealing with non-linear models, the quasi-Newton method has been slightly adapted to
21 include the constant scaling factor ε_i (with a value <1.0) to ensure that the algorithm will
22 converge.

23 Of course no inversion algorithm is perfect, and therefore it is possible that the true
24 “global” minimum of the cost function has not been found. Derivative methods in particular
25 can get stuck in so-called “local minima”, preventing the algorithm from finding the true
26 minimum. To address this issue we carry out a number of assimilations with different random
27 first guess points in the parameter space. If they all result in the same reduction in cost
28 function value, we can have more confidence that the true minimum has been found.

29 Once the minimum of the cost function has been found, the posterior parameter error
30 covariance can be approximated (using the linearity assumption) from the inverse Hessian of



1 the cost function around its minimum, which is calculated using the Jacobian of the model at
2 the minimum of $J(\mathbf{x})$ (for the set of optimized parameters), \mathbf{H}_∞ , following Tarantola (1987):

$$3 \quad \mathbf{A} = [\mathbf{H}_\infty^T \mathbf{R}^{-1} \mathbf{H}_\infty + \mathbf{B}^{-1}]^{-1} \quad (7)$$

4 Note that the posterior error covariance matrix can be propagated into the model space to
5 determine the posterior uncertainty on the simulated state variables as a result of the
6 parametric uncertainty (as shown in the coloured error bands in the time series plots – Figures
7 1 and 5) using the following matrix product and the hypothesis of local linearity (Tarantola,
8 1987):

$$9 \quad \mathbf{R}_{post} = \mathbf{H}_\infty \mathbf{A} \mathbf{H}_\infty^T \quad (8)$$

10 However, we do not detail the propagated posterior uncertainty on the state variables further
11 in this study; rather, we describe the impact of the optimisation on the model–data fit in terms
12 of the RMSE value and also in terms of the relative uncertainty reduction on the parameters.

13

14 2.1.4 Step-wise versus simultaneous assimilation

15 *Step-wise approach*

16 In the step-wise approach each data stream (in our cases s_1 and s_2 , see above) is
17 assimilated sequentially, and the posterior error covariance matrix of Eq. (7) is propagated to
18 the next step as the prior in Eq. (6). Note that the error covariance matrix can only be
19 propagated if it is calculated within the inversion algorithm, which is the case here but may
20 not be possible in other studies. The following details an example for two data streams.

21 Step 1: Assimilation of the first data stream, s_1 . The prior parameters, including their values
22 and error covariance (\mathbf{x}^b and \mathbf{B}), are optimised to produce a first set of posterior
23 optimised parameters \mathbf{x}_1 with error covariance \mathbf{A}_1 .

24 Step 2: Assimilation the second data stream, s_2 . The parameters, \mathbf{x}_1 , and their error
25 covariance, \mathbf{A}_1 , are used as a prior to the optimisation system and further optimised
26 to produce the second (and final) set of posterior optimised parameters, \mathbf{x}_{post} , and the
27 associated error covariance \mathbf{A} .



1 *Simultaneous approach*

2 Both data streams s_1 and s_2 are included in the optimisation and all parameters are optimised
3 at the same time. The prior parameters, including their values and error covariance (\mathbf{x}^b and \mathbf{B})
4 are optimised to produce the posterior parameter vector (\mathbf{x}_{post}) and associated uncertainties \mathbf{A} .

5

6 2.1.5 Optimisation set-up: parameter values and uncertainty, and generation 7 of synthetic observations

8 In this study we used synthetic observations that were generated by running the model
9 with known (or ‘true’) parameter values and adding random Gaussian noise corresponding to
10 the defined observation error for both s_1 and s_2 (see Table 1). The true values of all parameters
11 for both models are given in Table 1, together with their upper and lower bounds. The
12 parameter uncertainty (1 sigma) was set to 40% of the parameter range following recent
13 studies (e.g. Bacour et al., 2015). Prior values were chosen from a uniform random
14 distribution bounded by the parameter bounds. We assumed independence (i.e. uncorrelated
15 errors) for both the parameters and observation covariance matrices, thus the \mathbf{R} and \mathbf{B}
16 matrices were diagonal.

17

18 2.1.6 Experiments

19 Table 2 details the experiments that were carried out based on all possible combinations
20 for assimilating the two data streams. Three approaches were compared: i) separate – where
21 only one data stream was included in the optimisation; ii) step-wise – where each data stream
22 was assimilated sequentially; and iii) simultaneous – where both data streams were included in
23 the optimisation. All parameters for both models were optimised in all experiments, therefore
24 in the step-wise cases the parameters were optimised twice. Tests for the step-wise were also
25 carried out with and without the propagation of the full posterior parameter error covariance
26 matrix, \mathbf{A}_1 , in between steps 1 and 2 (test cases 2b and d – see Table 2) – i.e. for these tests
27 only the posterior variance was propagated. An additional test was included for the
28 simultaneous assimilation in order to test the impact of having a substantial difference in the
29 number of observations for the data stream included in the optimisation; therefore in test case
30 3b, only one observation was included for data stream s_2 .



1 The differences in the parameter values and the theoretical reduction in their uncertainty
2 $(1 - (\sigma_{\text{post}} / \sigma_{\text{prior}}))$ were examined for all eight test cases, as well as the fit (RMSE) to both
3 data streams after the optimisation. For the step-wise approach we investigated if the fit to the
4 first data stream is degraded in the second step by comparing the RMSE after each step. Note
5 that the reduction in uncertainty is a theoretical, or approximate, estimate of the real
6 uncertainty reduction because of the assumptions made in the inversion scheme.

7 In a second stage the impact of an unknown, un-accounted for bias in the model was
8 examined. This bias could be a systematic bias in the observations due to the algorithm used
9 for their derivation, the result of missing or incomplete processes in the model, or an
10 incompatibility between the observations and the model, for example due to differences in
11 spatial resolution or an inconsistent characterisation of a variable between the model and the
12 observations. To test the impact of such an occurrence, we introduced a uniform scalar bias
13 into the modelled s_2 variable with a value of 10 (i.e. twice the magnitude of the defined
14 observation uncertainty). All eight experiments were repeated, but a bias was introduced into
15 the model calculation of s_2 that was not accounted for in the cost function (i.e. the error
16 distributions retained a mean of zero). This was treated as an unknown bias, and therefore not
17 corrected or accounted for in the inversion scheme and the defined observation uncertainty
18 (Table 1) was not changed for this set of experiments.

19 In all experiments for both models twenty assimilations were performed starting from
20 different random “first guess” points in the parameter space. As discussed in Section 2.1.3
21 this was done to test the ability of the algorithm to converge to the true global minimum of the
22 cost function. Note that the global minimum and possible reduction in $J(x)$ will be different
23 for each experiment, as each is based on a different cost function.

24

25 **2.2 Results**

26 The twenty random first guess assimilations were examined for each set of experiments
27 for both models (before the results for each test were examined in more detail), in order to
28 check that the algorithm converged to a global minimum. As shown in the supplemental
29 information (Fig. S1), a high proportion of the twenty first guess assimilations across all test
30 cases for both models resulted in a similar reduction in $J(x)$, even though the overall
31 magnitude of the reduction was sometimes different between tests. This indicates that the



1 algorithm does not easily get stuck in any local minima (if they exist). The examples shown in
2 the results below were taken from one first guess parameter set for each model that belonged
3 to the cluster that had the highest cost function reduction. Any differences seen in the
4 parameter values, their posterior uncertainty or the resultant RMSE reduction described below
5 therefore are due to the specific details of each test and not the inability of the algorithm to
6 find the minimum.

7

8 2.2.1 Typical performance with a quasi-linear model and no bias

9 Figures 1a and b show the simple carbon model simulations for test case 3a (in which
10 both data streams are assimilated simultaneously) for the s_1 and s_2 variables. A large reduction
11 in RMSE is achieved after optimisation (blue curve) with respect to the observations (black
12 curve). Overall, there is a good reduction in RMSE for all test cases (including the individual
13 assimilations 1a and 1b) with a reduction of ~80% for s_1 and s_2 . In addition, the optimisation
14 of the s_1 and s_2 variables resulted in a good or moderate reduction in RMSE for variables not
15 included in any assimilation: ~60% for the litterfall (Eqn. 1) and ~16% for the heterotrophic
16 respiration (Rh – Eqn. 2) across all test cases (not shown), although there was already a good
17 prior fit to the data. As would be expected from these results, the parameter values and the
18 theoretical reduction in parameter uncertainty do not vary between the tests (Figures 2a and b
19 blue symbols), except for a slight difference in the value of the k_2 parameter in test cases 1a
20 and 3b, for which there is also a lower reduction in uncertainty (~82% compared to >95%).
21 Note that Fig. 2a shows the normalised parameter values to account for differences in the
22 magnitude of the different parameters and their range (the zero line represents the “true”
23 parameter value – see caption). In this situation therefore, where we have a relatively simple
24 linear model and two data streams to which the model parameters are highly sensitive, we see
25 that the differences between the step-wise and simultaneous approaches are minimal. This is
26 even the case when the error covariance is not propagated between the two steps (test cases 2b
27 and d), suggesting that under this assimilation set-up both s_1 and s_2 individually contain
28 enough information to retrieve the true values of all parameters, as we can see from the
29 separate test cases 1a and b.

30



1 2.2.2 Impact of unknown bias in one data stream – example with a simple 2 carbon model

3 In Section 2.2.1 we saw that there is little difference between a step-wise and
4 simultaneous optimisation if there is no bias in the model or observations, and if the model is
5 quasi-linear and therefore the critical assumptions behind the inversion approach were not
6 violated. However, it is not uncommon to have a bias between your observations and model
7 that is not obvious and therefore not accounted for in the optimisation, as the cost function
8 used in most inversion algorithms (and in this study) assume Gaussian error distributions with
9 zero mean. Note that this is also the case when defining a likelihood function for accepting or
10 rejecting parameter values in a global search method. To test the impact of a bias, we added a
11 uniform value to the simulated s_2 variable in a second test (see Section 2.1.6) that was treated
12 as an unknown bias, and therefore not corrected or accounted for in the inversion scheme. The
13 impact of this bias on s_1 and s_2 is shown in Figures 1c-d, and the reduction in RMSE between
14 the model and observations is seen in Fig. 3 for all variables (including Rh and litterfall). The
15 red symbols in Fig. 2 show the resultant parameter values and theoretical reduction in
16 uncertainty as a result of the bias. The inversion cannot accurately find the correct values for
17 all parameters in any test case and there are now considerable differences between the
18 simultaneous and step-wise approach. Furthermore the order in which the data streams are
19 assimilated in the step-wise cases also results in different posterior parameter values (test
20 cases 2a and b versus 2c and d in Fig. 2a and Fig. 3). Nevertheless the optimisation results in
21 a similar reduction in uncertainty on the parameters, except in test case 1b where only s_2 data
22 are assimilated (Fig. 2b).

23 The main impact of the bias in the modelled s_2 variable is on the value of k_2 parameter
24 (Fig. 2a), which is consistently offset from the true value (dashed line in Fig. 2a) in all test
25 cases. This was expected given that it is the parameter most directly related to the calculation
26 of s_2 . However, in test cases 2a and 3a, the values of p_1 and p_2 are also incorrect (and p_1 for
27 test case 2b). Note that these parameters only indirectly influence the s_2 pool in the model,
28 and therefore we might have expected that they would be less affected by the bias. This nicely
29 demonstrates one issue that could arise in all DA studies, where the bias in a particular
30 variable (in the observations or the model) is aliased onto another process in the model
31 (Wutzler and Carvalhais, 2014). Such an aliasing of bias onto indirectly related parameters is
32 even more evident when only s_2 is included in the assimilation and s_1 does not provide any



1 constraint (test case 1b) – in this case all parameters are incorrect but the p_2 parameter in
2 particular shows a strong deviation from the true value (Fig. 2a). As a result we see a
3 deterioration in the RMSE for the s_1 , litterfall and Rh variables in test case 1b and in the step-
4 wise cases where s_2 is assimilated in the second step (Figures 3a, c and d – test case 1b, 2a
5 and 2b). However, the RMSE reduction remains high for the s_2 variable for these test cases
6 (Fig. 3b), as the inversion has found a solution that accounts for the bias even though all
7 inferred parameter values are incorrect. The assimilation of s_1 in the second step lowers the
8 reduction in RMSE for s_2 gained in the first step to $\sim 70\%$, but it is not a considerable
9 degradation.

10 Even though the posterior parameter values are incorrect, and despite the fact that the
11 first step results in a degradation, the final reductions in RMSE are largely the same than the
12 situation with no bias for all variables when s_1 is included in a simultaneous assimilation or
13 optimised in the second step (test cases 2c, d and 3a in Fig. 3). This shows that the inclusion
14 of s_1 observations can find a solution to counter the bias in s_2 and prevents a degradation in
15 the fit to the data. If s_2 is assimilated in the second step there is a negative impact on all other
16 variables as discussed above, demonstrating again that the order of data stream assimilation
17 can matter when there are biases or inconsistencies between the data and the model.

18 The analysis of the impact of the bias presented here is specific to this model and the
19 type and magnitude of the bias that was added, but the broader findings can be generalised to
20 any situation in which there is a bias or inconsistency between a model and data that is not
21 accounted for in the assigned error distributions. Exactly what might constitute a bias or
22 inconsistency is discussed more in Section 3.2. Also note that it is important to examine the
23 impact on the other variables. For the separate test case 1b in which only s_2 data are used to
24 optimise the model, the negative impact on the other variables (Fig. 3) would have been
25 concealed if we had only examined the posterior reduction in RMSE for the s_2 variable. Again
26 this is a concern that is inherent to all DA experiments, whether single- or multi-data stream,
27 but we can see from these results (i.e. by comparing the separate test cases 1b with 2a and b)
28 that adding another data stream in a multi-constraint approach does not always reduce the
29 problem.

30



1 2.2.3 Difference between the step-wise and simultaneous approaches in the 2 presence of a non-linear model

3 As discussed in Section 2.2.1, there is little difference between the step-wise and the
4 simultaneous assimilation approaches for simple, relatively linear models, unless the
5 observation error (including measurement and model errors) distribution deviates strongly
6 from the Gaussian assumption. However in reality, large-scale, complex LSMs may contain
7 highly non-linear responses to certain model parameters. To demonstrate the impact of non-
8 linearity in a multiple data stream assimilation context, we used a non-physically based toy
9 model chosen for its non-linear characteristics (see Section 2.1.2).

10 Fig. 4a shows the posterior parameter values for both the a and b parameters of the
11 non-linear toy model for all test cases. The values were not normalised as both parameters
12 have the same range. The horizontal dashed line shows the “true” known values of the
13 parameters (both equal to 1.0) that were used to generate the synthetic observations. Note that
14 no bias has been introduced into the model in the results described here. The prior and
15 posterior model s_1 and s_2 simulations for the non-linear toy model are compared to the
16 synthetic observations in Fig. 5 for both step-wise cases in which the posterior error
17 covariance matrix from step 1 (A_1 – see section 2.1.4) was propagated to step 2 (experiments
18 2a and c – Fig. 5a-d) and both simultaneous cases 3a and b (Fig. 5 e-h). Finally Fig. 6
19 summarises the reduction in RMSE between the simulated and observed s_1 and s_2 variables
20 for the non-linear toy model for all test cases and, in the step-wise cases, the reduction in
21 RMSE after both the first and second steps (light versus dark green bars).

22 Assimilating each data stream individually (test cases 1a and b) does not result in an
23 accurate retrieval of the posterior parameters (Fig. 4a), nor in a strong constraint on either
24 parameter, as shown by the lack of theoretical reduction in the parameter uncertainty after the
25 optimisation (Fig. 4b). Despite this, there is a 91-92% reduction in RMSE for the data stream
26 that was included in the optimisation (i.e. for s_1 in test case 1a – Fig. 6a, and s_2 in test case 1b
27 – Fig. 6b). However, the improvement on the other data stream is much less (28% reduction
28 in RMSE for s_1 when s_2 is assimilated) or even results in a degradation compared to the prior
29 fit (e.g. in the case of s_2 when s_1 is assimilated – Fig. 6b). Lack of improvement, or even
30 degradation, in the RMSE of other variables in the model is a common issue for data
31 assimilation in general – one that is not often evaluated in model-data fusion studies.



1 Only the simultaneous case, in which all s_1 observations have been included in the cost
2 function (test case 3a), manages to retrieve the correct parameter values after the optimisation.
3 All other posterior parameter values are incorrect, and are considerably different between
4 each case, unlike for the simple carbon model (without a model bias). Most step-wise test
5 cases (particularly 2b-d) do not result in the same parameter values as the simultaneous test
6 case 3a in which all the observations are included (Fig. 4a), highlighting that strong non-
7 linearity in the model sensitivity to parameters together with the use of an algorithm that is
8 only adapted to weakly non-linear problems, as well as the assumption of linearity in
9 calculating the posterior error covariance matrix at the minimum of the cost function, can
10 result in differences between a step-wise and simultaneous approach in multiple – data stream
11 assimilation (see Section 1).

12 In the simultaneous optimisation in which all observations are included (test case 3a)
13 the posterior fit to the data dramatically improves for both the s_1 and s_2 data streams after the
14 assimilation (blue dashed line in Fig. 5e and f). This was expected given that the correct
15 values of the parameters were found. For the step-wise cases (test case 2a in Figures 5a and b,
16 and test case 2c in Fig. 5c and d), the black dashed line shows the prior, and the posterior after
17 step 1 is shown by green dashed line. In the step-wise assimilation we see two different
18 scenarios depending on which data stream was assimilated first. In the first step the results are
19 the same as the case where each individual data stream is assimilated separately. In both cases
20 the first step results in a good fit to the data that was included in the optimisation in that step.
21 When the s_1 data was assimilated in the first step (Fig. 5 first row), the fit to s_2 deteriorated
22 after the optimisation (Fig. 5b green dashed line and Fig. 6b – test case 2a_s1), but when the
23 s_2 data were assimilated first (Fig. 5 second row) the optimisation step did manage to achieve
24 an improvement in the s_1 data stream (Fig. 5c green dashed line and Fig. 6a – test case 2c_s1).

25 In the second step the optimisation of s_2 in test cases 2a and b does not degrade the fit
26 to s_1 when the full parameter error covariance matrix (\mathbf{A}_1) is propagated between step 1 and 2
27 (Figures 5a blue curve and 6a 2a_s2). Furthermore optimising s_2 in the second step reverses
28 the deterioration in s_2 caused by assimilating s_1 in the first step (Figures 5b blue curve and 6b
29 2a and b dark green bars). However, when s_1 data were assimilated in the second step (test
30 cases 2c and d), we found that the good fit achieved with s_2 observations in the first step was
31 effectively reversed (Fig. 5d blue curve). Therefore assimilating s_1 in the second step
32 degraded the fit to the s_2 observations, even compared to the prior case (Fig. 6b, dark green



1 bars for test cases 2c and d). This nicely highlights one of the main possible issues with a
2 step-wise assimilation framework.

3 The fact that the final reduction in RMSE values after both steps was ~90% for most
4 cases, even though the values were not correct for all but case 3a (Fig. 4), indicates that the
5 error correlation between the two parameters (~ -1.0 – calculated from the posterior error
6 covariance matrix but not shown) led to alternative sets of values that resulted in a similar
7 improvement to the data – a phenomenon known as model equifinality.

8

9 2.2.4 Order of assimilation of data streams and propagation of parameter 10 error covariance matrices in a step-wise approach

11 Comparing the step-wise cases 2a and b with 2c and d for the non-linear toy model
12 reveals that neither order in the assimilation, s_1 then s_2 , or s_2 then s_1 , results in the correct
13 posterior parameter values that match the simultaneous test case (Fig. 4a). This is not a result
14 that can be generalised to all step-wise assimilations as it will depend on the data stream
15 involved and whether they contain enough information to accurately constrain all the
16 parameters included in the optimisation, as well as any biases in the model or observations (as
17 discussed in Section 2.2.2) or model non-linearity (section 2.2.3). In the case of the non-linear
18 toy model, neither s_1 nor s_2 find the right parameter values when assimilated individually,
19 therefore it is not surprising that neither order manages to achieve the right posterior
20 parameter values. Nevertheless, the theoretical uncertainty of both parameters is reduced by
21 >95% for the step-wise cases in which \mathbf{A}_1 from step 1 is propagated between step 1 and 2 (test
22 cases 2a and c – Fig. 4b), even though the posterior values for the step-wise cases are
23 incorrect. This demonstrates that a good theoretical reduction in uncertainty is not always
24 indicative that the right parameters have been found by the optimisation. The lower
25 theoretical reduction in parametric uncertainty for cases 2b and d (Fig. 4b) demonstrates that
26 information is lost between the steps if the posterior error covariance terms of \mathbf{A}_1 after step 1
27 are not propagated to step 2, and therefore cannot be used to further constrain the
28 optimisation.

29 From a mathematical standpoint the most rigorous approach is to propagate the full
30 parameter error covariance matrices between each step. Without that constraint not only is
31 information lost in the second step, but the information contained in the second data stream



1 may have a stronger influence compared to a simultaneous or step-wise case with a
2 propagated error covariance matrix. The inversion may therefore be more vulnerable to any
3 strong biases or incompatibilities between the model and the observations of the second data
4 stream, or indeed the particular sensitivity of its corresponding model state variable to the
5 parameters. This is one possible explanation for the degradation seen in s_1 in the non-linear
6 toy model when s_2 is optimised in the second step and \mathbf{A}_1 is not propagated between the steps
7 (Fig. 6a test case 2b_s2). The same was also true for the simple carbon model for test case 2b
8 when a bias was introduced into the s_2 simulation (see Section 2.2.2 and Fig. 3a).

9 However, the reverse is also true – if the first data stream contains strong biases then
10 the associated error correlations will be also propagated with \mathbf{A}_1 . If autocorrelation in the
11 observation errors, or indeed correlation between the errors of the data streams, is not
12 accounted for it is likely that the posterior simulations are over-tuned, i.e. we will
13 overestimate the reduction in parameter uncertainty. If this is the case and the first step results
14 in incorrect parameter values, the propagation of \mathbf{A}_1 could restrict the parameter values to the
15 wrong location in the parameter space and thus inhibit the ability of the inversion to find the
16 correct global minimum. These issues are likely to be more considerable for non-linear
17 models, as seen by the lack of difference between test cases 2a-d in the simple carbon model
18 example (Fig. 2).

19

20 2.2.5 Lessons to be learned when dealing with non-linearity

21 Most optimisation studies with a large-scale LSM use derivative methods based on a
22 least-squares approach, and therefore rely on assumptions of Gaussian probability and linear
23 model sensitivity. However, if the model is weakly non-linear within the probability
24 distribution around the point in parameter space that is being analysed (see Tarantola, 1987,
25 p72), it is possible to use an iterative algorithm, such as the one described in Eq. (6), to find
26 the minimum of the cost function (i.e. the maximum likelihood of the posterior parameter
27 distribution). Furthermore a linearization of the model around the maximum likelihood
28 estimation (minimum of $J(\mathbf{x})$) of the parameters can be used to calculate the posterior error
29 covariance (see Eq. (6)). If the model is too strongly non-linear and therefore these
30 assumptions are not met, it may not be possible to find the true global minimum of the cost
31 function and the characterisation of the posterior probability distribution will be incorrect.



1 This is a particular problem if the posterior parameter error covariance matrix is then
2 propagated in a step-wise approach, although these issues are relevant to both step-wise and
3 simultaneous assimilation. Note that performing a number of tests starting from different
4 random “first guess” points in parameter space can help to diagnose if the global minimum
5 has been reached, as outlined in Section 2.1.6 and discussed at the beginning of the results
6 (Section 2.2).

7 It is possible to avoid dealing with issues related to non-linearity in the model
8 sensitivity and non-Gaussian error distributions by using a global search method (e.g. Markov
9 Chain Monte Carlo or a genetic algorithm) that randomly, but effectively, searches the entire
10 parameter space. However in large dimensional problems, as is likely the case when
11 optimising a LSM at large scales with multiple data streams, using a global search method is
12 likely not feasible due to computational time constraints. In these cases, a derivative method
13 is likely the only option.

14 An important finding of the results presented for the non-linear toy model in Section
15 2.2.3 is that degradation in another data stream is not necessarily the result of a bias or
16 incompatibility between the observations and the model. Rather if the model sensitivity to the
17 parameters is very non-linear, multiple combinations of parameter values may exist that result
18 in a similar reduction of the cost function (multiple minima), but provide a different fit to
19 each data stream. Even though all data streams may be sensitive to all the parameters, the
20 information content of each will not be the same. Finding the true global minimum in this
21 instance may require a bit more careful thought in planning the assimilation set-up, and may
22 depend on having a reasonable idea of the level of information each data stream can bring to
23 constrain each parameter. It may be the case that one data stream has a higher non-linear
24 sensitivity to the parameters and therefore may act as a “troublemaker” and pull the
25 parameters in a direction that results in a degradation to the other data streams, as seen in
26 Section 2.2.3. If a simultaneous optimisation is not possible, it may be useful under such
27 circumstances to identify any “troublemaker” data streams, and assimilate them in the first
28 step of the optimisation. In the second step “peacemaker” data streams, with a lower non-
29 linear sensitivity to the parameters, will then find a compromise set of parameter values that
30 can fit both data streams well, provided the full posterior parameter error covariance matrix is
31 propagated between the steps in order to retain all the information brought by the first data
32 stream. As discussed this could be an explanation for the results seen for the non-linear toy



1 model test case 2a where s_1 was assimilated prior to s_2 (Figures 6a and b) as discussed in
2 Section 2.2.3.

3

4 **3 Examples from existing carbon cycle data assimilation studies**

5 **3.1 Extra constraint from multiple data streams**

6 Most site-based carbon cycle data assimilation studies have used eddy covariance
7 measurements of NEE and LE fluxes to constrain the relevant parameters of ecosystem
8 models. However, a few studies have also made use of chamber flux soil respiration data and
9 field measurements of vegetation characteristics (e.g. tree height, budburst, LAI) or estimates
10 of litterfall and carbon stocks as ancillary information (e.g. Keenan et al., 2012; Thum et al.,
11 in review; Van Oijen et al., 2005; Richardson et al., 2010; Williams et al., 2005). Two recent
12 studies combined high-resolution satellite-derived FAPAR data and in-situ eddy covariance
13 measurements to optimize parameters related to carbon, water and energy cycles of the
14 ORCHIDEE and BETHY LSMs (Bacour et al., 2015; Kato et al., 2013, respectively).

15 At global scales the number of studies that use multiple data streams from satellites or
16 large-scale networks to optimise LSMs has been increasing in recent years, although this
17 remains a relatively new area of research. CCDAS-BETHY was the first global carbon cycle
18 data assimilation system (CCDAS) making use of the high-precision measurements of the
19 atmospheric CO₂ concentration flask sampling network (Rayner et al., 2005; Scholze, 2003)
20 to constrain process parameters of the prognostic terrestrial carbon cycle model BETHY
21 (Knorr, 2000). Since its first application assimilating atmospheric CO₂ concentration data
22 only, CCDAS-BETHY has been further developed to consistently assimilate multiple data
23 streams both at local and global scales. In particular, Kaminski et al. (2012) optimised 70
24 process parameters plus one initial condition by simultaneously assimilating a satellite-
25 derived FAPAR product derived from the Medium Resolution Imaging Spectrometer
26 (MERIS; Gobron et al., 2008) and flask samples of atmospheric CO₂ at two sites from the
27 GLOBALVIEW product (GLOBALVIEW-CO₂, 2008) on a coarse resolution. More recently,
28 Scholze et al. (2015) demonstrated the added value of assimilating remotely sensed soil
29 moisture data in addition to observations of atmospheric CO₂ concentration from the flask-
30 sampling network. They used the same coarse resolution set-up of CCDAS as Kaminski et al.
31 (2012) and CO₂ observations from 10 sites of the GLOBALVIEW product (GLOBALVIEW-



1 CO₂, 2012) together with the SMOS L3 daily soil moisture product (version 246; CATDS-
2 L3, 2012).

3 Two other global CCDAS based on different LSMs have been developed in recent years
4 (Peylin et al., 2016; Schürmann et al., 2016). Schürmann et al. (2016) optimized model
5 parameters and initial conditions of the land component JSBACH (Raddatz et al. 2007) of the
6 MPI Earth System Model (ESM) (Giorgetta et al. 2013) using atmospheric CO₂ concentration
7 data and the TIP-FAPAR product (Pinty et al., 2007) as joint constraints over a 5 year period
8 in addition to evaluating the mutual benefit of each data stream in a fully factorial design.
9 Peylin et al. (2016) used three different data streams as global constraints for the ORCHIDEE
10 LSM (Krinner et al., 2005), which forms the land surface component of the IPSL ESM
11 (Dufresne et al., 2013), in a multi-site step-wise assimilation approach. First, satellite-derived
12 vegetation index data (NDVI) from the MODIS instrument was used to constrain the
13 phenology parameters at 60 sites for the temperate and boreal deciduous PFTs, followed by
14 NEE and LE observations at 78 FLUXNET sites for 7 PFTs to optimise all the carbon-related
15 parameters, and finally atmospheric CO₂ concentration measurements from 53 sites in the
16 GLOBALVIEW network (GLOBALVIEW-CO₂, 2013), which predominantly provided a
17 constraint on the initial magnitude of the soil carbon reserves in the model. Atmospheric CO₂
18 concentration observations are one of the most accurate, long-term data sets in environmental
19 science and they provide important information about the global CO₂ sink capacity by land
20 and ocean. These three global multiple data stream CCDAS have allowed an improvement in
21 both the mean seasonal cycle as well as the trend of net land surface CO₂ exchange, especially
22 with the inclusion of the atmospheric CO₂ data (Kaminski et al., 2012; Peylin et al., 2016;
23 Schürmann et al., 2016).

24 Many of the aforementioned studies reported that adding extra data streams helped to
25 constrain unresolved sub-spaces of the total parameter space. Scholze et al. (2015) found that
26 adding SMOS soil moisture data to the assimilation simultaneously with CO₂ observations
27 reduced the ambiguity in the solution space when assimilating CO₂ only, and the multiple
28 data constraint was able to resolve a much larger sub-space in parameter space (about 30
29 parameters out of the 101 compared to 15 without SMOS data). Bacour et al. (2015) and
30 Schürmann et al. (2016) both reported that the addition of FAPAR data brought extra
31 information on the phenology-related processes in the model, and therefore retrieved different
32 posterior C flux-related parameter values than when assimilating NEE or atmospheric CO₂



1 data alone. An interesting aspect of the Kaminski et al. (2012) study was that the inclusion of
2 FAPAR in addition to atmospheric CO₂ concentration samples resulted in a particular
3 improvement for the hydrological fluxes in the model, thus demonstrating the importance of
4 assessing the potential benefit for model variables that may not have been the main target of
5 optimisation. Richardson et al. (2010) concluded that using ancillary information (e.g. woody
6 biomass increment, field-based LAI and chamber measurements of soil respiration) as
7 orthogonal constraints to NEE data provided a valuable added constraint on many model
8 parameters, which improved both the bias in model predictions and reduced the associated
9 uncertainties. Thum et al. (in review) also found that the addition of aboveground biomass
10 stocks brought a longer-term constraint on allocation parameters and mortality/turnover
11 processes. However, they noted an incompatibility when assimilating both annual increment
12 and total biomass data, as the total stocks take into account losses related to disturbance and
13 management (e.g. canopy thinning) – processes that were not included in that version of the
14 model. Keenan et al. (2012) also argued that the use of such ancillary constraints is essential
15 for a better partitioning of net carbon fluxes into their gross components. However, Williams
16 et al. (2005) observed that one-off, or rarely taken, measurements of carbon stocks were
17 unable to constrain components of the carbon cycle to which they were not directly related.

18 This calls into question the issue of the influence of different data streams in a joint
19 assimilation, especially if the number of observations for each is vastly different which is the
20 case when assimilating both half-hourly C flux data in addition to soil C stock observations
21 that are typically available at an annual time scale. The spatial coverage of each data stream is
22 also important, especially for heterogeneous landscapes (Barrett et al., 2005). Test case 3b, in
23 which only one observation was included for the s_2 data stream instead of the complete time-
24 series, shows that a substantial difference in number of observations between the data streams
25 can influence the resulting parameter values and posterior uncertainty (compare test cases 3a
26 and b in Fig. 2 for the simple C model and Fig. 4 for the non-linear toy model) as each data
27 stream will have a different overall “weight” in the cost function. However, the impact of
28 having a different number of observations for each data stream in the cost function also
29 depends strongly on the prescribed observation error and relative sensitivity of each
30 corresponding model variable to the model parameters. If one variable has a greater
31 sensitivity than the other, it will matter less if fewer observations of that variable are included
32 in the cost function.



1 Xu et al. (2006), among others, have mentioned the possible need to weight the cost
2 function for different data sets. Different arguments abound on this issue. Some contend that
3 the cost function should not be weighted by the number of observations because the error
4 covariance matrices (**B** and **R**) already define this weight in an objective way (e.g. Keenan et
5 al., 2013). Certainly it should not be necessary to weight by the number of observations in the
6 cost function if there is sufficient information to properly build the prior error covariance
7 matrices (**B** and **R**). On the other hand, it is a difficult task to characterise the model structural
8 uncertainty and the observation error correlations (see Kuppel et al., 2013 for practical
9 solutions). Given this, our expert knowledge on the workings of the model processes and the
10 sensitivity of the model to the parameters may permit us to specify a stronger weight to a data
11 stream that could help to constrain a particular section of the model, but for which there are
12 only a few data points. Clearly the definition of the prior error model, including for the
13 covariance between errors of the data streams, is of the utmost importance (Trudinger et al.,
14 2007) and merits close attention in future multiple data stream assimilation studies.

15 Although a number of multiple data stream assimilation studies exist at various scales,
16 very few studies have specifically investigated the added benefit of different combinations of
17 data streams, with a few notable exceptions (Barrett et al., 2005; Richardson et al., 2010; Kato
18 et al., 2013; Keenan et al., 2013; Bacour et al., 2015; Schürmann et al., 2016). Kato et al.
19 (2013) and Bacour et al. (2015) both evaluated the complementarity of eddy covariance and
20 FAPAR data streams at site level, i.e. the impact of assimilating one individual data stream on
21 the other model state variable, as well as when both data streams were included in the
22 optimization (see discussion in Section 3.2). The study of Keenan et al. (2013) was
23 particularly notable in its aim to quantify which data streams provide the most information
24 and how many data streams are actually needed to constrain the problem. They reported that
25 of the 17 field-based data streams available, projections of future carbon dynamics were well-
26 constrained with only 5 of the data sources, and crucially, not with eddy covariance NEE
27 measurements alone. These results may be specific to this site or type of ecosystem, but this
28 study highlights the need for further research in this area, and in particular, for synthetic data
29 experiments that allow us to understand which data will be the most useful for a given
30 scientific question. This will also enable researchers to plan more efficient measurement
31 campaigns with experimentalists, as also pointed out by Keenan et al. (2012).

32



1 **3.2 Issue of bias and inconsistencies between the observations and the** 2 **model**

3 Despite the theoretical benefit of adding data streams into an assimilation system as
4 orthogonal constraints, several of the aforementioned studies at both site and global scale
5 have reported a bias or inconsistency either between the different observation data streams, or
6 between the observations and the model. This is easily detected when the optimisation of one
7 data stream results in a worse fit than the prior in one or more of the other data streams, as
8 seen in Section 2.2.2. Kato et al. (2013) assimilated SeaWiFS FAPAR (Gobron et al., 2006)
9 and eddy covariance LE measurements at the FLUXNET site in Maun, Botswana. They
10 showed that the individual assimilation of each the two data streams resulted in a perfect (i.e.
11 within the observational uncertainty) fit to the assimilated data set, but a considerable
12 degradation of the fit to the non-assimilated data set compared to the prior. A comparison
13 against eddy covariance measurements of gross carbon uptake (gross primary production –
14 GPP) hinted to a bias problem with the FAPAR data because the fit to the independent GPP
15 data was degraded after assimilating FAPAR data only, while the fit improved after
16 assimilating the LE data only. Nevertheless, the simultaneous assimilation of both data
17 streams achieved a compromise between the two suboptimal states reached after assimilating
18 only one data stream. The calibration further limited the number of parameters with correlated
19 errors, and yielded a higher theoretical reduction in parameter uncertainty and a decrease in
20 the RMS difference by 16% for the GPP data compared to the prior.

21 Bacour et al. (2015) also noted that when assimilating both in-situ and satellite-derived
22 FAPAR data (from the SPOT and MERIS instruments) and in-situ NEE and LE flux data
23 from two French FLUXNET sites into the ORCHIDEE LSM both separately and together, the
24 posterior parameter values changed significantly for the photosynthesis and phenology-related
25 parameters, depending on the bias between the model and the observations and the correlation
26 between the parameter errors. If NEE data were assimilated alone there was an even stronger
27 positive bias (model–observations) in the start of leaf onset in the FAPAR data than in the
28 prior simulations, and no improvement in the maximum value. This was likely due to the fact
29 that there were enough degrees of freedom to fit the NEE without changing the phenology-
30 related parameters. Similarly, the fit to the NEE was degraded when the model was only
31 optimized with FAPAR data. The model was able to fit the maximum FAPAR but this
32 resulted in an adverse effect on the carbon assimilation capacity of the vegetation. The



1 authors argued this was related to incompatibilities between the FAPAR and both the model
2 and NEE measurements, possibly due to its larger spatial footprint of the satellite-derived
3 FAPAR data and/or inaccuracies in the retrieval algorithm. However, given that assimilating
4 in-situ FAPAR also degraded the fit to the NEE, another culprit may be an inconsistency
5 between the model and the data. The authors suggested this could be due to the different
6 assumptions or characterisation of a variable in a model compared to what is described in the
7 data. For example, satellite-derived greenness measures (FAPAR/NDVI) also contain
8 information on the non-green elements of vegetation, but the model only simulates green LAI.
9 Furthermore parameters and processes in models have been developed at certain temporal and
10 spatial scales. Vegetation is often simply represented as a “big leaf” model in LSMs, taking
11 no account of vertical canopy structure or the spatial heterogeneity in a scene, which is an
12 additional source of inconsistency with what is measured. The joint (simultaneous)
13 assimilation of all three data streams in Bacour et al. (2015) reconciled the different sources
14 of information, with an improvement in the model-data fit for NEE, LE and FAPAR.
15 However, the compromise achieved in the joint assimilation was only possible when the
16 FAPAR data were normalised to their maximum and minimum values, which thus partially
17 accounted for any bias in the magnitude of the FAPAR or inconsistency with the model.

18 The story of biases and apparent inconsistencies in FAPAR data doesn't end there. A
19 bias correction was also necessary in the study by Kaminski et al. (2012) with CCDAS-
20 BETHY using the MERIS FAPAR product in addition to atmospheric CO₂ data (see above).
21 They found that optimisation procedure failed when using the original FAPAR product
22 because the FAPAR values were biased towards higher values. Only after applying a bias
23 correction on the FAPAR data before the assimilation procedure was the optimisation
24 successful. Schürmann et al. (2016) also reported the need to reduce a prior model bias in
25 FAPAR. Even though the assimilation corrected successfully for this FAPAR bias, an imprint
26 of the prior bias was evident in the spatial patterns of the modelled heterotrophic respiration.
27 Assimilating FAPAR data alone therefore resulted in a slight degradation in the net C flux
28 and consequently led to incorrect simulations of the atmospheric CO₂ growth rate. The
29 addition of CO₂ as a constraint prevented this degradation and resulted in a compromise in
30 which FAPAR helped to disentangle these processes and find different parameter values
31 compared to the CO₂-only case, thus improving the fit to both data streams. Forkel et al.
32 (2014) discovered an apparent inconsistency between satellite-derived FAPAR and GPP data
33 in tundra regions when using these data (plus satellite-derived albedo) to optimise the LPJmL



1 LSM. They too speculated that the data might be positively biased, in this case due to issues
2 with satellite measurements taken at high sun zenith angles. However, they gave alternative
3 suggestions, one being that an inadequate model structure may be at fault – for example, the
4 LPJmL does not include vegetation classes corresponding to shrub, moss and lichen species
5 that are dominant in these ecosystems. They also noted that the GPP product they used, which
6 is based on a model tree ensemble up-scaling of FLUXNET data (Jung et al., 2011), might
7 contain representation-related biases, given that there are very few FLUXNET stations in
8 tundra regions. The issue of representation errors of sites has been touched upon before (e.g.
9 Raupach et al., 2005). Alton (2013), who performed a global multi-site optimisation of the
10 JULES LSM with a diverse range of data including satellite-derived LAI, FLUXNET, soil
11 respiration and global river discharge, raised the point that FLUXNET sites are known to be
12 large carbon sinks, which could potentially result in biased global NEE estimates. Resolving
13 these apparent inconsistencies was beyond the scope of most of these studies, aside from
14 applying a bias correction where one was evident. Nonetheless this issue clearly merits further
15 attention if the increasing number of available datasets is to be fully utilised.

16

17 **3.3 Step-wise versus simultaneous assimilation**

18 The paper by Alton (2013) documents the only previous study to have used a step-wise
19 assimilation approach with more than two data streams, stating that the final parameter values
20 were independent of the order of data streams assimilated. No studies in the LSM community
21 to date have explicitly examined a step-wise versus simultaneous assimilation framework
22 with the same optimisation system and model. The step-wise assimilation with the
23 ORCHIDEE-CCDAS detailed in Peylin et al. (2016) has been compared to a simultaneous
24 optimisation using the same three data streams as part of an on-going study. At each step, the
25 resulting simulations (using the posterior parameters) were compared to the data stream from
26 the previous steps. The fit to the MODIS NDVI (used in a similar manner to FAPAR as a
27 proxy for vegetation greenness) was unchanged after further optimization of the phenology-
28 related parameters in the second and third steps using in-situ flux and atmospheric CO₂
29 concentration data. In the simultaneous optimisation, the addition of NEE or atmospheric CO₂
30 concentration measurements resulted in a lower improvement to the fit to MODIS NDVI. As
31 the NDVI data were normalised this was not a result of a simple bias in the magnitude of the
32 data. Rather, it was likely due to inconsistencies between the model and data as discussed by



1 Bacour et al. (2015) and in Section 3.2. It is important to reiterate that there should be no
2 difference between the step-wise and the simultaneous given an adequate description of the
3 error covariance matrices and compliance with the assumptions associated with the inversion
4 algorithm used. However, in practice it is very difficult to define a PDF that properly
5 characterises the model structural uncertainty and observation errors accounting for biases and
6 non-Gaussian distributions. This leads to issues within a simultaneous assimilation,
7 particularly if the information content of one data stream is much higher, and a greater risk of
8 differences between a step-wise and simultaneous assimilation. As discussed in Section 2.2.5
9 a step-wise assimilation may be useful on a provisional basis for dealing with possible
10 inconsistencies. In the step-wise approach of Peylin et al. (2016) the error covariance of the
11 phenology-related parameters was strongly constrained by the satellite data in the first step
12 (and was propagated to the second step), the later assimilations with NEE and atmospheric
13 CO₂ data in steps 2 and 3 found alternative solutions for the C flux-related parameters that
14 provided a reasonable fit to all data streams. Wherever possible however, a simultaneous
15 optimisation is favourable because the strong parameter linkages between different processes
16 are maintained, and therefore biases and inconsistencies between the model and observations
17 should be addressed prior to optimisation.

18

19 **4 Advice for Land Surface Modellers**

20 Although it is clear that in many cases, increasing the number of observations in a model
21 optimisation provides additional, orthogonal constraints, challenges remain that should not be
22 ignored. Based on the simple toy model results presented in this study, in addition to lessons
23 learned from existing studies, we recommend the following points when carrying out multiple
24 data stream carbon cycle data assimilation experiments:

- 25 • Devote time to characterising the error structure for the observations and
26 parameter error distributions, including their correlations (Raupach et al., 2005).
27 For the observations this should include the model structural errors (Kuppel et al.
28 2013), the temporal or spatial autocorrelation and correlation between different
29 data streams.
- 30 • In the case of non-Gaussian error distributions consider performing a
31 transformation to make the distributions more Gaussian, or avoid a least squares



- 1 formulation and instead use a method that avoids outliers (e.g. absolute deviations
2 – Trudinger et al., 2007).
- 3 • Analyse and correct for biases in the observations, or approximately account for it
4 in the observation error covariance matrix, \mathbf{R} , using the off-diagonal terms or
5 inflated errors (Chevallier, 2007), or by using the prior model-data RMSE to
6 define the observation uncertainty.
 - 7 • Investigate potential incompatibilities between your model and data. Take time to
8 understand which physical quantities your data correspond to and whether that is
9 consistent with the description of the equivalent variable in the model. As for the
10 previous point, one way of attempting to account for unknown inconsistencies
11 between the model and data is to set the observation uncertainty, \mathbf{R} , the prior
12 RMSE between the model and the data.
 - 13 • Evaluate the impact on other model variables with independent observations, and
14 if the optimisation degrades the fit compared to the prior, investigate the reasons
15 behind the inconsistency and address them as above.
 - 16 • Assess the non-linearity of your model (multiple first guess tests can help in this
17 regard), and if strongly so, avoid a least squares formulation of the cost function or
18 use global search algorithms for the optimisation – although at the resolution of
19 typical LSM simulations ($\geq 0.5 \times 0.5^\circ$) this will likely only be computationally
20 feasible at site or multi-site scale.
 - 21 • Prior information is key in a Bayesian framework. Effort should be put into better
22 constraining the prior parameter bounds of all parameters based on literature
23 wherever possible.
 - 24 • Conduct preliminary sensitivity analyses to determine which parameters should be
25 constrained by each data stream.
 - 26 • Set up experiments with synthetic data, as in this study, to understand the
27 constraints posed by the different data streams you will include in the experiment.
 - 28 • If technical constraints require a step-wise approach is used it is preferable (from a
29 mathematical standpoint) to propagate the full parameter error covariance matrix



1 between each step, if it can be calculated, and carefully consider the order of the
2 assimilation of data streams (a synthetic experiment will aid in this regard).

- 3 • Be aware that a good theoretical reduction in model or parameter uncertainty can
4 be misleading, as it is not necessarily indicative that the right parameter values
5 have been found. If this is the case, it could impact predictions made outside the
6 spatio-temporal window included in the optimisation.

7

8 Many of these issues are relevant to any data assimilation study, including those only
9 using one data stream. However, most are more pertinent when considering more than one
10 source of data. The impact of bias in the magnitude of satellite-derived FAPAR data has
11 featured highly in past multiple data stream assimilation studies. Aside from simple
12 corrections, Quaife et al. (2008) and Zobitz et al. (2014) suggested that LSMs should be
13 coupled to radiative transfer models to provide a more realistic and mechanistic observation
14 operator between the quantities simulated by the model and the raw radiance measured by
15 satellite instruments. This proposition followed the experience gained in the case of
16 atmospheric models for several decades (Morcrette, 1991).

17 Other promising directions could also be considered to help constrain the problem of lack
18 of information in resolving the parameter space, including the use of other ecological and
19 dynamical “rules” that limit the optimisation (see for example Bloom and Williams, 2015), or
20 the addition of different timescales of information extracted from the data such as annual
21 sums (e.g. Keenan et al., 2012). Of course, optimising the parameters of the model will not
22 account for all the uncertainty in a model. Inaccurate or incomplete process representation is
23 likely a key factor that may also bias the posterior values retrieved in any optimisation.
24 Keenan et al. (2012) reflected that despite using multiple different constraints and different
25 time increments in the cost function, the inter-annual variability and long-term trend of carbon
26 uptake at Harvard forest FLUXNET site in the USA could not be reproduced without a
27 temporal variation of the parameters, suggesting a missing process in the model. However, as
28 this paper shows, the complexities of model-data fusion require that we continue to develop
29 DA techniques alongside development of LSMs, with the hope of converging upon more
30 reliable and accurate predictions of the global C budget in the near future. Finally we should
31 also seek to develop collaborations with researchers in other fields who may have advanced
32 further in a particular direction. Members of the atmospheric and hydrological modelling



1 communities, for example, have implemented techniques for inferring the properties of the
2 prior error covariance matrices, including the mean and variance, but also potential biases,
3 autocorrelation and heteroscedasticity, by including these terms as “hyper-parameters” within
4 the inversion (e.g. Michalak et al. 2005; Evin et al., 2014; Renard et al., 2010; Wu et al.
5 2013;). Of course this extends the parameter space – making the problem harder to solve
6 unless sufficient prior information is available (Renard et al., 2010), but such avenues are
7 worth exploring.

8

9 **5 Conclusions**

10 In this study we have attempted to highlight and discuss some of the challenges
11 associated with using multiple data streams to constrain the parameters of LSMs, with a
12 particular focus on the carbon cycle. We demonstrated some of the issues using two simple
13 models constrained with synthetic observations for which the ‘true’ parameters are known.
14 We performed a variety of tests in Section 2 to demonstrate the differences between
15 assimilating each data stream separately, sequentially (in a step-wise approach) and together
16 in the same assimilation (simultaneous approach). In particular we focused on difficulties that
17 may arise in the presence of biases or inconsistencies between the data and the model, as well
18 as non-linearity in the model equations. In Section 3 we discussed the experimental results
19 with reference to similar difficulties that have been documented in recent C cycle assimilation
20 studies.

21 Many of the issues faced are inherent to all optimisation experiments, including those in
22 which only one data stream is used. It is of utmost importance to determine if the
23 observations contain biases, and/or if inconsistencies or incompatibilities exist between the
24 model and the observations, and to correct for this or properly account for this in the error
25 covariance matrices. Secondly it is crucial to understand the assumptions and limitations
26 related to the inversion algorithm used. Without these two points being met, there is a greater
27 risk of obtaining incorrect parameter values, which may not be obvious by examining the
28 posterior uncertainty and model-data RMSE reduction. Furthermore it is more likely that the
29 implementation of a step-wise versus simultaneous approach will lead to different results.

30 This study was not able to examine an exhaustive list of all possible challenges that may
31 be faced when assimilating multiple data streams, but we hope that this tutorial style paper
32 will serve as a guide for those wishing to optimise the parameters of LSMs using the variety



1 of C cycle related observations that are available today. Furthermore we hope that by
2 increasing awareness about the possible difficulties of model-data integration, we can further
3 bring the modelling and experimental communities together to work more closely on these
4 issues.

5

6 **Code availability**

7 The model and inversion code will be made available via the ORCHIDAS website (upon
8 registration): https://orchidas.lsce.ipsl.fr/multi_data_stream.php.

9

10

11 **Acknowledgements**

12 We acknowledge the support from the International Space Science Institute (ISSI). This
13 publication is an outcome of the ISSI's Working Group on "Carbon Cycle Data Assimilation:
14 How to Consistently Assimilate Multiple Data Streams". N. MacBean was also funded by the
15 GEOCARBON Project (ENV.2011.4.1.1-1-283080) within the European Union's 7th
16 Framework Programme for Research and Development. The authors wish to thank colleagues
17 and collaborators in the atmospheric inversion and C cycle DA communities with whom they
18 have had numerous past conversations that have led to an improvement in their understanding
19 of the issues presented here.

20



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 Anav, A., Friedlingstein, P., Kidston, M., Bopp, L., Ciais, P., Cox, P., Jones, C., Jung, M.,
6 Myneni, R. and Zhu, Z.: Evaluating the land and ocean components of the global carbon cycle
7 in the CMIP5 earth system models, *J. Clim.*, 26(18), 6801–6843, doi:10.1175/JCLI-D-12-
8 00417.1, 2013.
- 9 Bacour, C., Peylin, P., MacBean, N., Rayner, P. J., Delage, F., Chevallier, F., Weiss, M.,
10 Demarty, J., Santaren, D., Baret, F., Berveiller, D., Dufrêne, E. and Prunet, P.: Joint
11 assimilation of eddy covariance flux measurements and FAPAR products over temperate
12 forests within a process-oriented biosphere model, *J. Geophys. Res. Biogeosciences*, 120,
13 1839–1857, doi:10.1002/2015JG002966. Received, 2015.
- 14 Barrett, D. J., Michael J Hill, I., Hutley, L. B., Beringer, J., Xu, J. H., Cook, G. D., Carter, J.
15 O. and Williams, R. J.: Prospects for improving savanna biophysical models by using
16 multiple-constraints model-data assimilation methods, *Aust. J. Bot.*, 53(7), 689–714,
17 doi:10.1071/BT04139, 2005.
- 18 Bloom, A. A. and Williams, M.: Constraining ecosystem carbon dynamics in a data-limited
19 world: integrating ecological ‘common sense’ in a model–data fusion framework,
20 *Biogeosciences*, 12(5), 1299–1315, doi:10.5194/bg-12-1299-2015, 2015.
- 21 Chevallier, F., 2007: Impact of correlated observation errors on inverted CO₂ surface fluxes
22 from OCO measurements, *Geophys. Res. Lett.*, 34, L24804, doi:10.1029/2007GL030463.
- 23 Dufresne, J. L., Foujols, M. a., Denvil, S., Caubel, a., Marti, O., Aumont, O., Balkanski, Y.,
24 Bekki, S., Bellenger, H., Benschila, R., Bony, S., Bopp, L., Braconnot, P., Brockmann, P.,
25 Cadule, P., Cheruy, F., Codron, F., Cozic, a., Cugnet, D., de Noblet, N., Duvel, J. P., Ethé,
26 C., Fairhead, L., Fichefet, T., Flavoni, S., Friedlingstein, P., Grandpeix, J. Y., Guez, L.,
27 Guilyardi, E., Hauglustaine, D., Hourdin, F., Idelkadi, a., Ghattas, J., Jousaume, S.,
28 Kageyama, M., Krinner, G., Labetoulle, S., Lahellec, a., Lefebvre, M. P., Lefevre, F., Levy,
29 C., Li, Z. X., Lloyd, J., Lott, F., Madec, G., Mancip, M., Marchand, M., Masson, S.,
30 Meurdesoif, Y., Mignot, J., Musat, I., Parouty, S., Polcher, J., Rio, C., Schulz, M.,
31 Swingedouw, D., Szopa, S., Talandier, C., Terray, P., Viovy, N. and Vuichard, N.: Climate



- 1 change projections using the IPSL-CM5 Earth System Model: From CMIP3 to CMIP5., 2013.
- 2 Evin, G., Thyer, M., Kavetski, D., McInerney, D. and Kuczera, G.: Comparison of joint
3 versus postprocessor approaches for hydrological uncertainty estimation accounting for error
4 autocorrelation and heteroscedasticity, *Water Resour. Res.*, 50(3), 2350–2375,
5 doi:10.1002/2013WR014185, 2014.
- 6 Forkel, M., Carvalhais, N., Schaphoff, S., v. Bloh, W., Migliavacca, M., Thurner, M. and
7 Thonicke, K.: Identifying environmental controls on vegetation greenness phenology through
8 model-data integration, *Biogeosciences*, 11, 7025–7050, doi:10.5194/bg-11-7025-2014, 2014.
- 9 Gobron, N., Pinty, B., Ausedat, O., Chen, J. M., Cohen, W. B., Fensholt, R., Gond, V.,
10 Huemmrich, K. F., Lavergne, T., Mélin, F., Privette, J. L., Sandholt, I., Taberner, M., Turner,
11 D. P., Verstraete, M. M. and Widlowski, J. L.: Evaluation of fraction of absorbed
12 photosynthetically active radiation products for different canopy radiation transfer regimes:
13 Methodology and results using Joint Research Center products derived from SeaWiFS against
14 ground-based estimations, *J. Geophys. Res.*, 111, D13110, doi:10.1029/2005JD006511, 2006.
- 15 Gobron, N., Pinty, B., Ausedat, O., Taberner, M., Faber, O., Mélin, F., Lavergne, T.,
16 Robustelli, M. and Snoeij, P.: Uncertainty estimates for the FAPAR operational products
17 derived from MERIS - Impact of top-of-atmosphere radiance uncertainties and validation
18 with field data, *Remote Sens. Environ.*, 112(4), 1871–1883, doi:10.1016/j.rse.2007.09.011,
19 2008.
- 20 Kaminski, T., Knorr, W., Scholze, M., Gobron, N., Pinty, B., Giering, R. and Mathieu, P. P.:
21 Consistent assimilation of MERIS FAPAR and atmospheric CO₂ into a terrestrial vegetation
22 model and interactive mission benefit analysis, *Biogeosciences*, 9(8), 3173–3184,
23 doi:10.5194/bg-9-3173-2012, 2012.
- 24 Kato, T., Knorr, W., Scholze, M., Veenendaal, E., Kaminski, T., Kattge, J. and Gobron, N.:
25 Simultaneous assimilation of satellite and eddy covariance data for improving terrestrial water
26 and carbon simulations at a semi-arid woodland site in Botswana, *Biogeosciences*, 10(2),
27 789–802, doi:10.5194/bg-10-789-2013, 2013.
- 28 Keenan, T. F., Davidson, E., Moffat, A. M., Munger, W. and Richardson, A. D.: Using
29 model-data fusion to interpret past trends, and quantify uncertainties in future projections, of
30 terrestrial ecosystem carbon cycling, *Glob. Chang. Biol.*, 18(8), 2555–2569,
31 doi:10.1111/j.1365-2486.2012.02684.x, 2012.



- 1 Keenan, T. F., Davidson, E. a., Munger, J. W. and Richardson, A. D.: Rate my data:
2 Quantifying the value of ecological data for the development of models of the terrestrial
3 carbon cycle, *Ecol. Appl.*, 23(1), 273–286, doi:10.1890/12-0747.1, 2013.
- 4 Knorr, W.: Annual and interannual CO₂ exchanges of the terrestrial biosphere: process-based
5 simulations and uncertainties, *Glob. Ecol. Biogeogr.*, 9(3), 225–252, doi:10.1046/j.1365-
6 2699.2000.00159.x, 2000.
- 7 Krinner, G., Viovy, N., de Noblet-Ducoudré, N., Ogée, J., Polcher, J., Friedlingstein, P.,
8 Ciais, P., Sitch, S. and Prentice, I. C.: A dynamic global vegetation model for studies of the
9 coupled atmosphere-biosphere system, *Global Biogeochem. Cycles*, 19(1), 1–33,
10 doi:10.1029/2003GB002199, 2005.
- 11 Kuppel, S., F. Chevallier and P. Peylin,: Quantifying the model structural error in Carbon
12 Cycle Data Assimilation Systems. *Geosci. Model Dev.*, 6, 45-55, doi:10.5194/gmd-6-45-
13 2013, 2013.
- 14 Michalak, A. M., Hirsch, A., Bruhwiler, L., Gurney, K. R., Peters, W. and co-authors:
15 Maximum likelihood estimation of covariance parameters for Bayesian atmospheric trace gas
16 surface flux inversions. *J. Geophys. Res.* 110, D24107. DOI: 10.1029/2005JD005970, 2005.
- 17 Morcrette, J.-J.: Evaluation of Model-generated Cloudiness: Satellite-observed and Model-
18 generated Diurnal Variability of Brightness Temperature. *Mon. Wea. Rev.*, **119**, 1205–1224,
19 1991.
- 20 van Oijen, M., Rougier, J. and Smith, R.: Bayesian calibration of process-based forest models:
21 bridging the gap between models and data, *Tree Physiol.*, 25(7), 915–927,
22 doi:10.1093/treephys/25.7.915, 2005.
- 23 Peylin, P., Bacour, C., MacBean, N., Leonard, S., Rayner, P. J., Kuppel, S., Koffi, E. N.,
24 Kane, A., Maignan, F., Chevallier, F., Ciais, P., and Prunet, P.: A new step-wise Carbon
25 Cycle Data Assimilation System using multiple data streams to constrain the simulated land
26 surface carbon cycle, *Geosci. Model Dev. Discuss.*, doi:10.5194/gmd-2016-13, in review,
27 2016.
- 28 Quaife, T., Lewis, P., De Kauwe, M., Williams, M., Law, B. E., Disney, M. and Bowyer, P.:
29 Assimilating canopy reflectance data into an ecosystem model with an Ensemble Kalman
30 Filter, *Remote Sens. Environ.*, 112(4), 1347–1364, doi:10.1016/j.rse.2007.05.020, 2008.



- 1 Raupach, M. R.: Dynamics of resource production and utilisation in two-component
2 biosphere-human and terrestrial carbon systems, *Hydrol. Earth Syst. Sci.*, 11, 875–889,
3 doi:10.5194/hess-11-875-2007, 2007.
- 4 Raupach, M. R., Rayner, P. J., Barrett, D. J., Defries, R. S., Heimann, M., Ojima, D. S.,
5 Quegan, S. and Schimmler, C. C.: Model-data synthesis in terrestrial carbon observation:
6 Methods, data requirements and data uncertainty specifications, *Glob. Chang. Biol.*, 11(3),
7 378–397, doi:10.1111/j.1365-2486.2005.00917.x, 2005.
- 8 Rayner, P. J., Scholze, M., Knorr, W., Kaminski, T., Giering, R. and Widmann, H.: Two
9 decades of terrestrial carbon fluxes from a carbon cycle data assimilation system (CCDAS), ,
10 19, doi:10.1029/2004GB002254, 2005.
- 11 Renard, B., Kavetski, D., Kuczera, G., Thyer, M. and Franks, S. W.: Understanding predictive
12 uncertainty in hydrologic modeling: The challenge of identifying input and structural errors,
13 *Water Resour. Res.*, 46(5), 1–22, doi:10.1029/2009WR008328, 2010.
- 14 Richardson, A. D., Williams, M., Hollinger, D. Y., Moore, D. J. P., Dail, D. B., Davidson, E.
15 a., Scott, N. a., Evans, R. S., Hughes, H., Lee, J. T., Rodrigues, C. and Savage, K.: Estimating
16 parameters of a forest ecosystem C model with measurements of stocks and fluxes as joint
17 constraints, *Oecologia*, 164(1), 25–40, doi:10.1007/s00442-010-1628-y, 2010.
- 18 Schürmann, G. J., Kaminski, T., Köstler, C., Carvalhais, N., Voßbeck, M., Kattge, J., Giering,
19 R., Rödenbeck, C., Heimann, M., and Zaehle, S.: Constraining a land surface model with
20 multiple observations by application of the MPI-Carbon Cycle Data Assimilation System,
21 *Geosci. Model Dev. Discuss.*, doi:10.5194/gmd-2015-263, in review, 2016.
- 22 Sitch, S., Friedlingstein, P., Gruber, N., Jones, S. D., Murray-Tortarolo, G., Ahlström, A.,
23 Doney, S. C., Graven, H., Heinze, C., Huntingford, C., Levis, S., Levy, P. E., Lomas, M.,
24 Poulter, B., Viovy, N., Zaehle, S., Zeng, N., Arneth, A., Bonan, G., Bopp, L., Canadell, J. G.,
25 Chevallier, F., Ciais, P., Ellis, R., Gloor, M., Peylin, P., Piao, S., Le Quéré, C., Smith, B.,
26 Zhu, Z. and Myneni, R.: Recent trends and drivers of regional sources and sinks of carbon
27 dioxide, *Biogeosciences*, 12, 653–679, doi:10.5194/bgd-12-653-2015, 2015.
- 28 Thum, T., N. MacBean, P. Peylin, C. Bacour, D. Santaren, B. Longdoz, D. Loustau and P.
29 Ciais, The potential benefit of using forest biomass data in addition to carbon and water flux
30 measurements to constrain ecosystem model parameters: case studies at two temperate forest
31 sites. In revision for *Agric. For. Meteorol.*



- 1 Trudinger, C. M., Raupach, M. R., Rayner, P. J., Kattge, J., Liu, Q., Park, B., Reichstein, M.,
2 Renzullo, L., Richardson, A. D., Roxburgh, S. H., Styles, J., Wang, Y. P., Briggs, P., Barrett,
3 D. and Nikolova, S.: OptIC project: An intercomparison of optimization techniques for
4 parameter estimation in terrestrial biogeochemical models, *J. Geophys. Res. Biogeosciences*,
5 112(2), doi:10.1029/2006JG000367, 2007.
- 6 Williams, M., Schwarz, P. a, Law, B. E., Irvine, J. and Kurpius, M. R.: An improved analysis
7 of forest carbon dynamics using data assimilation, *Glob. Chang. Biol.*, 11(1), 89–105,
8 doi:10.1111/j.1365-2486.2004.00891.x, 2005.
- 9 Wu, L., M. Bocquet, F. Chevallier, T. Lauvaux, and K. Davis, 2013: Hyperparameter
10 estimation for uncertainty quantification in mesoscale carbon dioxide inversions. *Tellus B*, 65,
11 doi:10.3402/tellusb.v65i0.20894.
- 12 Xu, T., White, L., Hui, D. and Luo, Y.: Probabilistic inversion of a terrestrial ecosystem
13 model: Analysis of uncertainty in parameter estimation and model prediction, *Global*
14 *Biogeochem. Cycles*, 20(2), 1–15, doi:10.1029/2005GB002468, 2006.
- 15 Zobitz, J. M., Moore, D. J. P., Quaife, T., Braswell, B. H., Bergeson, A., Anthony, J. a. and
16 Monson, R. K.: Joint data assimilation of satellite reflectance and net ecosystem exchange
17 data constrains ecosystem carbon fluxes at a high-elevation subalpine forest, *Agric. For.*
18 *Meteorol.*, 195–196, 73–88, doi:10.1016/j.agrformet.2014.04.011, 2014.

19

20

21

22

23

24

25

26

27

28

29



1 Table 1: The optimisation set-up for both models, including the true parameter values, their
 2 range and the observation uncertainty (1 sigma). The parameter uncertainty (1 sigma) was set
 3 to 40% of the range for each parameter.

4

Model	Parameter value (range)				Observation uncertainty	
Simple carbon model	p_1 1 (0.5,5)	p_2 1 (0.5,5)	k_1 0.2 (0.03,0.9)	k_2 0.1 (0.01,0.12)	s_1 0.5	s_2 5
Non-linear toy model	a 1 (0,2)		b 1 (0,2)		s_1 0.5	s_2 0.5

5

6

7

8

9

10

11

12

13

14

15

16

17

18

19

20

21



1 Table 2: List of experiments performed for both models with synthetic data. All parameters
 2 are optimised in all cases (therefore in both steps for the step-wise approach).

3

Test case	Step 1	Step 2	Parameter error covariance terms propagated in step 2?
<i>Separate</i>			
1a	s_1	-	-
1b	s_2	-	-
<i>Step-wise</i>			
2a	s_1	s_2	yes
2b	s_1	s_2	no
2c	s_2	s_1	yes
2d	s_2	s_1	no
<i>Simultaneous</i>			
3a	s_1 and s_2	-	-
3b	s_1 and only 1 obs for s_2	-	-

4

5

6

7

8

9

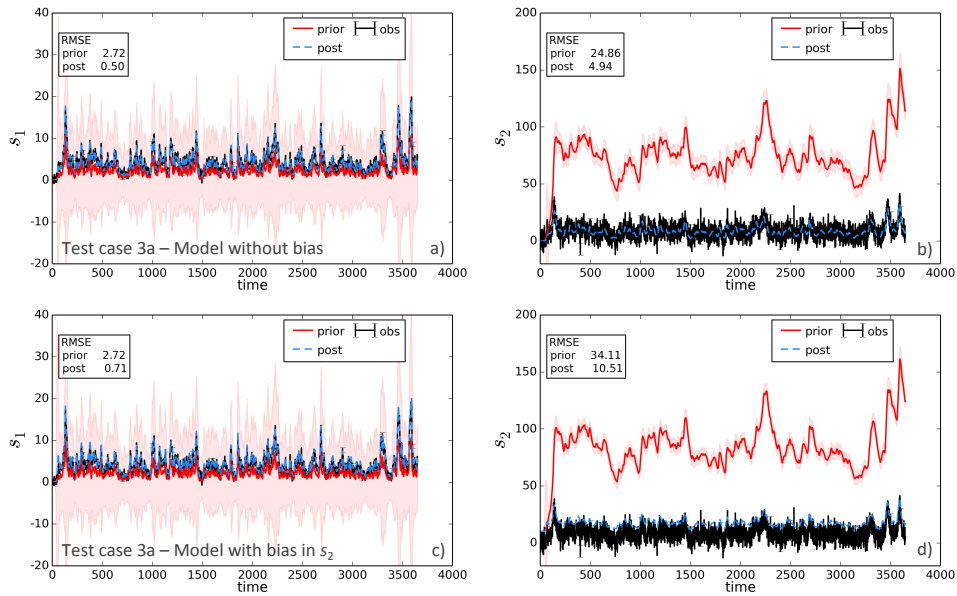
10

11

12

13

14



1

2 Figure 1: Prior and posterior model simulations compared to the synthetic observations for the
3 simple carbon model for test case 3a for a) s_1 and b) s_2 simulations without any model bias,
4 and c and d) with bias in the simulated s_2 variable. The coloured error band on the prior and
5 posterior represents the propagated parameter uncertainty (1 sigma) on the model state
6 variables (in the equivalent colour as the mean curve). This is mostly visible for the prior
7 model simulation (pink band) as there is a high reduction in model uncertainty reduction as a
8 result of the assimilation.

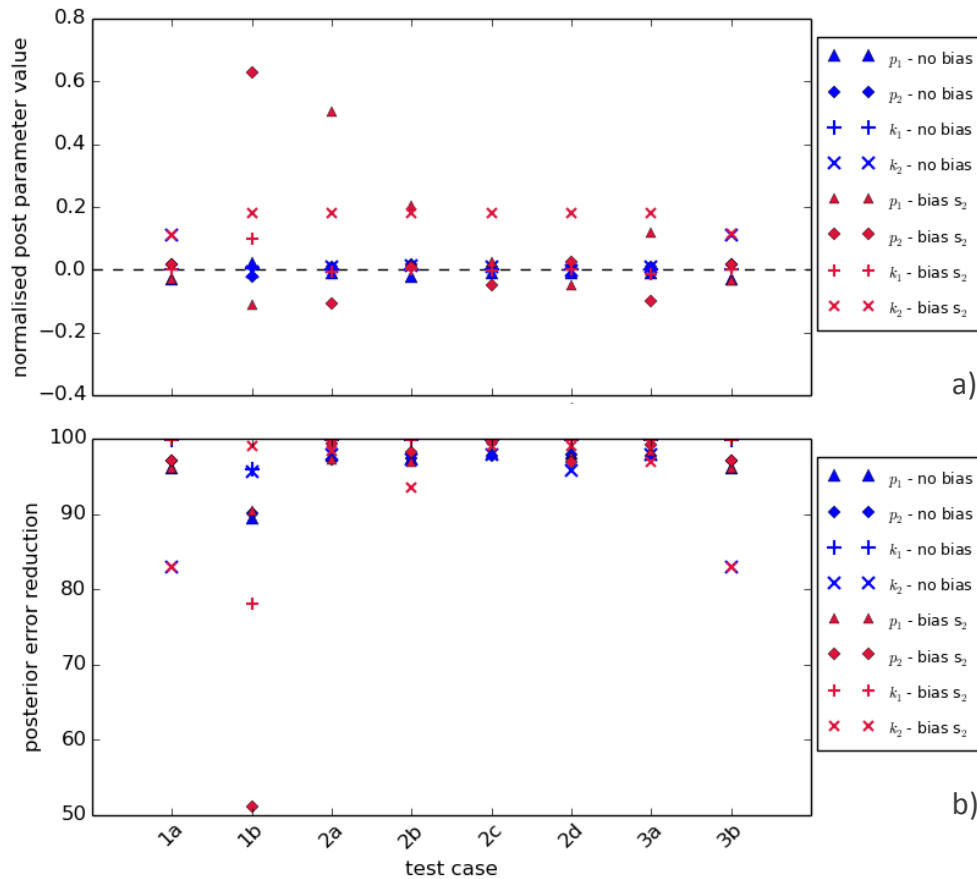
9

10

11

12

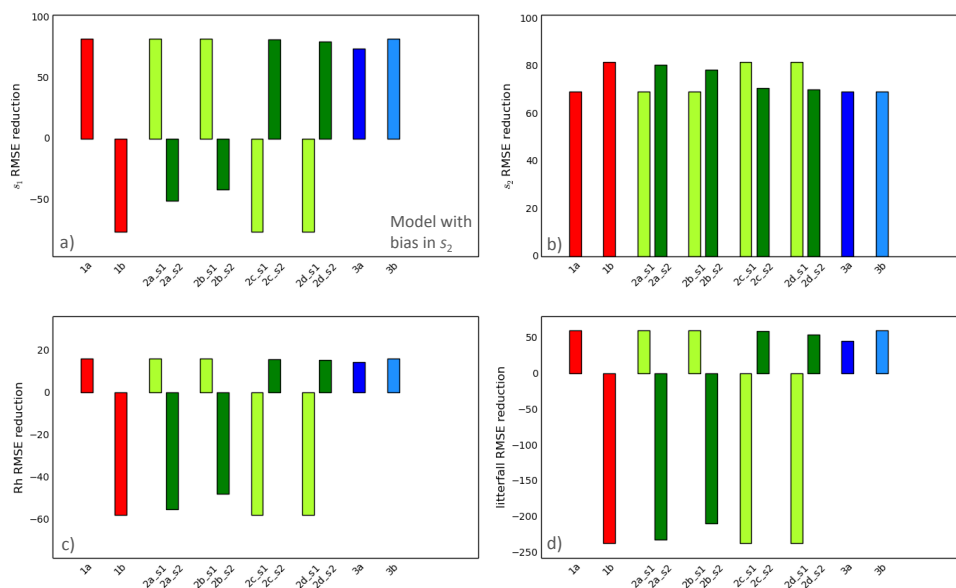
13



1

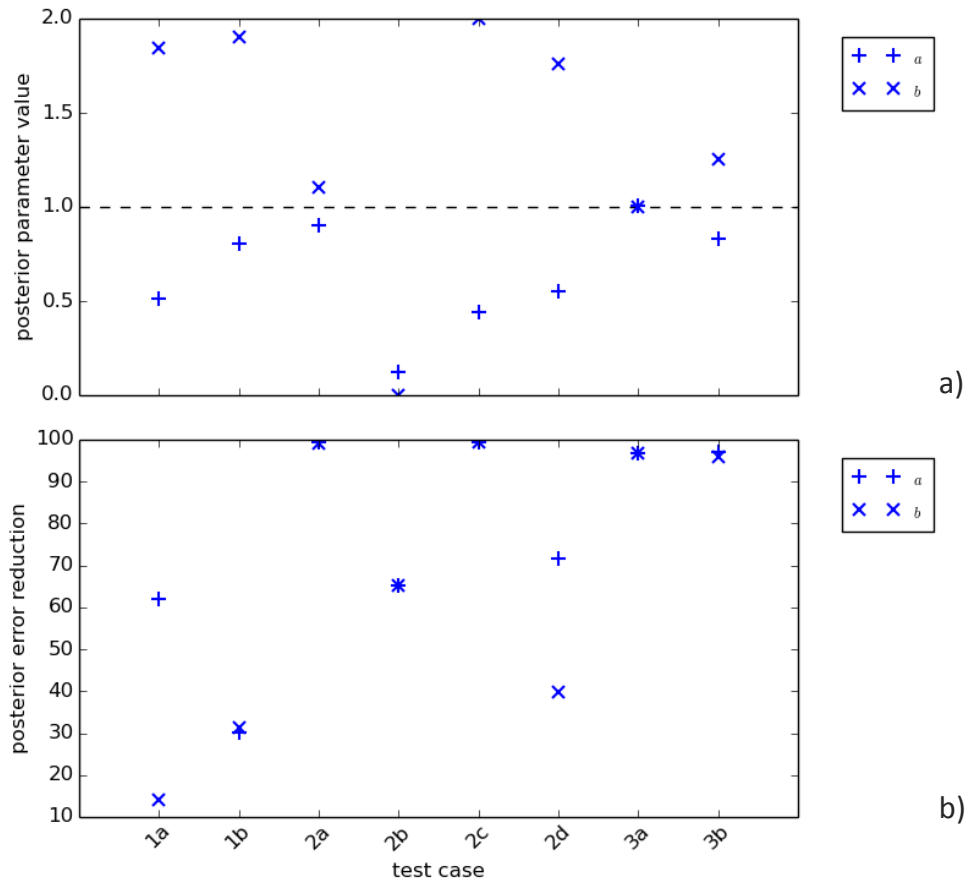
2 Figure 2: a) Normalised posterior parameter values and b) posterior parameter error reduction
 3 for all parameters of the simple carbon model for each test case, and for both the simulations
 4 with no bias (blue) and simulations with a bias in the s_2 variable that was not accounted for in
 5 the inversion (red). In a) parameters values were normalised to account for differences in the
 6 magnitude of the different parameters and their range, thus it is a measure of the distance
 7 from the true value as a fraction of the range and is calculated as: (posterior value – true value
 8 / max parameter value – minimum parameter value). The closer the value to the zero dashed
 9 line represents a better match to the “true” parameter value. To give an indication of the
 10 optimisation performance, the following are the normalised first guess parameter values for
 11 this particular example test (compare with posterior values in Fig. 2a): p_1 0.09, p_2 0.29, k_1 0.1,
 12 k_2 0.15.

13



1
 2 Figure 3: Reduction in RMSE for all test cases for simulations with a bias in the s_2 variable: a)
 3 s_1 , b) s_2 , c) litterfall and d) heterotrophic respiration (Rh). For the step-wise cases (2a, b, c and
 4 d) the reduction after both step 1 and step 2 are shown in light and dark green respectively,
 5 and are denoted in the x-axis labels with ‘_s1’ for step 1 and ‘_s2’ for step 2. The reduction
 6 (in %) is calculated as $1 - (\text{RMSE}_{\text{post}} / \text{RMSE}_{\text{prior}})$.

7
 8
 9
 10
 11
 12
 13
 14
 15
 16
 17



1

2 Figure 4: Posterior parameter values of both the non-linear toy model *a* and *b* parameters for
 3 each test case for the simulations with no model bias. The y-axis range corresponds to the
 4 parameter bounds and the dashed horizontal line represents the “true” known value of both
 5 parameters. To give an indication of the optimisation performance, the following are the first
 6 guess parameter values for this particular example test (compare with posterior values in Fig.
 7 4a): *a* 0.87, *b* 1.98. b) Posterior uncertainty reduction for both parameters for all test cases.

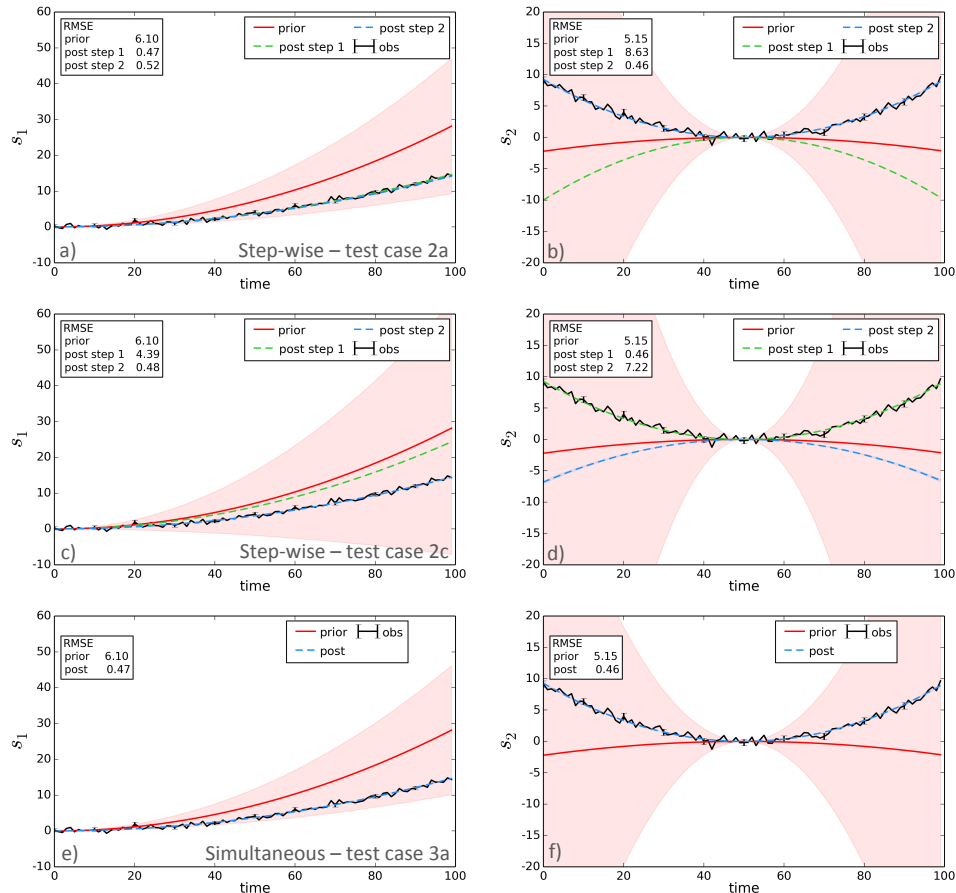
8

9

10

11

12



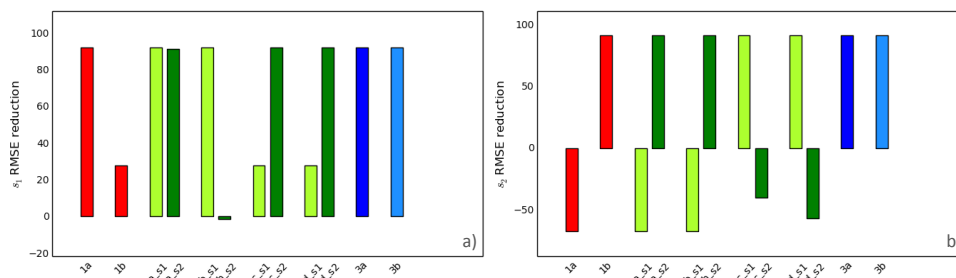
1

2 Figure 5: Prior and posterior model simulations compared to the synthetic observations for the
 3 non-linear toy model (with no bias) for both the s_1 (left column) and s_2 (right column)
 4 variables for a) and b) test case 2a (1st row) – step-wise approach with s_1 observations
 5 assimilated in the first step, followed by the s_2 observations in the second step; c) and d)
 6 test case 2c (2nd row) – step-wise approach with s_2 observations assimilated in the first step,
 7 followed by s_1 observations in the second step; and e) and f) test case 3a (3rd row) – the
 8 simultaneous case in which both data streams were included. For both step-wise examples A_1
 9 was propagated between the 1st and 2nd steps. The coloured error band on the prior and
 10 posterior represents the propagated parameter uncertainty (1 sigma) on the model state
 11 variables (in the equivalent colour as the mean curve). This is mostly visible for the prior



1 model simulation (pink band) as there is a high reduction in model uncertainty reduction as a
2 result of the assimilation.

3
4
5
6
7
8
9
10
11
12
13
14
15
16
17
18
19
20
21
22
23
24
25
26



1

2 Figure 6: Reduction in RMSE for all test cases for both a) s_1 and b) s_2 variables for the non-
 3 linear toy model simulations with no model bias. For the step-wise cases (2a, b, c and d) the
 4 reduction after both step 1 and step 2 are shown in light and dark green respectively, and are
 5 denoted in the x-axis labels with ‘_s1’ for step 1 and ‘_s2’ for step 2. The reduction (in %) is
 6 calculated as $1 - (\text{RMSE}_{\text{prior}} / \text{RMSE}_{\text{post}})$.

7

8

9

10

11