

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)

13 **Abstract**

14 Data assimilation methods provide a rigorous statistical framework for constraining
15 parametric uncertainty in land surface models (LSMs), which in turn helps to improve their
16 predictive capability and to identify areas in which the representation of physical processes is
17 inadequate. The increase in the number of available datasets in recent years allows us to
18 address different aspects of the model at a variety of spatial and temporal scales. However,
19 combining data streams in a DA system is not a trivial task. In this study we highlight some of
20 the challenges surrounding multiple data stream assimilation for the carbon cycle component
21 of LSMs. We give particular consideration to the assumptions associated with the type of
22 inversion algorithm that are typically used when optimising global LSMs – namely, Gaussian
23 error distributions and linearity in the model dynamics. We explore the effect of biases and
24 inconsistencies between the observations and the model (resulting in non-Gaussian error
25 distributions), and we examine the difference between a simultaneous assimilation (in which
26 all data streams are included in one optimisation) and a step-wise approach (in which each
27 data stream is assimilated sequentially) in the presence of non-linear model dynamics. In
28 addition, we perform a preliminary investigation into the impact of correlated errors between
29 two data streams for two cases, both when the correlated observation errors are included in

1 the prior observation error covariance matrix, and when the correlated errors are ignored. We
2 demonstrate these challenges by assimilating synthetic observations into two simple models:
3 the first a simplified version of the carbon cycle processes represented in many LSMs, and the
4 second a non-linear toy model. Finally, we provide some perspectives and advice to other
5 land surface modellers wishing to use multiple data streams to constrain their model
6 parameters.

7

8 Keywords: data assimilation, parameter optimisation, carbon cycle, biogeochemical cycles,
9 land surface model.

10

11 **1 Introduction**

12 The carbon cycle is an important component of the Earth system, especially when
13 considering the climatic impact of rising greenhouse gas concentrations from fossil fuel
14 emissions and land use change. It is estimated that the oceans and land surface absorb
15 approximately half of the CO₂ emissions due to anthropogenic activity, but uncertainties
16 remain in the strength and location of sources and sinks, as well as in predictions of future
17 trends (Ciais et al., 2013). Observations allow us to understand the system up until the present
18 day and provide inference about how ecosystems may respond to future change. However,
19 their use in estimating model state variables and boundary conditions is limited beyond
20 diagnostic purposes, and they can be restricted in their spatial coverage. They also do not
21 contain all the information we may need to distinguish between the complex interactions that
22 may occur between many different processes. Incorporating our current knowledge of
23 physical mechanisms of biogeochemical cycles, including carbon, C, dynamics, into land
24 surface models (LSMs) represents a promising approach for analysing these interacting
25 effects, upscaling observations to larger regions, and making future predictions. However, the
26 models can be limited by the lack of process representation, either due to gaps in our
27 knowledge, or in our technical and computing capability. As a result, model evaluations
28 reveal that not all variables are well-captured by the model under current conditions (Anav et
29 al., 2013), and the spread between model projections is still very large (Sitch et al., 2015).

30 Aside from model structural and forcing errors, one source of uncertainty is related to the
31 parameter (i.e. fixed) values of a model. Model-data fusion, or data assimilation (DA), allows

1 the calibration, or optimisation, of these values by minimising a cost function that quantifies
2 the model-data misfit, while accounting for the uncertainties inherent in both the model and
3 data in a statistically rigorous framework. The C cycle component of most LSMs is complex
4 and contains a large number of parameters; luckily however, there are an increasing number
5 of in-situ and remote sensing-based data streams that can be used for parameter optimisation.
6 These data bring information on different spatial and temporal scales, such as:

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

29

1 Researchers are increasingly attempting to bring these sources of data together to
2 constrain different parts of a model at different spatio-temporal scales within a multiple data
3 stream assimilation framework (e.g. Richardson et al., 2010; Keenan et al., 2012; Kaminski et
4 al., 2012; Forkel et al., 2014; Bacour et al., 2015). However, whilst the potential benefit of
5 adding in extra data streams to constrain the C cycle of LSMs is clear, multiple data stream
6 assimilation is not as simple as it may seem. This is particularly true when considering a
7 regional-to-global scale, multiple site optimisation of a complex LSM that contains many
8 parameters, and which typically takes on the order of minutes to an hour to run a one year
9 simulation. When using more than one data stream there is the option to include all data
10 streams together in the same optimisation (simultaneous approach), or to take a sequential
11 (step-wise) approach. Mathematically, the optimal approach is the simultaneous, but
12 complications may arise due computational constraints related to the inversion of large
13 matrices or the requirement of numerous simulations, particularly for global datasets (e.g.
14 Peylin et al., 2016), and/or due to the “weight” of different data streams in the optimisation
15 (e.g. Wutzler and Carvalhais, 2014). On the other hand, in a step-wise assimilation the
16 parameter error covariance matrix has to be propagated at each step, which implies that it can
17 be computed. If the parameter error covariance matrix can be properly estimated and is
18 propagated between each step, the step-wise approach should be mathematically equal to
19 simultaneous. However, many inversion algorithms (e.g. derivative-based methods that use
20 the gradient of the cost function to find its minimum) require assumptions of model (quasi-)
21 linearity and Gaussian parameter and observation error distributions (Tarantola, 1987, p195).
22 If these assumptions are violated, or the error distributions are poorly defined, it is likely that
23 the step-wise will not be equal to the simultaneous, because information will be lost at each
24 step due to an incorrect calculation of the posterior error covariance matrix at the end of each
25 step. An incorrect description of the observation (– model) error distribution could result from
26 i) the wrong assumption about the distribution of the residuals between the observation and
27 the model, ii) a poor characterisation of the error correlations, iii) an incompatibility between
28 the model and the data (possibly due to a model structural issue or differences in how a
29 variable is characterised), or iv) a bias in the observations that is not unaccounted for (i.e. is
30 treated as a random error). As mentioned, whilst a simultaneous optimisation is
31 mathematically more rigorous, in the sense that the error correlations are treated within the
32 same inversion, if the prior distributions are not properly characterised any bias may be

1 aliased to the wrong parameters (Wutzler and Carvalhais, 2014), and possibly more so than in
2 a step-wise approach.

3 This tutorial-style paper highlights some of the challenges of multiple data stream
4 optimisation of carbon cycle models discussed above. Note that we do not aim to explore all
5 possible issues related to a DA system, for example the choice of the cost function,
6 minimization algorithm, or the characterization of the prior error distributions; indeed,
7 previous studies have investigated such aspects at length (e.g. Fox et al., 2009; Trudinger et
8 al., 2007), and therefore we refer the reader to these papers for more information. Section 2
9 reviews recent carbon cycle multiple data stream assimilation studies with reference to some
10 of the aforementioned challenges. Section 3 demonstrates some these issues related to
11 multiple data stream assimilation with synthetic experiments using two simple models: one a
12 simplified version of the carbon dynamics included in many LSMs, and the other a “toy”
13 model designed to demonstrate the issues that arise with complex, non-linear models. Finally
14 Section 4 provides some advice to land surface modellers wishing to carry out multiple data
15 stream assimilation to constrain the parameters of their model.

16

17

18 **2 Review of existing multiple data stream carbon cycle data assimilation** 19 **studies**

20 **2.1 Extra constraint from multiple data streams**

21 Most site-based carbon cycle data assimilation studies have used eddy covariance
22 measurements of NEE and LE fluxes to constrain the relevant parameters of ecosystem
23 models. However, a few studies have also made use of chamber flux soil respiration data and
24 field measurements of vegetation characteristics (e.g. tree height, budburst, LAI) or estimates
25 of litterfall and carbon stocks as ancillary information (e.g. Fox et al., 2009; Keenan et al.,
26 2012; Thum et al., in review; Van Oijen et al., 2005; Richardson et al., 2010; Williams et al.,
27 2005). Two recent studies combined high-resolution satellite-derived FAPAR data with in-
28 situ eddy covariance measurements to optimize parameters related to carbon, water and
29 energy cycles of the ORCHIDEE and BETHY LSMs at a couple of sites (Bacour et al., 2015;
30 Kato et al., 2013, respectively).

1 At global scales the number of studies that use multiple data streams from satellites or
2 large-scale networks to optimise LSMs has been increasing in recent years, although this
3 remains a relatively new area of research. CCDAS-BETHY was the first global carbon cycle
4 data assimilation system (CCDAS) to make use of the high-precision measurements of the
5 atmospheric CO₂ concentration flask sampling network (Rayner et al., 2005; Scholze, 2003)
6 to constrain parameters of the terrestrial carbon cycle model BETHY (Knorr, 2000). Since its
7 first application using only atmospheric CO₂ concentration data, CCDAS-BETHY has been
8 further developed to consistently assimilate multiple data streams both at local and global
9 scales. In particular, Kaminski et al. (2012) optimised 70 process parameters, plus one initial
10 condition, by simultaneously assimilating a satellite-derived FAPAR product derived from the
11 Medium Resolution Imaging Spectrometer (MERIS; Gobron et al., 2008) and flask samples
12 of atmospheric CO₂ at two sites from the GLOBALVIEW product (GLOBALVIEW-CO₂,
13 2008) at coarse resolution. More recently, Scholze et al. (2016) demonstrated the added value
14 of assimilating remotely sensed soil moisture data in addition to atmospheric CO₂
15 concentration data. They used the same coarse resolution set-up of CCDAS-BETHY as
16 Kaminski et al. (2012) and CO₂ observations from 10 sites of the GLOBALVIEW product
17 (GLOBALVIEW-CO₂, 2012) together with the SMOS L3 daily soil moisture product
18 (version 246; CATDS-L3, 2012).

19 Three other global CCDAS based on LSMs that are part of earth system models (ESMs)
20 have been developed in recent years (Peylin et al., 2016; Raoult et al., 2016; Schürmann et al.,
21 2016). Two of these used multiple data streams as constraints. Schürmann et al. (2016)
22 optimized model parameters and initial conditions of the land component, JSBACH (Raddatz
23 et al. 2007), of the MPI ESM (Giorgetta et al. 2013) using atmospheric CO₂ concentration
24 data from 28 sites and the TIP-FAPAR product (Pinty et al., 2007) as joint constraints over a
25 5-year period. As part of their study they evaluated the mutual benefit of each data stream in a
26 fully factorial design. Peylin et al. (2016) used three different data streams as global
27 constraints for the ORCHIDEE LSM (Krinner et al., 2005), which forms the land surface
28 component of the IPSL ESM (Dufresne et al., 2013), in a multi-site, step-wise assimilation
29 approach. First, satellite-derived NDVI data from the MODIS instrument were used in a
30 similar manner to FAPAR as a proxy for vegetation greenness, in order to constrain the
31 phenology parameters at 60 sites for 4 temperate and boreal deciduous PFTs (MacBean et al.,
32 2015), followed by NEE and LE observations at 78 FLUXNET sites for 7 PFTs to optimise
33 all the carbon-related parameters (Kuppel et al., 2014), and finally atmospheric CO₂

1 concentration measurements from 53 sites in the GLOBALVIEW network (GLOBALVIEW-
2 CO₂, 2013), which predominantly provided a constraint on the initial magnitude of the soil
3 carbon reserves in the model. The three global multiple data stream CCDAS have allowed an
4 improvement in both the mean seasonal cycle as well as the trend of net land surface CO₂
5 exchange, especially with the inclusion of the atmospheric CO₂ data (Kaminski et al., 2012;
6 Peylin et al., 2016; Schürmann et al., 2016). Atmospheric CO₂ concentration observations are
7 one of the most accurate, long-term data sets in environmental science and they provide
8 important information about the global CO₂ sink capacity by the land and ocean.

9 Many of the aforementioned studies reported that adding extra data streams helped to
10 constrain unresolved sub-spaces of the total parameter space. Richardson et al. (2010) and
11 Keenan et al. (2012) concluded that using ancillary information (e.g. woody biomass
12 increment, field-based LAI and chamber measurements of soil respiration) as in addition to
13 NEE data provided a valuable extra constraint on many model parameters, which improved
14 both the bias in model predictions and reduced the associated uncertainties. The results of the
15 REFLEX model-data fusion inter-comparison project also indicated that observations of the
16 different carbon pools would help to constrain parameters such as root allocation and woody
17 turnover that were not well resolved using NEE and LAI data alone (Fox et al., 2009).
18 Similarly at global scale, Scholze et al. (2016) found that assimilating SMOS soil moisture
19 data in addition CO₂ observations reduced the ambiguity in the solution space compared to
20 assimilating CO₂ alone; about 30 parameters out of the 101 were resolved compared to 15
21 without SMOS data. Bacour et al. (2015) and Schürmann et al. (2016) both reported that the
22 addition of FAPAR data brought extra information on the phenology-related processes in the
23 model, and therefore retrieved different posterior C flux-related parameter values than when
24 assimilating NEE or atmospheric CO₂ data alone. An interesting aspect of the Kaminski et al.
25 (2012) study was that the inclusion of FAPAR in addition to atmospheric CO₂ concentration
26 samples resulted in a particular improvement for the hydrological fluxes in the model, thus
27 demonstrating the importance of assessing the potential benefit for model variables that may
28 not have been the main target of optimisation.

29 On the other hand, Williams et al. (2005) observed that one-off, or rarely taken,
30 measurements of carbon stocks were unable to constrain components of the carbon cycle to
31 which they were not directly related. This raises the issue of the relative influence of different
32 data streams in a joint assimilation, particularly if the number of observations for each is

1 vastly different, which will be the case when assimilating both half-hourly C flux data in
2 addition to C stock observations that are typically available at an annual time scale or greater.
3 The spatial distribution of each data stream is also important, especially for heterogeneous
4 landscapes (Barrett et al., 2005; Alton, 2013).

5 Although a number of multiple data stream assimilation studies exist at various scales,
6 very few studies have specifically investigated the added benefit of different combinations of
7 data streams in a factorial study, with a few notable exceptions (Barrett et al., 2005;
8 Richardson et al., 2010; Kato et al., 2013; Keenan et al., 2013; Bacour et al., 2015;
9 Schürmann et al., 2016). Kato et al. (2013) and Bacour et al. (2015) both evaluated the
10 complementarity of eddy covariance and FAPAR data streams at site level, i.e. the impact of
11 assimilating one individual data stream on the other model state variable, as well as when
12 both data streams were included in the optimization (see discussion in Section 2.2). The study
13 of Keenan et al. (2013) was particularly notable in its aim to quantify which data streams
14 provide the most information (in terms of model-data mismatch) and how many data streams
15 are actually needed to constrain the problem. They reported that of the 17 field-based data
16 streams available, projections of future carbon dynamics were well-constrained with only 5 of
17 the data sources, and crucially, not with eddy covariance NEE measurements alone. These
18 results may be specific to this site or type of ecosystem, but their study highlights the need for
19 further research in this area, and in particular, for synthetic data experiments that allow us to
20 understand which data will be the most useful for a given scientific question. This will also
21 enable researchers to plan more efficient measurement campaigns with experimentalists, as
22 also pointed out by Keenan et al. (2012).

23

24 **2.2 Issue of bias and inconsistencies between the observations and the** 25 **model**

26 Despite the theoretical benefit of adding data streams into an assimilation system as
27 additional constraints, several of the aforementioned studies at both site and global scale have
28 reported a bias or inconsistency either between the different observation data streams, or
29 between the observations and the model. This is easily detected when the optimisation of one
30 data stream results in a worse fit than the prior in one or more of the other data streams. Thum
31 et al. (in review) found that the addition of aboveground biomass stocks brought a longer-

1 term constraint on allocation parameters, but they noted an incompatibility when assimilating
2 both annual increment and total biomass data to optimise the longer timescale
3 mortality/turnover parameter. This was due to the fact the total stocks take into account losses
4 related to disturbance and management (e.g. canopy thinning) – processes that were not
5 included in that version of the model.

6 Kato et al. (2013) assimilated SeaWiFS FAPAR (Gobron et al., 2006) and eddy
7 covariance LE measurements at the FLUXNET site in Maun, Botswana. They showed that
8 the individual assimilation of each the two data streams resulted in a perfect (i.e. within the
9 observational uncertainty) fit to the assimilated data set, but a considerable degradation of the
10 fit to the non-assimilated data set compared to the prior. A comparison against eddy
11 covariance measurements of gross carbon uptake (gross primary production – GPP) pointed to
12 a bias in the FAPAR data because the fit to the independent GPP data was degraded after
13 assimilating FAPAR data only, while the fit improved after assimilating the LE data only.
14 Nevertheless, the simultaneous assimilation of both data streams achieved a compromise
15 between the two suboptimal results achieved after assimilating only one data stream. The
16 calibration further limited the number of parameters with correlated errors, and yielded a
17 higher theoretical reduction in parameter uncertainty and a decrease in the RMS difference by
18 16% for the GPP data compared to the prior.

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

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

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

1 one being that an inadequate model structure may be at fault – for example, LPJmL does not
2 include vegetation classes corresponding to shrub, moss and lichen species that are dominant
3 in these ecosystems. They also noted that the GPP product they used, which is based on a
4 model tree ensemble up-scaling of FLUXNET data (Jung et al., 2011), might contain
5 representation-related biases, given that there are very few FLUXNET stations in tundra
6 regions. The issue of representation errors of sites has been touched upon before (e.g.
7 Raupach et al., 2005). Alton (2013), who performed a global multi-site optimisation of the
8 JULES LSM with a diverse range of data including satellite-derived LAI, FLUXNET, soil
9 respiration and global river discharge, raised the point that FLUXNET sites are known to be
10 large carbon sinks, which could potentially result in biased global NEE estimates.

11 Resolving these apparent inconsistencies was beyond the scope of most of these
12 studies, aside from applying a bias correction where one was evident. Aside from simple
13 corrections, Quaife et al. (2008) and Zobitz et al. (2014) suggested that LSMs should be
14 coupled to radiative transfer models to provide a more realistic and mechanistic observation
15 operator between the quantities simulated by the model and the raw radiance measured by
16 satellite instruments. This proposition follows experience gained in the case of atmospheric
17 models for several decades (Morcrette, 1991).

18

19 **2.3 Step-wise versus simultaneous assimilation**

20 The paper by Alton (2013) documents the only previous study to have used a step-wise
21 assimilation approach with more than two data streams, and they found that the final
22 parameter values were independent of the order of data streams assimilated. No studies in the
23 LSM community to date have explicitly examined a step-wise versus simultaneous
24 assimilation framework with the same optimisation system and model. The step-wise
25 assimilation with the ORCHIDEE-CCDAS detailed in Peylin et al. (2016) has been compared
26 to a simultaneous optimisation using the same three data streams (as well as the same model
27 and inversion algorithm) as part of an on-going study. In the simultaneous optimisation, the
28 addition of NEE or atmospheric CO₂ concentration measurements resulted in a smaller
29 reduction in the fit to MODIS NDVI compared to the step-wise approach presented in Peylin
30 et al. (2016). As the NDVI data were normalised to the 95th percentile range this was not a
31 result of a simple bias in the magnitude of the data. Rather, it was likely due to

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

19

20

21 **3 Demonstration with two simple models and synthetic data**

22 The three sub-sections in Section 2 highlight examples within a carbon cycle modeling
23 context of the three main challenges faced when performing a multiple data stream
24 assimilation, namely, i) the possible negative influence of including additional data streams
25 on other model variables; ii) the impact of bias in the observations, missing model processes
26 or inconsistency between the observations and model (as discussed in Section 2.2), and iii) the
27 difference between a step-wise and simultaneous optimization (and the order of data stream
28 assimilation) if the assumptions of the inversion algorithm are violated, which is more likely
29 to be the case with non-linear models when using derivative-based algorithms and least-
30 squares formulation of the cost function (as discussed in Section 2.3). The latter point is
31 important because derivative methods (compared to global search) are the only viable option
32 for large-scale, complex LSMs given the time taken to run a simulation. In addition to the

1 above three challenges we have performed a preliminary investigation into the impact of
2 correlated errors between the data streams, which is a topic that has not yet been studied in the
3 context of carbon cycle models to our knowledge.

4 This section aims to demonstrate these challenges using simple toy models and synthetic
5 experiments where the true values of the parameters are known. Thus the following sections
6 include a description of the toy models together with the derivation of synthetic observations,
7 the inversion algorithm used to optimise the model parameters and the experiments
8 performed, followed by the results for each test case.

9 **3.1 Methods**

10 **3.1.1 Simple carbon model**

11 To demonstrate the challenges of multiple data stream assimilation in a carbon cycle
12 context, we have chosen a test model that represents a simplified version of the carbon cycle
13 dynamics typically implemented in most LSMs. The model has been well-documented in
14 Raupach (2007) and has been used previously in the OptIC DA inter-comparison project
15 (Trudinger et al., 2007). It is based on two equations that describe the temporal evolution (on
16 a daily time step) of two living biomass (carbon) stores, s_1 and s_2 , and the biomass fluxes
17 between these two stores:

$$18 \quad \frac{ds_1}{dt} = F(t) \left(\frac{s_1}{p_1 + s_1} \right) \left(\frac{s_2}{p_2 + s_2} \right) - k_1 s_1 + s_0 \quad (1)$$

$$19 \quad \frac{ds_2}{dt} = k_1 s_1 - k_2 s_2 \quad (2)$$

20 In this model formulation, s_1 and s_2 are approximately equivalent to above- and belowground
21 biomass stocks. The unknown parameters p_1 , p_2 , k_1 and k_2 will be optimised in the inversions.
22 The first term on the right-hand side of Eq. (1) corresponds to the Net Primary Production
23 (NPP) i.e. the carbon input to the system as a function of time, represented by $F(t)$, weighted
24 by factors (the two fractions in parentheses) that account for the size of both pools, in order to
25 introduce a limitation on NPP. The $F(t)$ forcing term is a random function of time (“log-
26 Markovian” random process) representing the effect of fluctuating light and water availability
27 due to climate on the NPP (Raupach, 2007 – Section 5.3). The litterfall is an output of s_1
28 (aboveground biomass store) and an input to s_2 (belowground biomass store) and is calculated
29 as a constant fraction (k_1) of s_1 (defined by $k_1 s_1$). Heterotrophic respiration (Rh) is a constant

1 fraction (k_2) of the belowground carbon reserve s_2 and is represented k_2s_2 . The constant s_0 is a
2 “seed production” term set to 0.01 (i.e. not optimised) to ensure the model does not verge
3 towards zero. A more detailed description of the properties of the model is given in Trudinger
4 et al. (2007 – Section 2.1) and an in-depth analysis of the dynamical behaviour of the model is
5 provided in Raupach (2007). Synthetic observations of both s_1 and s_2 variables were used to
6 optimise all the unknown parameters in the model (see Section 3.1.5).

7

8 3.1.2 Non-linear toy model

9 Although the simple carbon model contains a non-linear term it is essentially still a
10 quasi-linear model. In order to illustrate the challenges associated with multiple data stream
11 data assimilation for more complex non-linear models, especially when using derivative
12 methods, we defined a simple non-linear toy model based on two equations with two
13 unknown parameters:

$$14 \quad s_1 = a \exp^b + at^2 \quad (3)$$

$$15 \quad s_2 = \sin(10a + 10b) + 10t^2 \quad (4)$$

16 where s_1 and s_2 also correspond to two model state variables (as for the simple C model), a
17 and b are the unknown parameters included in the optimisation, and t is the independent
18 variable, which could represent time in a real-world scenario. Note that this model is not
19 based on any particular physical process associated with land surface biogeochemical cycles,
20 but it does contain typical mathematical functions that are observed in reality and
21 implemented in LSMs. For example, the sinusoidal function (Eq. (4)) could represent diurnal
22 variations of various processes such as photosynthesis and respiration. Exponential response
23 functions (such as in Eq. (3)) are also observed for certain processes, including the
24 temperature sensitivity of soil microbial decomposition. As for the simple carbon model,
25 synthetic observations corresponding to the s_1 and s_2 variables were used to optimise both
26 parameters (see Section 3.1.5).

27

28 3.1.3 Bayesian inversion algorithm

29 Most data assimilation approaches follow a Bayesian formalism which, simply put,
30 allows prior knowledge of a system (in this case the model parameters) to be updated, or

1 optimised, based on new information (from the observations). In order to achieve this we
 2 define a “cost function” that describes the misfit between the data and the model, taking into
 3 account their respective uncertainties, as well as the uncertainty on the prior information. If
 4 we follow a Bayesian formalism and least-squares minimisation approach, and assume
 5 Gaussian probability distributions for the model parameter and observation error
 6 variance/covariance, we derive the following cost-function (Tarantola, 1987):

$$7 \quad J(\mathbf{x}) = \frac{1}{2}[(H(\mathbf{x}) - \mathbf{y})^T \cdot \mathbf{R}^{-1} \cdot (H(\mathbf{x}) - \mathbf{y}) + (\mathbf{x} - \mathbf{x}^b)^T \cdot \mathbf{B}^{-1} (\mathbf{x} - \mathbf{x}^b)] \quad (5)$$

8 where \mathbf{y} is the observation vector, $H(\mathbf{x})$ the model outputs given parameter vector \mathbf{x} , \mathbf{R} the
 9 observation error covariance matrix (including measurement and model errors), \mathbf{x}^b the a priori
 10 parameter values, and \mathbf{B} the prior parameter error covariance matrix. This framework leads to
 11 a Gaussian posterior parameter probability distribution function and requires that the model
 12 and its observation operator are linear.

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

$$27 \quad \mathbf{x}_{i+1} = \mathbf{x}_i - \varepsilon_i (\mathbf{H}^T \mathbf{R}^{-1} \mathbf{H} + \mathbf{B}^{-1})^{-1} (\mathbf{H}^T \mathbf{R}^{-1} (\mathbf{y} - H(\mathbf{x})) + \mathbf{B}^{-1} (\mathbf{x}_i - \mathbf{x}^b)) \quad (6)$$

28 where i is the iteration number and \mathbf{H} is the Jacobian, or first-order derivatives, of H , which in
 29 this study is determined using a finite difference method. Note that as we are potentially
 30 dealing with non-linear models, the quasi-Newton method has been slightly adapted to

1 include the constant scaling factor ε_i (with a value <1.0) to ensure that the algorithm will
2 converge.

3 Of course no inversion algorithm is perfect, and therefore if the characterisation of the
4 error distribution is inaccurate, or when optimising strictly non-linear models, it is possible
5 that the true “global” minimum of the cost function has not been found. Derivative methods in
6 particular can get stuck in so-called “local minima”, preventing the algorithm from finding the
7 true minimum. To address this issue we carry out a number of assimilations with different
8 random first guess points in the parameter space. If they all result in the same reduction in
9 cost function value, we can have more confidence that the true minimum has been found.

10 Once the minimum of the cost function has been found, the posterior parameter error
11 covariance can be approximated (using the linearity assumption) from the inverse Hessian of
12 the cost function around its minimum, which is calculated using the Jacobian of the model at
13 the minimum of $J(\mathbf{x})$ (for the set of optimized parameters), \mathbf{H}_∞ , following Tarantola (1987):

$$14 \quad \mathbf{A} = [\mathbf{H}_\infty^T \mathbf{R}^{-1} \mathbf{H}_\infty + \mathbf{B}^{-1}]^{-1} \quad (7)$$

15 Note that the posterior error covariance matrix can be propagated into the model space to
16 determine the posterior uncertainty on the simulated state variables as a result of the
17 parametric uncertainty (as shown in the coloured error bands in the time series plots – Figures
18 1 and 5) using the following matrix product and the hypothesis of local linearity (Tarantola,
19 1987):

$$20 \quad \mathbf{R}_{post} = \mathbf{H}_\infty \cdot \mathbf{A} \cdot \mathbf{H}_\infty^T \quad (8)$$

21

22 3.1.4 Step-wise versus simultaneous assimilation

23 *Step-wise approach*

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

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

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

8 *Simultaneous approach*

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

12

13 3.1.5 Optimisation set-up: parameter values and uncertainty, and generation 14 of synthetic observations

15 In this study we used synthetic observations that were generated by running the model
16 with known (or ‘true’) parameter values and adding random Gaussian noise corresponding to
17 the defined observation error for both s_1 and s_2 (see Table 1). We optimised a ten-year time
18 window for the simple carbon model, in order to capture the dynamics of the s_1 and s_2 pools
19 over a time period compatible with typically available observations. For the non-linear toy
20 model, which did not correspond to physical processes in the terrestrial biosphere, we ran a
21 simulation over a window of 100 integrations (steps) of the equations. The observation
22 frequency was daily, corresponding to the time-step of the simple carbon model (a value of 1
23 for the non-linear toy model), and the observation error was set to 10% of the mean value for
24 each set of pseudo-observations derived from multiple first guesses of the model.

25 The true values of all parameters for both models are given in Table 1, together with
26 their upper and lower bounds (following Trudinger et al., 2007). We have not performed a
27 prior sensitivity analysis to decide to which parameters are important to include in the
28 optimisation, as the model variables are sensitive to all of the (small set of) parameters.
29 However, in the case of a more complex, large-scale LSM it is advisable to carry out such an
30 analysis, particularly given the computational burden of optimising many parameters. In this

1 study the parameter uncertainty (1 sigma) was set to 40% of the parameter range following
2 recent studies (e.g. Bacour et al., 2015). Prior values were chosen from a uniform random
3 distribution bounded by the parameter bounds.

4

5 3.1.6 Experiments

6 The specific objective of the following experiments was to test the impact of a bias in the
7 observations that is not accounted for in the R matrix, and the impact of using derivative
8 methods with non-linear models (as may be necessary with large-scale LSMs), particularly in
9 reference to the differences that may arise between step-wise and simultaneous optimisations.

10 Table 2 details the experiments that were carried out based on all possible combinations
11 for assimilating the two data streams. Three approaches were compared: i) separate – where
12 only one data stream was included in the optimisation; ii) step-wise – where each data stream
13 was assimilated sequentially; and ii) simultaneous – where both data streams were included in
14 the optimisation. All parameters for both models were optimised in all experiments, therefore
15 in the step-wise cases the parameters were optimised twice. The step-wise assimilations were
16 also carried out with and without the propagation of the full posterior parameter error
17 covariance matrix, \mathbf{A}_1 , in between steps 1 and 2 (test cases 2b and d – see Table 2); only the
18 For the tests in which the full posterior covariance matrix was not propagated only the
19 posterior variance was propagated. An additional test was included for the simultaneous
20 assimilation in order to test the impact of having a substantial difference in the number of
21 observations for the data stream included in the optimisation, as may be the case for
22 belowground (e.g. soil) biomass observations in reality. Therefore in test case 3b, only one
23 observation was included for data stream s_2 .

24 The differences in the parameter values and the theoretical reduction in their uncertainty
25 ($1 - (\sigma_{\text{post}} / \sigma_{\text{prior}})$) were examined for all eight test cases, as well as the fit (RMSE) to both
26 data streams after the optimisation. For the step-wise approach we investigated if the fit to the
27 first data stream is degraded in the second step by comparing the RMSE after each step. Note
28 that the reduction in uncertainty is a theoretical, or approximate, estimate of the real
29 uncertainty reduction because of the assumptions made in the inversion scheme.

30 In a second stage the impact of an unknown, un-accounted for bias in the model was
31 examined. This bias could be a systematic bias in the observations due to the algorithm used

1 for their derivation, the result of missing or incomplete processes in the model, or an
2 incompatibility between the observations and the model, for example due to differences in
3 spatial resolution or an inconsistent characterisation of a variable between the model and the
4 observations. To test the impact of such an occurrence, we introduced a constant scalar bias
5 into the modelled s_2 variable with a value of 10 (i.e. twice the magnitude of the defined
6 observation uncertainty). All eight experiments were repeated, but a bias was introduced into
7 the model calculation of s_2 that was not accounted for in the cost function (i.e. the error
8 distributions retained a mean of zero). This was treated as an unknown bias, and therefore not
9 corrected or accounted for in the inversion scheme and the defined observation uncertainty
10 (Table 1) was not changed for this set of experiments.

11 In all experiments for both models, we used fifteen iterations of the inversion algorithm,
12 and twenty assimilations were performed starting from different random “first guess” points
13 in the parameter space. As discussed in Section 3.1.3, this was done to test the ability of the
14 algorithm to converge to the global minimum of the cost function. Note that the global
15 minimum and possible reduction in $J(x)$ will be different for each experiment, as each is based
16 on a different cost function.

17 For all the above tests we assumed independence (i.e. uncorrelated errors) for both the
18 observation and parameter prior error covariance matrices, thus the \mathbf{R} and \mathbf{B} matrices were
19 diagonal. In a final test we performed a simultaneous optimisation to examine the impact of
20 having correlated errors between the s_1 and s_2 observations. Thus the random Gaussian noise
21 added to s_1 for each time step was correlated to the noise added to s_2 . The correlated
22 observation errors were generated following the method used by Trudinger et al. (2007 –
23 paragraph 22). We first defined the covariance matrix, \mathbf{R} , using the prescribed observation
24 error and correlation between s_1 and s_2 . The correlated error that is added to the synthetic
25 observations is then a multiplication of a vector of Gaussian random noise (variance of 1) by a
26 matrix, \mathbf{X} , that corresponds to the Cholesky decomposition of \mathbf{R} (so that $\mathbf{R} = \mathbf{X}^T \mathbf{X}$). The added
27 noise was time invariant, i.e. there was no correlation between one time step and the next, as
28 we were specifically looking at correlations between the two data streams (see Pinnington et
29 al. (2016) for an analysis of the impact of correlations in the matrix and temporal error
30 correlation in the observations). We tested both accounting for the correlated errors by
31 populating the corresponding off-diagonal elements of the \mathbf{R} (observation error covariance)
32 matrix, and ignoring the correlated errors by keeping \mathbf{R} diagonal. The reason for performing
33 both tests was to demonstrate the possible real world scenario where correlated observation

1 errors exist but that information is not included in the optimisation, likely due to a lack of
2 knowledge as to how to characterise the errors. For both tests we performed optimisations
3 using a combination of different of observation error and correlation magnitudes
4 (observations errors between 0.05 and 20 in 9 uneven intervals, and observation correlations
5 between -0.9 and 0.9 with an interval of 0.4), in order to test the hypothesis that observations
6 with lower uncertainty (therefore higher information content) were less affected by the
7 presence of error correlations. As in the above experiments, 20 random first guesses in the
8 parameter space were used and 15 iterations of the inversion algorithm were performed.

9

10 **3.2 Results**

11 The 20 random first guess assimilations were examined for each set of experiments for
12 both models (before the results for each test were examined in more detail), in order to check
13 that the algorithm converged to a global minimum. As shown in the supplementary
14 information (Fig. S1), a high proportion of the 20 first guess assimilations across all test cases
15 for both models resulted in a similar reduction in $J(x)$, even though the overall magnitude of
16 the reduction was sometimes different between tests. This indicates that the algorithm does
17 not easily get stuck in any local minima (if they exist). The examples shown in the results
18 below were taken from one first guess parameter set for each model that belonged to the
19 cluster that had the highest cost function reduction. Any differences seen in the parameter
20 values, their posterior uncertainty or the resultant RMSE reduction described below therefore
21 are due to the specific details of each test and not the inability of the algorithm to find the
22 minimum.

23

24 **3.2.1 Typical performance with a quasi-linear model and no bias**

25 Figures 1a and b show the simple carbon model simulations for test case 3a (in which
26 both data streams are assimilated simultaneously) for the s_1 and s_2 variables. A large reduction
27 in RMSE is achieved after optimisation (blue curve) with respect to the observations (black
28 curve). Overall, there is a good reduction in RMSE for all test cases (including the individual
29 assimilations 1a and 1b) with a reduction of ~80% for s_1 and s_2 . In addition, the optimisation
30 of the s_1 and s_2 variables resulted in a good or moderate reduction in RMSE for variables not
31 included in any assimilation: ~60% for the litterfall (Eqn. 1) and ~16% for the heterotrophic

1 respiration (R_h – Eqn. 2) across all test cases (not shown), although there was already a good
2 prior fit to the data. As would be expected from these results, the parameter values and the
3 theoretical reduction in parameter uncertainty do not vary between the tests (Figures 2a and b
4 blue symbols), except for a slight difference in the value of the k_2 parameter in test cases 1a
5 and 3b, for which there is also a lower reduction in uncertainty (~82% compared to >95%).
6 Note that Fig. 2a shows the normalised parameter values to account for differences in the
7 magnitude of the different parameters and their range (the zero line represents the “true”
8 parameter value – see caption). In this situation therefore, where we have a relatively simple
9 linear model and two data streams to which the model parameters are highly sensitive, we see
10 that the differences between the step-wise and simultaneous approaches are minimal. This is
11 even the case when the error covariance is not propagated between the two steps (test cases 2b
12 and d), suggesting that under this assimilation set-up with this model both s_1 and s_2
13 individually contain enough spatio-temporal information to retrieve the true values of all
14 parameters, as we can see from the separate test cases 1a and b. However, we cannot
15 definitively say whether this is due to the simplicity or relative linearity of the model – it is
16 possible that observations of variables in more complex linear model would not be able to
17 retrieve the true values of all parameters.

18

19 3.2.2 Impact of unknown bias in one data stream – example with a simple 20 carbon model

21 In Section 3.2.1 we saw that there is little difference between a step-wise and
22 simultaneous optimisation if there is no bias in the model or observations, and if the model is
23 quasi-linear and therefore the critical assumptions behind the inversion approach were not
24 violated. However, it is not uncommon to have a bias between your observations and model
25 that is not obvious, and therefore not accounted for in the optimisation, as the cost function
26 used in most inversion algorithms (and in this study) assume Gaussian error distributions with
27 a mean of zero. Note that this is also the case when defining a likelihood function for
28 accepting or rejecting parameter values in a global search method. To test the impact of a
29 bias, we added a constant value to the simulated s_2 variable in a second test (see Section 3.1.6)
30 that was treated as an unknown bias, and therefore not corrected or accounted for in the
31 inversion scheme. The impact of this bias on s_1 and s_2 is shown in Figures 1c-d, and the
32 reduction in RMSE between the model and observations is seen in Fig. 3 for all variables

1 (including Rh and litterfall). The red symbols in Fig. 2 show the resultant parameter values
2 and theoretical reduction in uncertainty as a result of the bias. The inversion cannot accurately
3 find the correct values for all parameters in any test case and there are now considerable
4 differences between the simultaneous and step-wise approach. Furthermore the order in which
5 the data streams are assimilated in the step-wise cases also results in different posterior
6 parameter values (test cases 2a and b versus 2c and d in Fig. 2a and Fig. 3). Nevertheless the
7 optimisation results in a similar reduction in uncertainty on the parameters, except in test case
8 1b where only s_2 data are assimilated (Fig. 2b).

9 The main impact of the bias in the modelled s_2 variable is on the value of k_2 parameter
10 (Fig. 2a), which is consistently offset from the true value (dashed line in Fig. 2a) in all test
11 cases. This was expected given that it is the parameter most directly related to the calculation
12 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
13 test case 2b). Note that these parameters only indirectly influence the s_2 pool in the model,
14 and therefore we might have expected that they would be less affected by the bias. This nicely
15 demonstrates one issue that could arise in all DA studies, where the bias in a particular
16 variable (in the observations or the model) is aliased onto another process in the model
17 (Wutzler and Carvalhais, 2014). Such an aliasing of bias onto indirectly related parameters is
18 even more evident when only s_2 is included in the assimilation and s_1 does not provide any
19 constraint (test case 1b) – in this case all parameters are incorrect but the p_2 parameter in
20 particular shows a strong deviation from the true value (Fig. 2a). As a result we see a
21 deterioration in the RMSE for the s_1 , litterfall and Rh variables in test case 1b and in the step-
22 wise cases where s_2 is assimilated in the second step (Figures 3a, c and d – test case 1b, 2a
23 and 2b). However, the RMSE reduction remains high for the s_2 variable for these test cases
24 (Fig. 3b), as the inversion has found a solution that accounts for the bias even though all
25 inferred parameter values are incorrect. The assimilation of s_1 in the second step lowers the
26 reduction in RMSE for s_2 gained in the first step to $\sim 70\%$, but it is not a considerable
27 degradation.

28 Even though the posterior parameter values are incorrect, and despite the fact that the
29 first step results in a degradation, the final reductions in RMSE are largely the same as the
30 situation with no bias for all variables when s_1 is included in a simultaneous assimilation or
31 optimised in the second step (test cases 2c, d and 3a in Fig. 3). This shows that the inclusion
32 of s_1 observations can find a solution to counter the bias in s_2 and prevents a degradation in

1 the fit to the data. If s_2 is assimilated in the second step there is a negative impact on all other
2 variables, as discussed above, demonstrating again that the order of data stream assimilation
3 can matter when biases or inconsistencies between the data and the model are present.

4 The analysis of the impact of the bias presented here is specific to this model and the
5 type and magnitude of the bias that was added, but the broader findings can be generalised to
6 any situation in which there is a bias or inconsistency between a model and data that is not
7 accounted for in the assigned error distributions. Exactly what might constitute a bias or
8 inconsistency is discussed more in Section 2.2. Note also that it is important to examine the
9 impact on the other variables. For the separate test case 1b in which only s_2 data are used to
10 optimise the model, the negative impact on the other variables (Fig. 3) would have been
11 concealed if we had only examined the posterior reduction in RMSE for the s_2 variable.
12 Again, this is a concern that is inherent to all DA experiments, whether single- or multiple-
13 data stream, but we can see from these results (i.e. by comparing the separate test cases 1b
14 with 2a and b) that adding another data stream in a multiple-constraint approach does not
15 always reduce the problem.

17 3.2.3 Difference between the step-wise and simultaneous approaches in the 18 presence of a non-linear model

19 As discussed in Section 3.2.1, there is little difference between the step-wise and the
20 simultaneous assimilation approaches for simple, relatively linear models, unless the
21 observation error (including measurement and model errors) distribution deviates strongly
22 from the Gaussian assumption. However in reality, large-scale, complex LSMs may contain
23 highly non-linear responses to certain model parameters. To demonstrate the impact of non-
24 linearity in a multiple data stream assimilation context we used a non-physically based toy
25 model chosen for its non-linear characteristics (see Section 3.1.2).

26 Fig. 4a shows the posterior parameter values for both the a and b parameters of the
27 non-linear toy model for all test cases. The values were not normalised as both parameters had
28 the same range. The horizontal dashed line shows the “true” known values of the parameters
29 (both equal to 1.0) that were used to generate the synthetic observations. Note that no bias has
30 been introduced into the model in the results described here. The prior and posterior model s_1
31 and s_2 simulations for the non-linear toy model are compared to the synthetic observations in

1 Fig. 5 for both step-wise cases in which the posterior error covariance matrix from step 1 (A_1
2 – see section 3.1.4) was propagated to step 2 (experiments 2a and c – Fig. 5a-d) and both
3 simultaneous cases 3a and b (Fig. 5 e-h). Finally Fig. 6 summarises the reduction in RMSE
4 between the simulated and observed s_1 and s_2 variables for the non-linear toy model for all
5 test cases and, in the step-wise cases, the reduction in RMSE after both the first and second
6 steps (light versus dark green bars).

7 Assimilating each data stream individually (test cases 1a and b) does not result in an
8 accurate retrieval of the posterior parameters (Fig. 4a), nor in a strong constraint on either
9 parameter, as shown by the lack of theoretical reduction in the parameter uncertainty after the
10 optimisation (Fig. 4b). Despite this, there is a ~90% reduction in RMSE for the data stream
11 that was included in the optimisation (i.e. for s_1 in test case 1a – Fig. 6a, and s_2 in test case 1b
12 – Fig. 6b). However, the improvement on the other data stream is much less (28% reduction
13 in RMSE for s_1 when s_2 is assimilated) or even results in a degradation compared to the prior
14 fit (e.g. in the case of s_2 when s_1 is assimilated – Fig. 6b). Lack of improvement, or even
15 degradation, in the RMSE of other variables in the model is a common issue for data
16 assimilation in general, one that is not often evaluated in model-data fusion studies. It is also
17 is not necessarily the result of a bias or incompatibility between the observations and the
18 model.

19 Only the simultaneous case, in which all s_1 observations have been included in the cost
20 function (test case 3a), manages to retrieve the correct parameter values after the optimisation.
21 The posterior parameter values for all other test cases are incorrect, and are considerably
22 different between each case, unlike for the simple carbon model (without a model bias). Most
23 step-wise test cases (particularly 2b-d) do not result in the same parameter values as the
24 simultaneous test case 3a in which all the observations are included (Fig. 4a). This highlights
25 that strong non-linearity in the model sensitivity to parameters, together with the use of an
26 algorithm that is only adapted to weakly non-linear problems, can result in differences
27 between a step-wise and simultaneous approach in multiple – data stream assimilation (see
28 Section 1).

29 In the simultaneous optimisation in which all observations are included (test case 3a),
30 the posterior fit to the data dramatically improves for both the s_1 and s_2 data streams after the
31 assimilation (blue dashed line in Fig. 5e and f). This was expected given that the correct
32 values of the parameters were found. For the step-wise cases (test case 2a in Figures 5a and b,

1 and test case 2c in Fig. 5c and d), the black dashed line shows the prior, and the posterior after
2 step 1 is shown by green dashed line. In the step-wise assimilation we see two different
3 scenarios depending on which data stream was assimilated first. In the first step the results are
4 the same as the case where each individual data stream is assimilated separately. In both cases
5 the first step results in a good fit to the data that was included in the optimisation in that step.
6 When the s_1 data was assimilated in the first step (Fig. 5 first row), the fit to s_2 deteriorated
7 after the optimisation (Fig. 5b green dashed line and Fig. 6b – test case 2a_s1), but when the
8 s_2 data were assimilated first (Fig. 5 second row) the optimisation step did manage to achieve
9 an improvement in the s_1 data stream (Fig. 5c green dashed line and Fig. 6a – test case 2c_s1).

10 In the second step the optimisation of s_2 in test cases 2a and b does not degrade the fit
11 to s_1 when the full parameter error covariance matrix (\mathbf{A}_1) is propagated between step 1 and 2
12 (Figures 5a blue curve and 6a 2a_s2). Furthermore optimising s_2 in the second step reverses
13 the deterioration in s_2 caused by assimilating s_1 in the first step (Figures 5b blue curve and 6b
14 2a and b dark green bars). However, when s_1 data were assimilated in the second step (test
15 cases 2c and d), we found that the good fit achieved with s_2 observations in the first step was
16 effectively reversed (Fig. 5d blue curve). Therefore assimilating s_1 in the second step
17 degraded the fit to the s_2 observations, even compared to the prior case (Fig. 6b, dark green
18 bars for test cases 2c and d). This nicely highlights one of the main possible issues with a
19 step-wise assimilation framework.

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

25

26 3.2.4 Order of assimilation of data streams and propagation of parameter 27 error covariance matrices in a step-wise approach

28 Comparing the step-wise cases 2a and b with 2c and d for the non-linear toy model
29 reveals that neither order in the assimilation, s_1 then s_2 , or s_2 then s_1 , results in the correct
30 posterior parameter values that match the simultaneous test case (Fig. 4a). This is not a result
31 that can be generalised to all step-wise assimilations as it will depend on the data stream

1 involved and whether they contain enough spatio-temporal information to accurately
2 constrain all the parameters included in the optimisation, as well as any biases in the model or
3 observations (as discussed in Section 3.2.2) or model non-linearity (section 3.2.3). In the case
4 of the non-linear toy model, neither s_1 nor s_2 find the right parameter values when assimilated
5 individually, therefore it is not surprising that neither order manages to achieve the right
6 posterior parameter values. Nevertheless, the theoretical uncertainty of both parameters is
7 reduced by >95% for the step-wise cases in which \mathbf{A}_1 from step 1 is propagated between step
8 1 and 2 (test cases 2a and c – Fig. 4b), even though the posterior values for the step-wise
9 cases are incorrect. This demonstrates that a good theoretical reduction in uncertainty is not
10 always indicative that the right parameters have been found by the optimisation. The lower
11 theoretical reduction in parametric uncertainty for cases 2b and d (Fig. 4b) demonstrates that
12 information is lost between the steps if the posterior error covariance terms of \mathbf{A}_1 are not
13 propagated to step 2.

14 From a mathematical standpoint the most rigorous approach is to propagate the full
15 parameter error covariance matrices between each step. Without that constraint not only is
16 information lost in the second step, but the information contained in the second data stream
17 may have a stronger influence compared to a simultaneous assimilation, or step-wise case
18 with a propagated error covariance matrix. The inversion may therefore be more vulnerable to
19 any strong biases or incompatibilities between the model and the observations of the second
20 data stream, or indeed the particular sensitivity of its corresponding model state variable to
21 the parameters. This is one possible explanation for the degradation seen in s_1 in the non-
22 linear toy model when s_2 is optimised in the second step and \mathbf{A}_1 is not propagated between the
23 steps (Fig. 6a test case 2b_s2). The same was also true for the simple carbon model for test
24 case 2b when a bias was introduced into the s_2 simulation (see Section 3.2.2 and Fig. 3a).

25 However, the reverse is also true – if the first data stream contains strong biases then
26 the associated error correlations will be also propagated with \mathbf{A}_1 . If autocorrelation in the
27 observation errors, or indeed correlation between the errors of the data streams, is not
28 accounted for, it is likely that the posterior simulations are over-tuned, i.e. we will
29 overestimate the reduction in parameter uncertainty. If this is the case and the first step results
30 in incorrect parameter values, the propagation of \mathbf{A}_1 could restrict the parameter values to the
31 wrong location in the parameter space, and thus inhibit the ability of the inversion to find the
32 correct global minimum. These issues are likely to be more considerable for non-linear

1 models, as seen by the lack of difference between test cases 2a-d in the simple carbon model
2 example (Fig. 2).

3

4 3.2.5 Impact of accounting for correlated observation errors in the prior 5 observation error covariance matrix

6 In a final test we introduced time invariant correlated noise between the two data streams
7 (see Section 3.1.6). We investigated the impact of ignoring cross-correlation between two
8 data streams by comparing the results of i) an optimisation in which the correlated errors were
9 included in the off-diagonal elements of the prior observation error covariance matrix, \mathbf{R} , to
10 ii) an optimisation in which the correlated observation errors were excluded (i.e. \mathbf{R} was kept
11 diagonal). Note that this experiment is only relevant to simultaneous multiple data stream
12 assimilation, as it is not possible to account for cross-correlation between data streams when
13 one is assimilated after the other in a step-wise approach.

14 The presence of correlated errors increases observation redundancy in the inversion,
15 which would therefore reduce the expected theoretical error reduction compared to
16 uncorrelated observations (experiments not shown). We would expect a further limitation on
17 the expected error reduction with a sub-optimal system, as represented by optimisation ii) in
18 which there was cross-correlation between the data streams, but the correlated observation
19 errors were ignored in the \mathbf{R} matrix (as seen in Chevallier, 2007).

20 Figure 7 shows the difference between the two optimisations, (i.e. including off-diagonal
21 elements in the \mathbf{R} matrix minus only diagonal elements in the \mathbf{R} matrix), for the reduction in
22 the cost function value (Figures 7a and d) and posterior s_1 and s_2 observation errors (1 sigma –
23 Figures 7b, c, e and f), for both the simple C model (top row) and the non-linear toy model
24 (bottom row) and for a range of observation error and correlation. The plot shows the median
25 difference across all twenty random first guess parameters, and the reduction is calculated as
26 $1 - (\text{posterior}/\text{prior})$.

27 At low observation error there is no discernible difference between accounting for the
28 correlated observation errors in the \mathbf{R} matrix or not. This is likely because there is enough
29 information in the observations to find the global minimum of the cost function. Trudinger et
30 al. (2007) also found that similar posterior values were obtained when comparing
31 observations with correlated and uncorrelated Gaussian errors. However, at a certain point as

1 observation error increases along the x-axis (i.e. decreasing information content) there is a
2 difference in the cost function and observation error reduction between the two optimisations
3 for both models (Figure 7). As expected, the optimal optimisation that includes off-diagonal
4 correlated errors in R results in a higher reduction (blue cells in Figure 7) in the cost function
5 and posterior observation error than the sub-optimal optimisation (in which the correlated
6 errors are ignored) in all cases except for the s_1 data stream in the simple C model (see
7 below). Furthermore we see a pattern emerging suggesting that the difference between the
8 two optimisations increases with higher observation correlation for the same error magnitude.
9 However, for some combinations of observation error and correlation, the pattern is opposite
10 to what we expect (red cells in Figure 7), particularly for the s_1 data stream in the simple C
11 model (Figure 7b). This is likely because the accuracy of the solution becomes limited by
12 observation uncertainty at higher observation errors, and also due to presence of model non-
13 linearity, which prevents a fully accurate characterisation of the posterior error covariance
14 matrix with the inversion algorithm we have used.

15 The key finding of this preliminary investigation into the impact of correlated
16 observation errors is that it becomes increasingly important to properly characterise and
17 account for correlations between data streams if the observations do not contain enough
18 information (i.e. high observation uncertainty or a limited number of observations). However,
19 this is a wide topic that has received little-to-no attention in the carbon cycle data assimilation
20 literature to date, aside from the 2 out of 21 experiments in the wider-ranging study of
21 Trudinger et al. (2007). We therefore suggest that an investigation such as this should be
22 extended in order to fully understand the impact of cross-correlation between data streams;
23 however, this is beyond the scope of this paper.

24

25

26 **4 Perspectives and advice for Land Surface Modellers**

27 Although it is clear that in many cases the addition of different observations in a model
28 optimisation provides additional constraints, challenges remain that need to be addressed.
29 Many of the issues that we have investigated are relevant to any data assimilation study,
30 including those only using one data stream. However, most are more pertinent when
31 considering more than one source of data. Based on the simple toy model results presented

1 here, in addition to lessons learned from existing studies, we recommend the following points
2 when carrying out multiple data stream carbon cycle data assimilation experiments:

- 3 • If technical constraints require that a step-wise approach be used, it is preferable
4 (from a mathematical standpoint) to propagate the full parameter error covariance
5 matrix between each step. Furthermore, it is important to check that the order of
6 assimilation of observations does not affect the final posterior parameter values,
7 and that the fit to the observations included in the previous steps is not degraded
8 after the final step (e.g. Peylin et al., 2016).
- 9 • Devote time to carefully characterising the parameter and observation error
10 covariance matrices, including their correlations (Raupach et al., 2005), although
11 we appreciate this is not an easy task (but see Kuppel et al., 2013 for practical
12 solutions). In the context of multiple data stream assimilation, accounting for the
13 error correlations between data streams is increasingly important with higher
14 observational uncertainty (or a limited number of observations), though note that
15 this is not possible in a step-wise assimilation.
- 16 • The presence of a bias in a data stream, or an incompatibility between the
17 observations and the model, will limit the utility of using multiple observation
18 types in an assimilation framework. Therefore it is imperative to analyse and
19 correct for biases in the observations and to determine if there is an incompatibility
20 or inconsistency between the model and data. Alternatively, it may be possible
21 account for any possible bias/inconsistency in the observation error covariance
22 matrix, \mathbf{R} , using the off-diagonal terms or inflated errors (see Chevallier, 2007), or
23 by using the prior model-data RMSE to define the observation uncertainty.
- 24 • Most optimisation studies with a large-scale LSM require the use of derivative-
25 based algorithms based on a least-squares formulation of the cost function, and
26 therefore rely on assumptions of Gaussian error distributions and quasi-model
27 linearity. However, if these assumptions are not met it may not be possible to
28 find the true global minimum of the cost function and the resultant calculation of
29 the posterior probability distribution will be incorrect. This is a particular problem
30 if the posterior parameter error covariance matrix is propagated multiple times in a
31 step-wise approach, although these issues are relevant to both step-wise and
32 simultaneous assimilation. Therefore it is important to assess the non-linearity of

1 your model, and if the model is strongly non-linear, use global search algorithms
2 for the optimisation – although at the resolution of typical LSM simulations
3 ($\geq 0.5 \times 0.5^\circ$) this will likely only be computationally feasible at site or multi-site
4 scale. Note also that performing a number of tests starting from different random
5 “first guess” points in parameter space can help to diagnose if the global minimum
6 has been reached, (as outlined in Section 3.1.6 and discussed at the beginning of
7 Section 3.2), and therefore whether the chosen inversion algorithm is appropriate
8 for optimising your model.

9
10 In addition to the above points we note the following related to a situation in which
11 there is a considerable difference in the number of observations for each data stream. We
12 investigated such a situation in this study with test case 3b, in which only one observation was
13 included for the s_2 data stream instead of the complete time-series. For both models, test case
14 3b showed that a substantial difference in number of observations between the data streams
15 could influence the resulting parameter values and posterior uncertainty (compare test cases
16 3a and b in Fig. 2 for the simple C model and Fig. 4 for the non-linear toy model) as each data
17 stream will have a different overall “weight” in the cost function. Different arguments abound
18 on this issue. Some authors have mentioned the possible need to weight different observation
19 terms in the cost function to increase the influence of data streams with a limited number of
20 observations (e.g. Xu et al., 2006), while others contend that the cost function should not be
21 weighted by the number of observations because the error covariance matrices (**B** and **R**)
22 should already define this weight in an objective way (e.g. Keenan et al., 2013); we would
23 agree with this assertion. Indeed Wutzler et al. (2014) showed that this approach could lead to
24 an overestimate of the posterior uncertainty. As an alternative they proposed a “parameter
25 block” approach in which each data stream only optimises the parameters to which they are
26 most sensitive. We therefore advise modellers not to weight the cost function by the number
27 of observations; instead we suggest adopting an approach such as proposed in Wutzler et al.
28 (2014) and/or ensuring that **B** and **R** matrices are adequately defined. It should not be
29 necessary to weight by the number of observations in the cost function if there is sufficient
30 information to properly build the prior error covariance matrices

31 Several diagnostic tests exist to help infer the relative level of constraint brought about
32 by different data streams, including the observation influence and degrees of freedom of

1 signal metrics (Cardinali et al., 2004). Performing these tests was beyond the scope of this
2 study, particularly given that the simple toy models contained so few parameters, but such
3 tests may be instructive when optimising many hundreds of parameters in a large-scale LSM
4 with a number of different data streams. Furthermore, we strongly suggest performing
5 synthetic experiments with pseudo observations, as in this study, as such tests can help
6 determine the possible constraint brought by different data streams, and the impact of a
7 possible bias and observation or observation–model inconsistency.

8 Aside from multiple data stream assimilation, other promising directions could also be
9 considered to constrain the problem of lack of information in resolving the parameter space
10 within a data assimilation framework, including the use of other ecological and dynamical
11 “rules” that limit the optimisation (see for example Bloom and Williams, 2015), or the
12 addition of different timescales of information extracted from the data such as annual sums
13 (e.g. Keenan et al., 2012). Finally we should also seek to develop collaborations with
14 researchers in other fields who may have advanced further in a particular direction. Members
15 of the atmospheric and hydrological modelling communities, for example, have implemented
16 techniques for inferring the properties of the prior error covariance matrices, including the
17 mean and variance, but also potential biases, autocorrelation and heteroscedasticity, by
18 including these terms as “hyper-parameters” within the inversion (e.g. Michalak et al. 2005;
19 Evin et al., 2014; Renard et al., 2010; Wu et al. 2013). Of course this extends the parameter
20 space – making the problem harder to solve unless sufficient prior information is available
21 (Renard et al., 2010), but such avenues are worth exploring.

22

23 **5 Conclusions**

24 In this study we have attempted to highlight and discuss some of the challenges
25 associated with using multiple data streams to constrain the parameters of LSMs, with a
26 particular focus on the carbon cycle. We demonstrated some of the issues using two simple
27 models constrained with synthetic observations for which the ‘true’ parameters are known.
28 We performed a variety of tests in Section 3 to demonstrate the differences between
29 assimilating each data stream separately, sequentially (in a step-wise approach) and together
30 in the same assimilation (simultaneous approach). In particular we focused on difficulties that
31 may arise in the presence of biases or inconsistencies between the data and the model, as well

1 as non-linearity in the model equations and the importance of accounting of observation error
2 correlations.

3 Many of the issues faced are inherent to all optimisation experiments, including those in
4 which only one data stream is used. It is of utmost importance to determine if the
5 observations contain biases, and/or if inconsistencies or incompatibilities exist between the
6 model and the observations, and to correct for this or properly account for this in the error
7 covariance matrices. Secondly it is crucial to understand the assumptions and limitations
8 related to the inversion algorithm used. Without these two points being met, there is a greater
9 risk of obtaining incorrect parameter values, which may not be obvious by examining the
10 posterior uncertainty and model-data RMSE reduction. Furthermore it is more likely that the
11 implementation of a step-wise versus simultaneous approach will lead to different results.
12 Finally, we note that the consequence of not accounting for cross-correlation between data
13 streams in the prior error covariance matrix becomes more critical with higher observation
14 uncertainty.

15 This study was not able to examine an exhaustive list of all possible challenges that may
16 be faced when assimilating multiple data streams, but we hope that this tutorial style paper
17 will serve as a guide for those wishing to optimise the parameters of LSMs using the variety
18 of C cycle related observations that are available today. We also hope that by increasing
19 awareness about the possible difficulties of model-data integration we can bring the modelling
20 and experimental communities more closely together to focus on these issues.

21

22 **Code availability**

23 The model and inversion code is available via the ORCHIDEE LSM Data Assimilation
24 System (ORCHIDAS) website: https://orchidas.lsce.ipsl.fr/multi_data_stream.php.

25

26 **Acknowledgements**

27 We acknowledge the support from the International Space Science Institute (ISSI). This
28 publication is an outcome of the ISSI's Working Group on "Carbon Cycle Data Assimilation:
29 How to Consistently Assimilate Multiple Data Streams". N. MacBean was also funded by the
30 GEOCARBON Project (ENV.2011.4.1.1-1-283080) within the European Union's 7th
31 Framework Programme for Research and Development. The authors wish to thank colleagues

1 and collaborators in the atmospheric inversion and carbon cycle DA communities with whom
2 they have had numerous past conversations that have led to an improvement in their
3 understanding of the issues presented here.

4

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 Cardinali, C., S. Pezzulli, E. Andersson (2004), Influence-matrix diagnostic of a data
22 assimilation system, *Q. J. R. Meteorol. Soc.*, 130: 2767–2786, doi: 10.1256/qj.03.205
- 23 Chevallier, F., 2007: Impact of correlated observation errors on inverted CO₂ surface fluxes
24 from OCO measurements, *Geophys. Res. Lett.*, 34, L24804, doi:10.1029/2007GL030463.
- 25 Dufresne, J. L., Foujols, M. a., Denvil, S., Caubel, a., Marti, O., Aumont, O., Balkanski, Y.,
26 Bekki, S., Bellenger, H., Benshila, R., Bony, S., Bopp, L., Braconnot, P., Brockmann, P.,
27 Cadule, P., Cheruy, F., Codron, F., Cozic, a., Cugnet, D., de Noblet, N., Duvel, J. P., Ethé,
28 C., Fairhead, L., Fichet, T., Flavoni, S., Friedlingstein, P., Grandpeix, J. Y., Guez, L.,
29 Guilyardi, E., Hauglustaine, D., Hourdin, F., Idelkadi, a., Ghattas, J., Jousaume, S.,
30 Kageyama, M., Krinner, G., Labetoulle, S., Lahellec, a., Lefebvre, M. P., Lefevre, F., Levy,

1 C., Li, Z. X., Lloyd, J., Lott, F., Madec, G., Mancip, M., Marchand, M., Masson, S.,
2 Meurdesoif, Y., Mignot, J., Musat, I., Parouty, S., Polcher, J., Rio, C., Schulz, M.,
3 Swingedouw, D., Szopa, S., Talandier, C., Terray, P., Viovy, N. and Vuichard, N.: Climate
4 change projections using the IPSL-CM5 Earth System Model: From CMIP3 to CMIP5., 2013.

5 Evin, G., Thyer, M., Kavetski, D., McInerney, D. and Kuczera, G.: Comparison of joint
6 versus postprocessor approaches for hydrological uncertainty estimation accounting for error
7 autocorrelation and heteroscedasticity, *Water Resour. Res.*, 50(3), 2350–2375,
8 doi:10.1002/2013WR014185, 2014.

9 Forkel, M., Carvalhais, N., Schaphoff, S., v. Bloh, W., Migliavacca, M., Thurner, M. and
10 Thonicke, K.: Identifying environmental controls on vegetation greenness phenology through
11 model-data integration, *Biogeosciences*, 11, 7025–7050, doi:10.5194/bg-11-7025-2014, 2014.

12 Gobron, N., Pinty, B., Ausedat, O., Chen, J. M., Cohen, W. B., Fensholt, R., Gond, V.,
13 Huemmrich, K. F., Lavergne, T., Mélin, F., Privette, J. L., Sandholt, I., Taberner, M., Turner,
14 D. P., Verstraete, M. M. and Widlowski, J. L.: Evaluation of fraction of absorbed
15 photosynthetically active radiation products for different canopy radiation transfer regimes:
16 Methodology and results using Joint Research Center products derived from SeaWiFS against
17 ground-based estimations, *J. Geophys. Res.*, 111, D13110, doi:10.1029/2005JD006511, 2006.

18 Gobron, N., Pinty, B., Ausedat, O., Taberner, M., Faber, O., Mélin, F., Lavergne, T.,
19 Robustelli, M. and Snoeij, P.: Uncertainty estimates for the FAPAR operational products
20 derived from MERIS - Impact of top-of-atmosphere radiance uncertainties and validation
21 with field data, *Remote Sens. Environ.*, 112(4), 1871–1883, doi:10.1016/j.rse.2007.09.011,
22 2008.

23 Kaminski, T., Knorr, W., Scholze, M., Gobron, N., Pinty, B., Giering, R. and Mathieu, P. P.:
24 Consistent assimilation of MERIS FAPAR and atmospheric CO₂ into a terrestrial vegetation
25 model and interactive mission benefit analysis, *Biogeosciences*, 9(8), 3173–3184,
26 doi:10.5194/bg-9-3173-2012, 2012.

27 Kato, T., Knorr, W., Scholze, M., Veenendaal, E., Kaminski, T., Kattge, J. and Gobron, N.:
28 Simultaneous assimilation of satellite and eddy covariance data for improving terrestrial water
29 and carbon simulations at a semi-arid woodland site in Botswana, *Biogeosciences*, 10(2),
30 789–802, doi:10.5194/bg-10-789-2013, 2013.

31 Keenan, T. F., Davidson, E., Moffat, A. M., Munger, W. and Richardson, A. D.: Using

1 model-data fusion to interpret past trends, and quantify uncertainties in future projections, of
2 terrestrial ecosystem carbon cycling, *Glob. Chang. Biol.*, 18(8), 2555–2569,
3 doi:10.1111/j.1365-2486.2012.02684.x, 2012.

4 Keenan, T. F., Davidson, E. a., Munger, J. W. and Richardson, A. D.: Rate my data:
5 Quantifying the value of ecological data for the development of models of the terrestrial
6 carbon cycle, *Ecol. Appl.*, 23(1), 273–286, doi:10.1890/12-0747.1, 2013.

7 Knorr, W.: Annual and interannual CO₂ exchanges of the terrestrial biosphere: process-based
8 simulations and uncertainties, *Glob. Ecol. Biogeogr.*, 9(3), 225–252, doi:10.1046/j.1365-
9 2699.2000.00159.x, 2000.

10 Krinner, G., Viovy, N., de Noblet-Ducoudré, N., Ogée, J., Polcher, J., Friedlingstein, P.,
11 Ciais, P., Sitch, S. and Prentice, I. C.: A dynamic global vegetation model for studies of the
12 coupled atmosphere-biosphere system, *Global Biogeochem. Cycles*, 19(1), 1–33,
13 doi:10.1029/2003GB002199, 2005.

14 Kuppel, S., F. Chevallier and P. Peylin,: Quantifying the model structural error in Carbon
15 Cycle Data Assimilation Systems. *Geosci. Model Dev.*, 6, 45-55, doi:10.5194/gmd-6-45-
16 2013, 2013.

17 Kuppel, S., Peylin, P., Maignan, F., Chevallier, F., Kiely, G., Montagnani, L., and Cescatti,
18 A.: Model–data fusion across ecosystems: from multisite optimizations to global simulations,
19 *Geosci. Model Dev.*, 7, 2581-2597, doi:10.5194/gmd-7-2581-2014, 2014.

20 MacBean, N., Maignan, F., Peylin, P., Bacour, C., Bréon, F.-M., and Ciais, P.: Using satellite
21 data to improve the leaf phenology of a global terrestrial biosphere model, *Biogeosciences*,
22 12, 7185-7208, doi:10.5194/bg-12-7185-2015, 2015.

23 Michalak, A. M., Hirsch, A., Bruhwiler, L., Gurney, K. R., Peters, W. and co-authors:
24 Maximum likelihood estimation of covariance parameters for Bayesian atmospheric trace gas
25 surface flux inversions. *J. Geophys. Res.* 110, D24107. DOI: 10.1029/2005JD005970, 2005.

26 Morcrette, J.-J.: Evaluation of Model-generated Cloudiness: Satellite-observed and Model-
27 generated Diurnal Variability of Brightness Temperature. *Mon. Wea. Rev.*, **119**, 1205–1224,
28 1991.

29 van Oijen, M., Rougier, J. and Smith, R.: Bayesian calibration of process-based forest models:
30 bridging the gap between models and data, *Tree Physiol.*, 25(7), 915–927,

1 doi:10.1093/treephys/25.7.915, 2005.

2 Peylin, P., Bacour, C., MacBean, N., Leonard, S., Rayner, P. J., Kuppel, S., Koffi, E. N.,
3 Kane, A., Maignan, F., Chevallier, F., Ciais, P., and Prunet, P.: A new step-wise Carbon
4 Cycle Data Assimilation System using multiple data streams to constrain the simulated land
5 surface carbon cycle, *Geosci. Model Dev. Discuss.*, doi:10.5194/gmd-2016-13, in review,
6 2016.

7 Pinnington, E. M., Casella, E., Dance, S. L., Lawless, A. S., Morison, J. I., Nichols, N. K., ...
8 & Quaife, T. L.: Investigating the role of prior and observation error correlations in improving
9 a model forecast of forest carbon balance using Four-dimensional Variational data
10 assimilation, *Agricultural and Forest Meteorology*, 228, 299-314, doi:
11 10.1016/j.agrformet.2016.07.006, 2016.

12 Quaife, T., Lewis, P., De Kauwe, M., Williams, M., Law, B. E., Disney, M. and Bowyer, P.:
13 Assimilating canopy reflectance data into an ecosystem model with an Ensemble Kalman
14 Filter, *Remote Sens. Environ.*, 112(4), 1347–1364, doi:10.1016/j.rse.2007.05.020, 2008.

15 Raoult, N. M., Jupp, T. E., Cox, P. M., and Luke, C. M.: Land-surface parameter optimisation
16 using data assimilation techniques: the adJULES system V1.0, *Geosci. Model Dev.*, 9, 2833-
17 2852, doi:10.5194/gmd-9-2833-2016, 2016

18 Raupach, M. R.: Dynamics of resource production and utilisation in two-component
19 biosphere-human and terrestrial carbon systems, *Hydrol. Earth Syst. Sci.*, 11, 875–889,
20 doi:10.5194/hess-11-875-2007, 2007.

21 Raupach, M. R., Rayner, P. J., Barrett, D. J., Defries, R. S., Heimann, M., Ojima, D. S.,
22 Quegan, S. and Schimmler, C. C.: Model-data synthesis in terrestrial carbon observation:
23 Methods, data requirements and data uncertainty specifications, *Glob. Chang. Biol.*, 11(3),
24 378–397, doi:10.1111/j.1365-2486.2005.00917.x, 2005.

25 Rayner, P. J., Scholze, M., Knorr, W., Kaminski, T., Giering, R. and Widmann, H.: Two
26 decades of terrestrial carbon fluxes from a carbon cycle data assimilation system (CCDAS), ,
27 19, doi:10.1029/2004GB002254, 2005.

28 Renard, B., Kavetski, D., Kuczera, G., Thyer, M. and Franks, S. W.: Understanding predictive
29 uncertainty in hydrologic modeling: The challenge of identifying input and structural errors,
30 *Water Resour. Res.*, 46(5), 1–22, doi:10.1029/2009WR008328, 2010.

1 Richardson, A. D., Williams, M., Hollinger, D. Y., Moore, D. J. P., Dail, D. B., Davidson, E.
2 a., Scott, N. a., Evans, R. S., Hughes, H., Lee, J. T., Rodrigues, C. and Savage, K.: Estimating
3 parameters of a forest ecosystem C model with measurements of stocks and fluxes as joint
4 constraints, *Oecologia*, 164(1), 25–40, doi:10.1007/s00442-010-1628-y, 2010.

5 Schürmann, G. J., Kaminski, T., Köstler, C., Carvalhais, N., Voßbeck, M., Kattge, J., Giering,
6 R., Rödenbeck, C., Heimann, M., and Zaehle, S.: Constraining a land surface model with
7 multiple observations by application of the MPI-Carbon Cycle Data Assimilation System,
8 *Geosci. Model Dev. Discuss.*, doi:10.5194/gmd-2015-263, in review, 2016.

9 Scholze, M., T. Kaminski, W. Knorr, S. Blessing, M. Vossbeck, J.p. Grant, and K. Scipal.
10 "Simultaneous Assimilation of SMOS Soil Moisture and Atmospheric CO2 In-situ
11 Observations to Constrain the Global Terrestrial Carbon Cycle." *Remote Sensing of*
12 *Environment* 180, 334-45, 2016.

13 Sitch, S., Friedlingstein, P., Gruber, N., Jones, S. D., Murray-Tortarolo, G., Ahlström, A.,
14 Doney, S. C., Graven, H., Heinze, C., Huntingford, C., Levis, S., Levy, P. E., Lomas, M.,
15 Poulter, B., Viovy, N., Zaehle, S., Zeng, N., Arneeth, A., Bonan, G., Bopp, L., Canadell, J. G.,
16 Chevallier, F., Ciais, P., Ellis, R., Gloor, M., Peylin, P., Piao, S., Le Quéré, C., Smith, B.,
17 Zhu, Z. and Myneni, R.: Recent trends and drivers of regional sources and sinks of carbon
18 dioxide, *Biogeosciences*, 12, 653–679, doi:10.5194/bgd-12-653-2015, 2015.

19 Thum, T., N. MacBean, P. Peylin, C. Bacour, D. Santaren, B. Longdoz, D. Loustau and P.
20 Ciais, The potential benefit of using forest biomass data in addition to carbon and water flux
21 measurements to constrain ecosystem model parameters: case studies at two temperate forest
22 sites. In revision for *Agric. For. Meteorol.*

23 Trudinger, C. M., Raupach, M. R., Rayner, P. J., Kattge, J., Liu, Q., Park, B., Reichstein, M.,
24 Renzullo, L., Richardson, A. D., Roxburgh, S. H., Styles, J., Wang, Y. P., Briggs, P., Barrett,
25 D. and Nikolova, S.: OptIC project: An intercomparison of optimization techniques for
26 parameter estimation in terrestrial biogeochemical models, *J. Geophys. Res. Biogeosciences*,
27 112(2), doi:10.1029/2006JG000367, 2007.

28 Williams, M., Schwarz, P. a, Law, B. E., Irvine, J. and Kurpius, M. R.: An improved analysis
29 of forest carbon dynamics using data assimilation, *Glob. Chang. Biol.*, 11(1), 89–105,
30 doi:10.1111/j.1365-2486.2004.00891.x, 2005.

31 Wu, L., M. Bocquet, F. Chevallier, T. Lauvaux, and K. Davis, 2013: Hyperparameter

1 estimation for uncertainty quantification in mesoscale carbon dioxide inversions. *Tellus B*, 65,
2 doi:10.3402/tellusb.v65i0.20894.

3 Xu, T., White, L., Hui, D. and Luo, Y.: Probabilistic inversion of a terrestrial ecosystem
4 model: Analysis of uncertainty in parameter estimation and model prediction, *Global*
5 *Biogeochem. Cycles*, 20(2), 1–15, doi:10.1029/2005GB002468, 2006.

6 Zobitz, J. M., Moore, D. J. P., Quaipe, T., Braswell, B. H., Bergeson, A., Anthony, J. a. and
7 Monson, R. K.: Joint data assimilation of satellite reflectance and net ecosystem exchange
8 data constrains ecosystem carbon fluxes at a high-elevation subalpine forest, *Agric. For.*
9 *Meteorol.*, 195-196, 73–88, doi:10.1016/j.agrformet.2014.04.011, 2014.

10

11

12

13

14

15

16

17

18

19

20

21

22

23

24

25

26

27

1 Table 1: The optimisation set-up for both models, including the true parameter values, their
 2 range and the observation uncertainty (1 sigma), which was set to 10% of the mean value for
 3 each set of pseudo-observations derived from multiple first guesses of the model. The
 4 parameter uncertainty (1 sigma) was set to 40% of the range for each parameter.

5

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

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

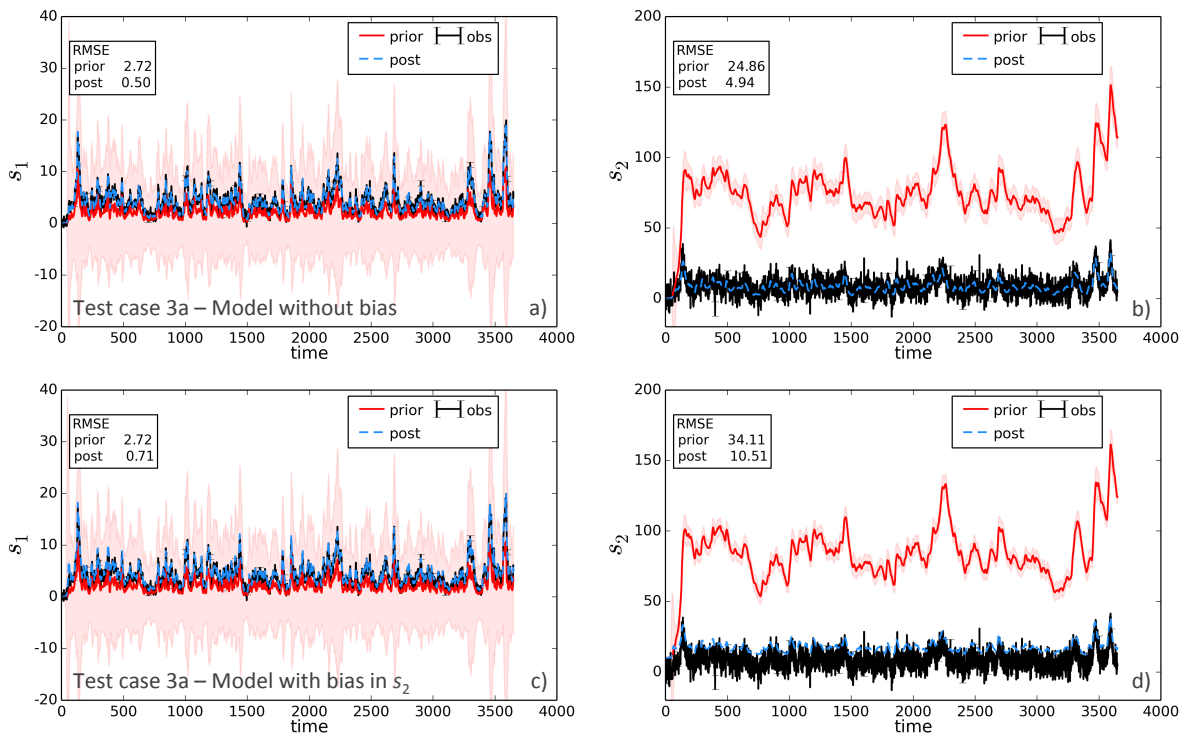
10

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.

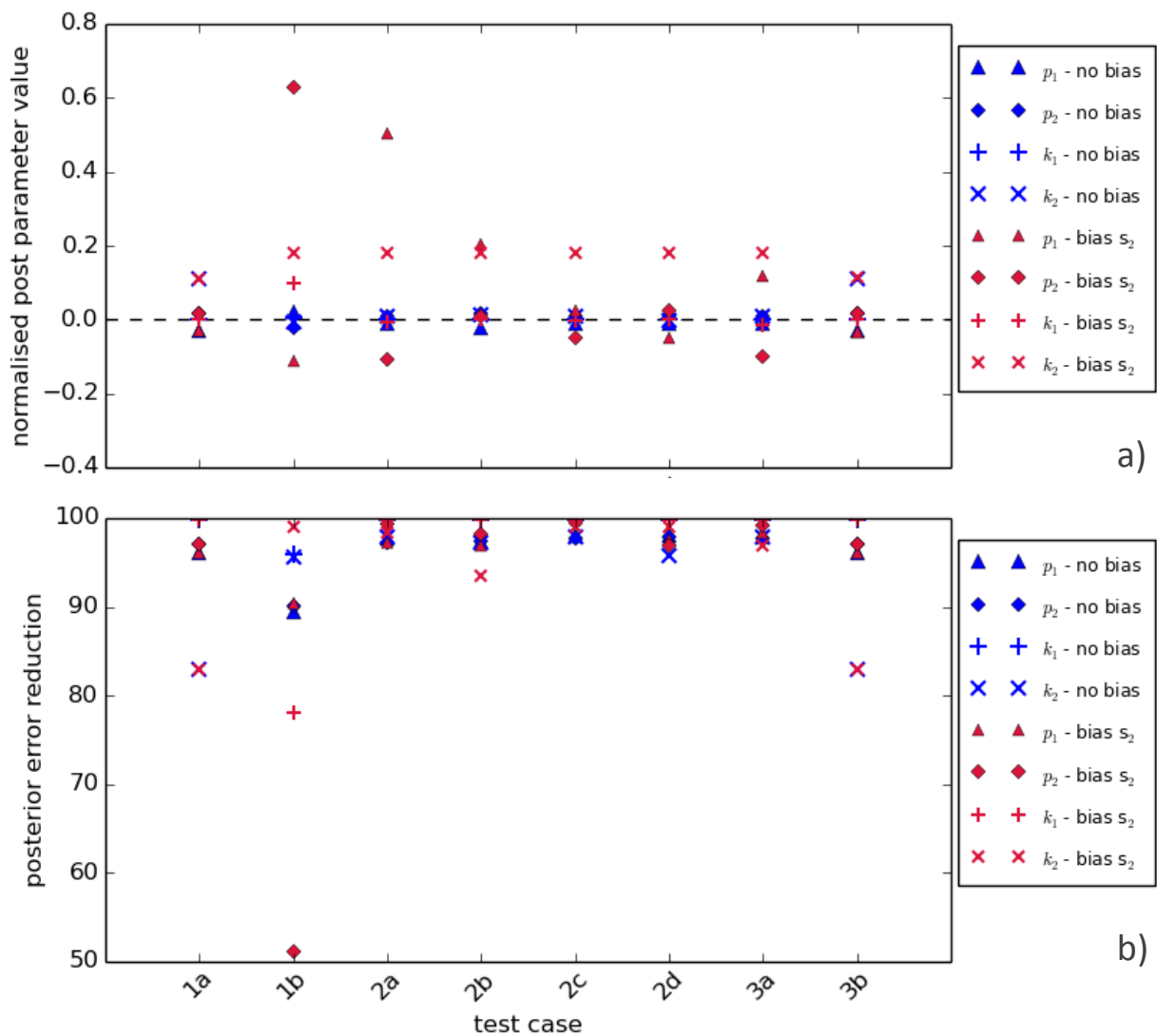
9

10

11

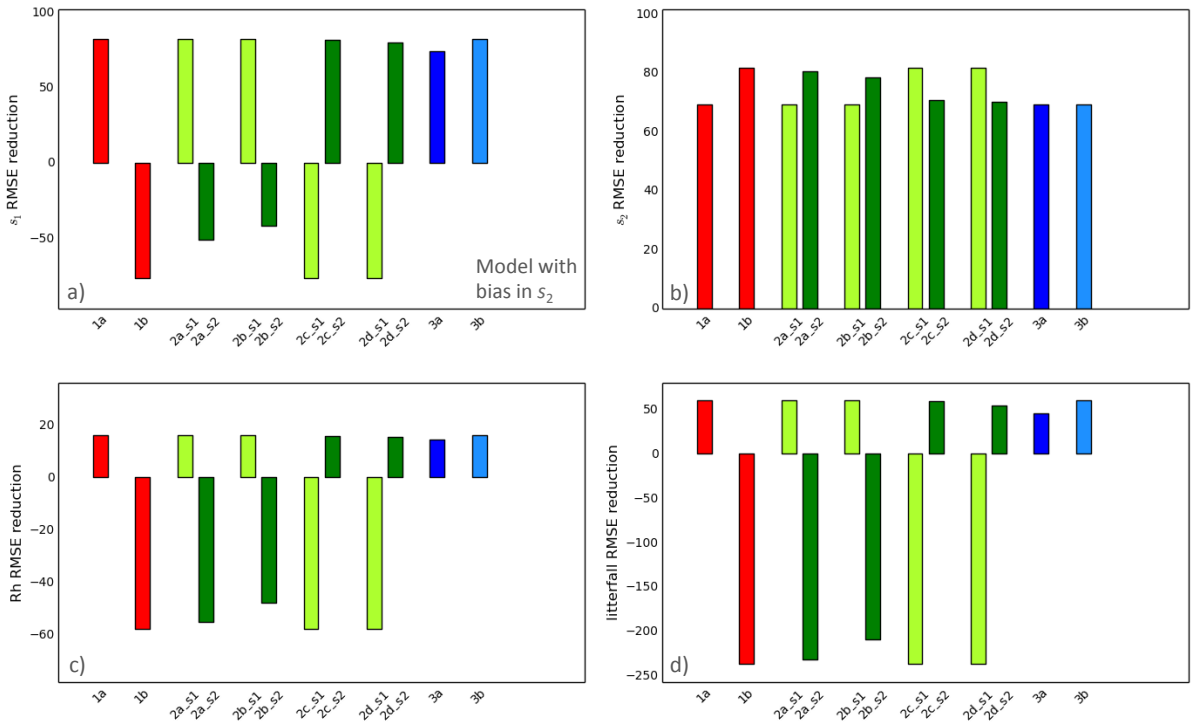
12

13



1
2
3
4
5
6
7
8
9
10
11
12
13

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



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

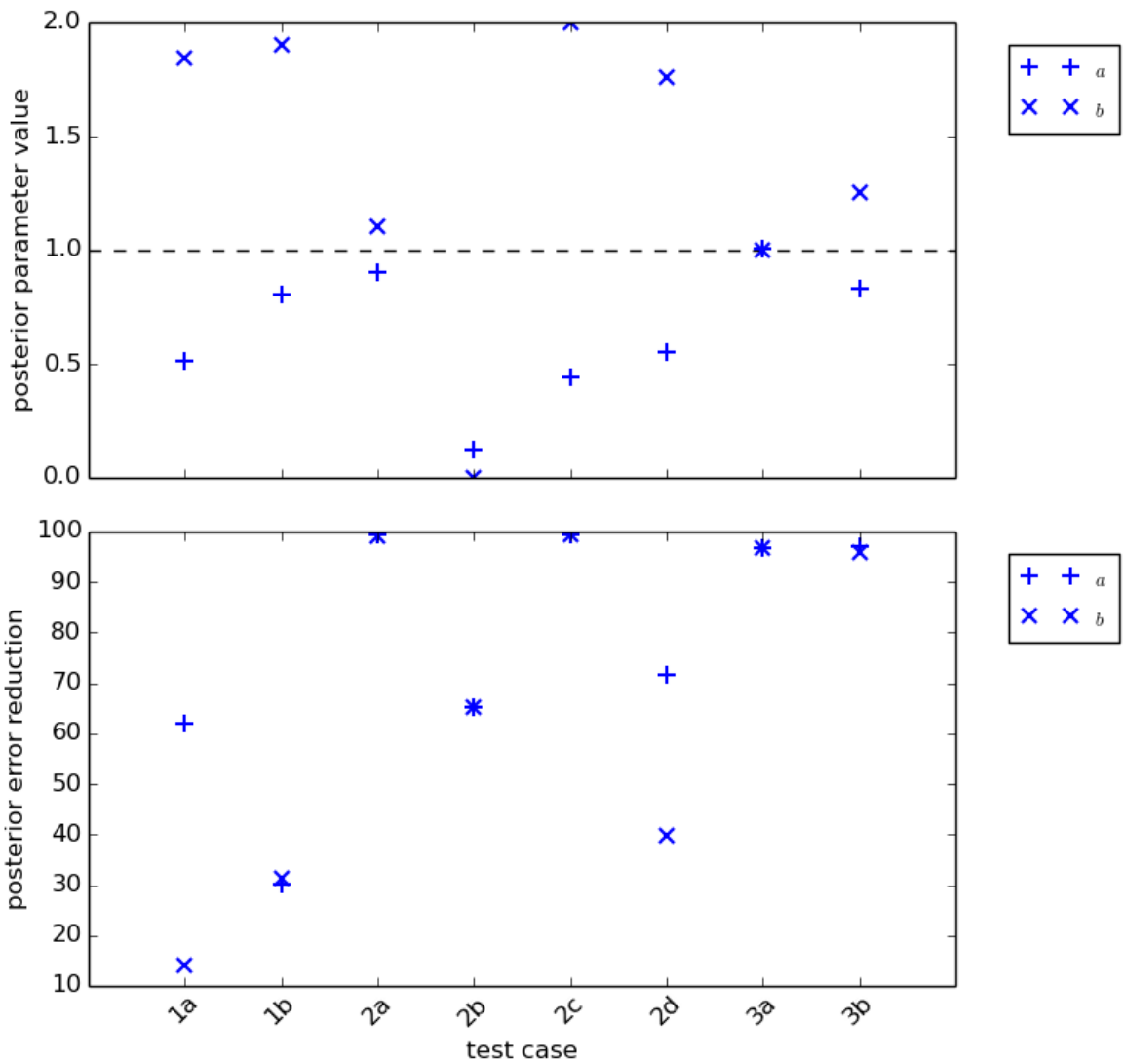
13

14

15

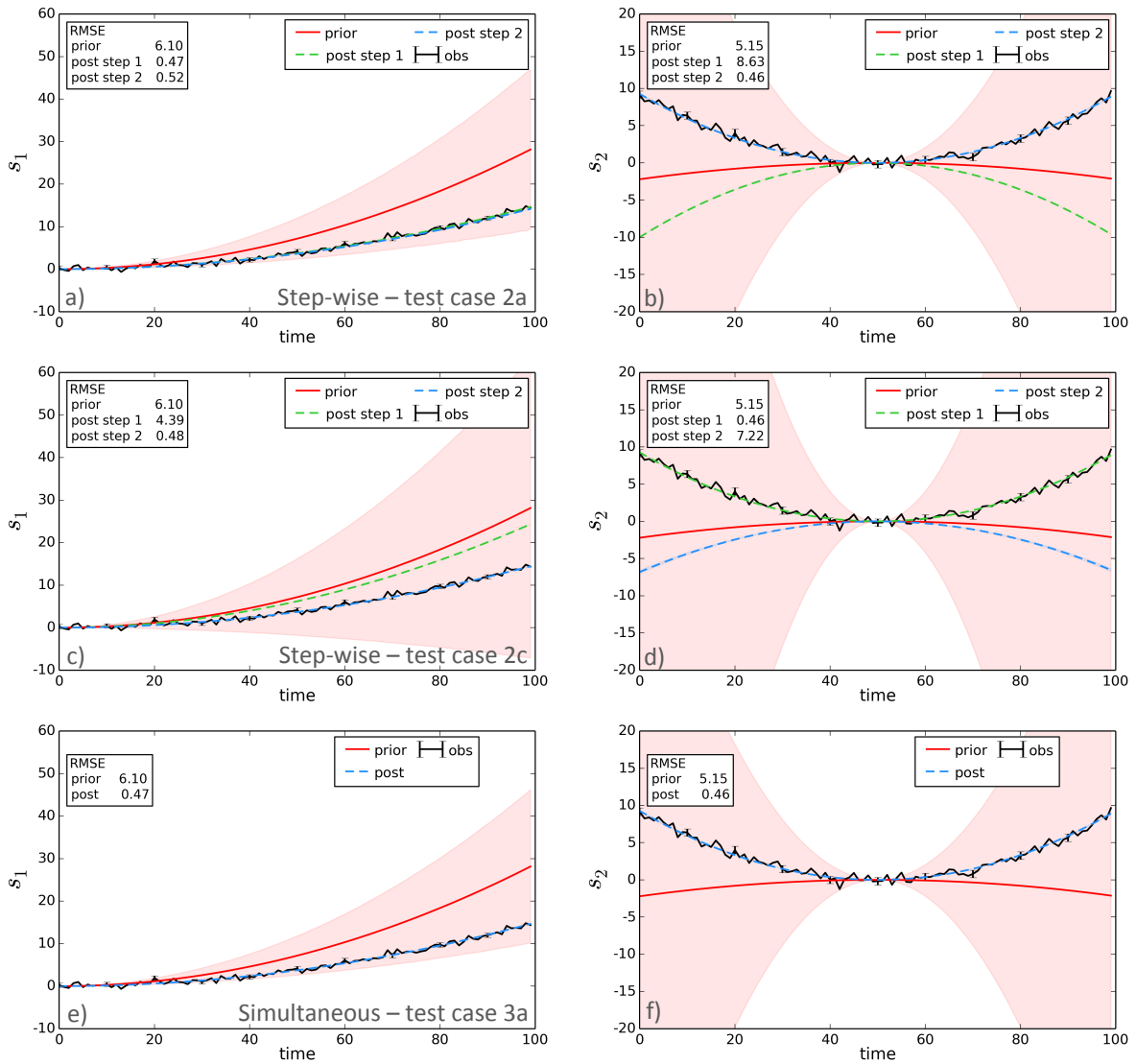
16

17



1
2
3
4
5
6
7
8
9
10
11
12

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



1

2

3

4

5

6

7

8

9

10

11

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

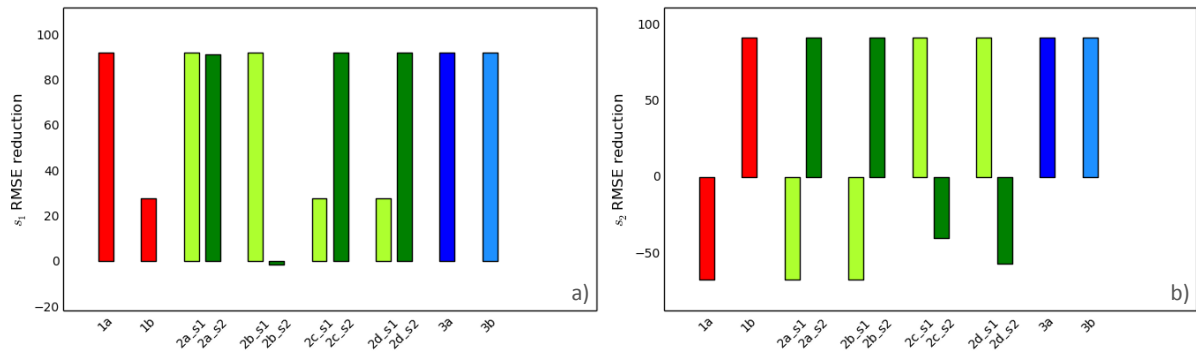
22

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

12

13

14

15

16

17

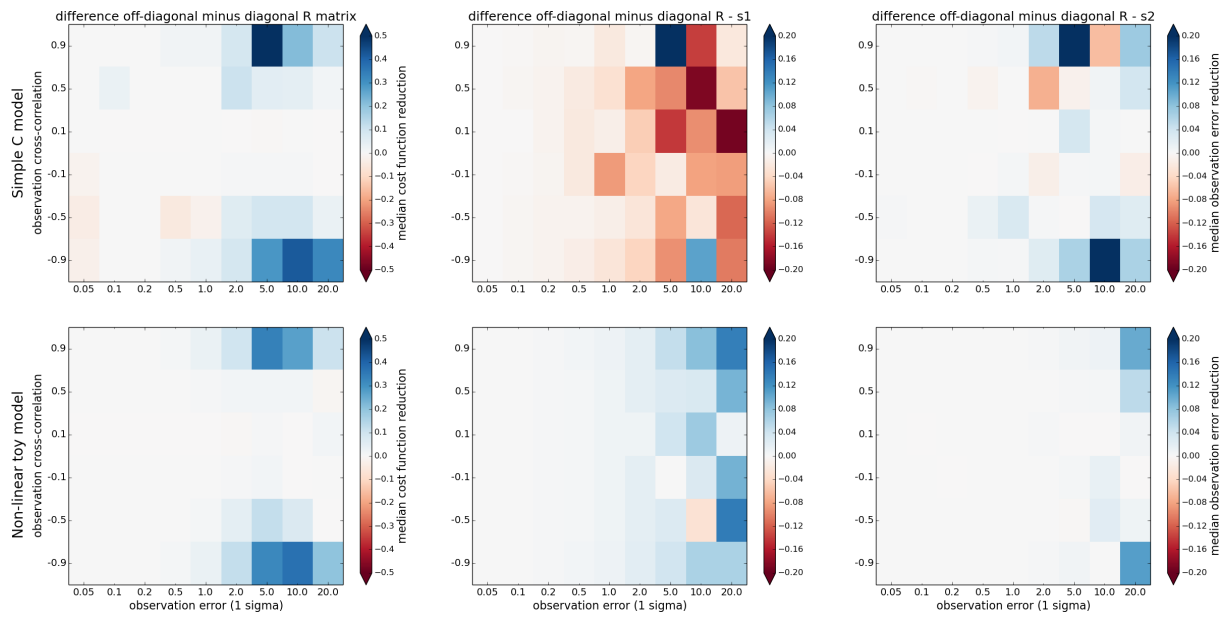
18

19

20

21

22



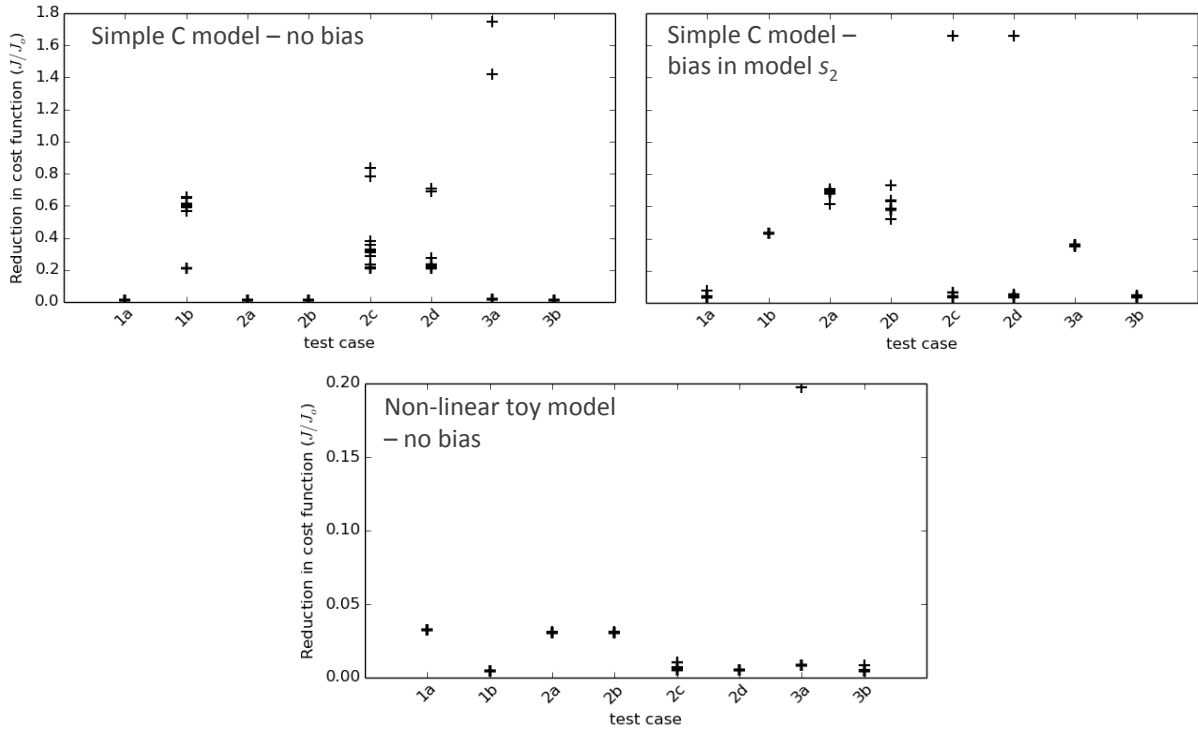
1

2 Figure 7: Median difference (across 20 first guess parameters) between including correlated
 3 observation errors in the R matrix (off-diagonal elements) minus ignoring the correlated
 4 observation errors (keeping R diagonal) for the reduction cost function (a and d: left column)
 5 and the reduction in s_1 and s_2 observation errors (b, c, e and f: middle and right columns), for
 6 both the simple C model (a, b and c: top row) and the non-linear toy model (d, e and f: bottom
 7 row) for a range of observation errors (x-axes) and correlation (y-axes) – see Section 3.1.6.
 8 The reduction is calculated as $1 - (\text{posterior}/\text{prior})$.

9

10

1 Supplementary material



2

3 Figure S1: Reduction in the cost function (J/J_0) for each model and each test for all 20
 4 assimilations with different random “first guess” points in the parameter space (i.e. each cross
 5 represents the 20 random “first guess” tests). Top panel – simple C model without bias (left)
 6 and with bias added to the simulated s_2 variable (right). Bottom panel – non-linear toy model
 7 with no added bias. Note that the majority of the random “first guess” assimilations achieve
 8 the same reduction in the cost function even though the final value is different for each test,
 9 which is to be expected as each test (for each model) has a different cost function.