

Bold black font: Topical editor and reviewer comments

Black font: Author response

Blue font: Verbatim copy and paste from the revised manuscript

Topical Editor Decision: Publish subject to minor revisions (review by editor) (08 Apr 2019) by Christoph Müller

Comments to the Author:

Dear Dr. Elshall,

the original reviewers have seen your revised paper again and suggest further amendments to the paper. I think what they suggest is easy enough to implement.

I look forward to your revised paper.

Cheers

Christoph

Thank you very much for handling the original and the revised submissions. We have responded to all the Editorial Support and reviewers' comments, and revised the manuscript accordingly.

Editorial Support

Besides adjustments requested by the Topical Editor or Referees, please check your manuscript carefully for typos, missing co-authors and their affiliations, terminology, updates of data in tables, or updates of variables in equations.

We updated the affiliation of the second author and the funding information. We improved the article keywords. We carefully checked the manuscript, and corrected several typos and grammatical errors, as shown in the *marked-up manuscript version*. We also slightly improved the writing style in several parts of the article and added more clarifications as shown in the *marked-up manuscript version*.

Anonymous Referee #1

I revisit some of my previous comments, considering response from the authors in the discussion forum.

1. Contribution: the authors have now better articulated their contribution. The paper is somewhat incremental, since the methodology is not new and the findings are similar to applications in other disciplines. Nevertheless it is still useful to have an explicit evaluation and comparison of the impacts of statistical assumptions in soil respiration models.

Thank you for taking the time to review the manuscript and for your valuable feedback.

2. Problems with residual autocorrelation: I agree that the authors do not need to solve this issue in this paper. However, the alternative suggested approach (Evin et al) could be easily tested (simply swap the order in which correlation and heteroscedasticity are applied in the likelihood function). In their extended response in the discussion forum the authors seem to indicate that they have tested this. Why then not add it to the paper? If not in the results then in the discussion part of the paper ("preliminary

results show..."). My suggestion is to include the entire extended response posted in the discussion forum (under comment 2 of RC1) in the paper itself, as it provides more information that is useful for readers.

We prefer to keep the manuscript in its current form since we are afraid that we cannot adequately address this autocorrelation problem in this manuscript. We prefer not to present incomplete or preliminary results in a published paper. However, we improved the extended response, and add it to the manuscript as suggested by the reviewer. The added part reads:

4.3 Discussion on handling residual correlation

Accounting for autocorrelation can lead to biased parameter estimation (Figure 5) and poor predictive performance (Figure 10). Auto-correlated residuals may be attributed to model discrepancy, as shown in Lu et al. (2013). The most obvious solution to handle the autocorrelation is to reduce the autocorrelation by improving the soil respiration model. If model improvement is difficult for practical reasons, we can improve the data model to better characterize the autocorrelation. Addressing autocorrelation in a data model is challenging since it involves several interlinked factors as follows:

- (1) Non-stationarity due to wet-dry periods could be a reason for this problem. By drawing on similarity from surface hydrology, the study of Ammann et al. (2018) suggests that auto-correlated residuals might be attributed to non-stationarity due to wet-dry periods with half-hourly data. Accounting for non-stationarity could better address the problem of auto-correlated residuals (Ammann et al., 2018; Smith et al., 2010b).
- (2) The way of implementing autocorrelation could have an impact. Autocorrelation could be applied to raw residuals directly (e.g., Li et al., 2015), to transformed residuals based on covariance matrix of residuals $L(\mathbf{e})$ (e.g., Lu et al., 2013), or to normalized residuals $L(\mathbf{a})$ (e.g., Schoups and Vrugt, 2010; Evin et al., 2013). Note that \mathbf{e} is a vector of transformed residuals, while \mathbf{a} denotes a vector of independent and identically distributed random errors with zero mean and unit standard deviation. The $L(\mathbf{e})$ approach based on covariance matrix of residuals is generally limited to Gaussian data models (e.g. Lu et al., 2013), while the $L(\mathbf{a})$ approach for normalized residuals can be readily adopted for non-Gaussian data models.
- (3) The autocorrelation model could have an impact. Using an autoregressive model is a popular technique to account for auto-correlated residuals. However, using an autoregressive model with either joint inversion approach (e.g., this study and Schoups and Vrugt, 2010) or sequential approaches (e.g., Evin et al., 2013, 2014; Lu et al., 2013) removes correlation errors through a filter approach, which can lead to a loss of information content. As this may cause overcorrection of prediction especially at surge events, Li et al. (2015) developed a restricted autoregressive model to overcome this adverse effect. Other autocorrelation models include moving average model and mixed autoregressive-moving averaging model (Chatfield, 2004).
- (4) Joint versus sequential inversion for autocorrelation could have an impact. Sequential inversion approaches include two-step procedures (e.g. Evin et al., 2013, 2014; Lu et al., 2013) or the multi-step procedure (Li et al., 2016a). These sequential approach estimates the autoregressive parameters sequentially in a later step after estimating the physical model parameters and other data model parameters. Evin et al. (2013, 2014) used a sequential approach to avoid the interaction between the parameters of the heteroscedasticity model and the autocorrelation model. In addition, the autoregressive model parameters can be deterministically calculated as an internal variables of the data model similar to Lu et al. (2013), and not as calibration parameters (e.g. Schoups and Vrugt; Evin et al. 2013; 2014). While the first step in the sequential approach would avoid the biased parameter estimation (Figure 10a-d), the second step can still lead a poor predicative performance since we are essentially using a filter approach to remove

residual correlation. To address this problem, Li et al. (2016) multi-step procedure that is based on Gaussian data model uses restricted autoregressive model. Generally, Ammann et al. (2018) states that the joint inversion is still preferred, and understanding the conditions where accounting for auto-correlation can be achieved remain poorly understood.

In addition, the text about autocorrelation in the conclusions section was accordingly shortened. The revised manuscript reads “While the reasons remain poorly understood (Ammann et al., 2018), it might be attributed to non-stationarity due to wet-dry periods with half-hourly data (Ammann et al., 2018) or to the method of handling autocorrelation (e.g., Schoups and Vrugt, 2010, Evin et al., 2013; 2014; Lu et al., 2013; Li et al., 2015, 2016a; Ammann et al. 2018). Further investigation to address autocorrelation in soil respiration modeling is warranted in a future study.”

3. Grammatical/spelling errors: the authors state that they have "corrected several other grammatical errors", but it's not clear what was corrected exactly.

Sorry for not listing the grammatical errors and typos we corrected in the previous submission, which are as follows:

- Line 123: was used to select ~~and~~ the best model -> was used to select the best model
- Line 192: Laplace distribution ~~used by~~ (van Wijk et al., 2008) ~~and~~ (Ricciuto et al., 2011) -> Laplace distribution (van Wijk et al., 2008; Ricciuto et al., 2011).
- Line 207: ~~to split~~ a dataset of CO₂ -> ~~by splitting the~~ dataset of CO₂
- Line 329: model parameters ~~are obtained~~ jointly with -> model parameters jointly ~~estimated~~ with
- Line 420: each ~~criteria~~ -> each ~~criterion~~
- Line 658: Accounting for [^] error term ~~e~~ -> Accounting for ~~the~~ error term ~~e~~
- Line 688: the rest [^] six data models -> the rest ~~of the~~ six data models
- Line 722: the residuals, [^] thus resulting -> the residuals, ~~and~~ thus resulting
- Line 743: We tested eight data modeling -> We tested eight data models
- Line 807: The [^] conclusions ~~above~~ -> The ~~above~~ conclusions

Additional comment:

- the abstract states that "not accounting for heteroscedasticity ... will definitely underestimate uncertainty". It is not clear why this would be the case. Perhaps you mean to say that not accounting for any residual error beyond parameter uncertainty will underestimate predictive uncertainty (as in section 4.2)? Although that is not really a significant finding.

Comparing Figures 10a-b (with the SLS data model that does not account for residual error beyond parameter uncertainty) with Figures 10c-d (with the SEP data model that accounts for residual error beyond parameter uncertainty) shows that using an SEP data model with two additional parameters did not significantly impact the uncertainty. This is not the case for Figures 10e-h for data models WLS and WSEP that account for heteroscedasticity. Visual comparison of Figures 10a-d with Figures 10e-h and examination of the sharpness and predictive coverage metrics indicate that not accounting for heteroscedasticity will underestimate uncertainty. Accordingly, accounting for heteroscedasticity with WLS (Figures 10e-f) or WSELP (Figures 10g-h) makes the predictions more sensitive to peak carbon effluxes. We clarified this point, and the revised manuscript reads “Not accounting for heteroscedasticity will underestimate the predication uncertainty (Figure 10b and Figure 10d). This is mainly because the variance of the efflux residuals increases with the magnitude of the carbon effluxes (Figure 3a), and thus

assuming constant variance is not representative. Accordingly, accounting for heteroscedasticity using WLS (Figure 10e) or WSEP (Figure 10h) will make the predictions more sensitive to peck carbon effluxes. This and will generally improve the predictive coverage on the expense of sharpness and the central mean tendency.”

Anonymous Referee #2

The revised version submitted by Elshall et al., generally answers quite well all my previous comment.

Thank you for taking the time to review the original and the revise manuscripts.

Nevertheless I still disagree with using a CUE maximum value of 0.6 and I would like to see some kind of sensitivity analysis to better understand the impact of such assumption. If the effect of changing the max CUE value is limited then the paper could be published in its current version. If not some discussions should be added.

We probably did not make ourselves clear. For the inverse modeling with MCMC sampling we did not assume CUE maximum value of 0.6. Thus, for parameter estimation and predictive performance, we did not impose this constraint. We merely obtained this CUE maximum value of 0.6 from literature, which is based on thermodynamic calculation (Fernández-Martínez et al., 2014; Li et al., 2014; Sinsabaugh, et al., 2013), to evaluate whether the posterior parameter distributions of CUE under different data models and different soil respiration models are within this physically reasonable range of 0 ~ 0.6 or beyond. Thus, this assumption has no impact on the results. We clarified this in the revised manuscript as follows “Note that, for inverse modeling with MCMC sampling, we did not assume CUE maximum value of 0.6. In other words, for parameter estimation and predictive performance we did not impose the constraint that CUE is less than 0.6. We merely use this CUE maximum value of 0.6 to evaluate whether the posterior CUE parameter samples obtained using different data models and different soil respiration models are within the physically reasonable range of 0 ~ 0.6.”

Minor corrections:

L376 you mean soil moisture I guess?

L424 you mean Y_i and not X_i , right?

Thank you for these two typos and we corrected both of them.

Bayesian Inference and Predictive Performance of Soil Respiration Models in the Presence of Model Discrepancy

Ahmed S. Elshall^{1,2}, Ming Ye^{3,*}, Guo-Yue Niu^{4,5} and Greg A. Barron-Gafford^{4,6}

¹ Department of Earth Sciences, University of Hawai'i Manoa, Honolulu, Hawaii, USA

² Water Resources Research Center, University of Hawai'i Manoa, Honolulu, Hawaii, USA

³ Department of [Earth, Ocean, and Atmospheric Science](#)~~Scientific Computing~~, Florida State University, Tallahassee, Florida, [USA](#)

⁴ Biosphere 2, University of Arizona, Tucson, Arizona, [USA](#)

⁵ Department of Hydrology and Water Resources, University of Arizona, Tucson, Arizona, [USA](#)

⁶ School of Geography and Development, University of Arizona, Tucson, Arizona, [USA](#)

*Corresponding Author: Ming Ye, Telephone: (850) 644-4587, Email: mye@fsu.edu

Submitted for publication in Geoscientific Model Development

[March](#)~~April~~, 2019

Key Points

- (1) Bayesian inference and prediction are useful to evaluate multiple soil respiration models with different levels of model complexity.
- (2) Data models used in Bayesian inference have substantial impacts on model parameter distributions and subsequently model predictions.
- (3) Using exponential power distribution and considering heteroscedasticity in data models improves Bayesian inference and prediction.

Keywords: Soil respiration, [modeling](#), Bayesian, likelihood function, data model, autocorrelation, heteroscedasticity, skew exponential power distribution, cross-validation, [relative model scorescoring rule](#)

Abstract

Bayesian inference of microbial soil respiration models is often based on the assumptions that the residuals are independent (i.e. no temporal or spatial correlation), identically distributed (i.e. Gaussian noise) and with constant variance (i.e. homoscedastic). In the presence of model discrepancy, since no model is perfect, this study shows that these assumptions are generally invalid in soil respiration modeling such that residuals have high temporal correlation, an increasing variance with increasing magnitude of CO₂ efflux, and non-Gaussian distribution. Relaxing these three assumptions stepwise results in eight data models. Data models are the basis of formulating likelihood functions of Bayesian inference. This study presents a systematic and comprehensive investigation of the impacts of data model selection on Bayesian inference and predictive performance. We use three mechanistic soil respiration models with different levels of model fidelity (i.e. model discrepancy) with respect to number of carbon pools and explicit representations of soil moisture controls on carbon degradation, and accordingly have different levels of model complexity with respect to the number of model parameters. The study shows data models have substantial impacts on Bayesian inference and predictive performance of the soil respiration models such that: (i) the level of complexity of the best model is generally justified by the cross-validation results for different data models; (ii) not accounting for heteroscedasticity and autocorrelation might not necessarily result in biased parameter estimates or predictions, but will definitely underestimate uncertainty; (iii) using a non-Gaussian data model improves the parameter

estimates and the predictive performance; and (iv) separate accounting for autocorrelation or joint inversion of correlation and heteroscedasticity can be problematic and requires special treatment. Although the conclusions of this study are empirical, the analysis may provide insights for selecting appropriate data models for soil respiration modeling.

1 Introduction

Developing accurate soil respiration models is important for realistic projection of global carbon [C] cycle, as global soils store 2,300Pg carbon, an amount more than 3 times that of the atmosphere (Schmidt et al., 2011) and release 60–75 Pg C/yr, about 7 times more CO₂ to the atmosphere than all human-caused emissions (Le Quéré et al., 2014). The major work on soil respiration modeling has been focused on advancing knowledge about model inputs and calibration data (e.g. Janssens et al., 2003; Peters et al., 2007; Scott et al., 2009; Barron-Gafford et al., 2011; Hilton et al., 2014) and on developing more advanced models for better representing soil microbial processes (e.g. Schimel and Weintraub, 2003; Allison et al., 2010; Davidson et al., 2011; Wieder et al., 2013, 2015; Xu et al., 2014; Zhang et al., 2014) . Integration of data and models is indispensable for improving predictability of the terrestrial carbon cycle, and statistical modeling is a vital tool for the model-data integration (Luo et al., 2011, 2014; Wieder et al., 2015). In addition, use of state-of-the-art statistical methods is necessary to accurately quantify uncertainty in parameters and structures of soil respiration models for improvement and practical uses of the models (Katz et al., 2013). A data model that is also known as a residuals model or an error model is used to characterize residuals (i.e., the difference between data and corresponding model simulations). While a large number of data models have been used (e.g. Elshall et al., 2018; Scholz et al., 2018) to our knowledge comprehensive and systematic evaluation of data models for soil respiration modeling has not been reported in literature.

The ~~objectivesgoal~~ of this study ~~areis~~ to evaluate the impacts of data models on Bayesian inference and predictive performance of three mechanistic soil respiration models, and ~~to use these~~ ~~evaluation results findings~~ to make broader recommendations. The three models were developed by Zhang et al. (2014) to simulate the Birch effect (the peak soil microbial respiration pulses in response to episodic rainfall pulses) at a site scale and a short temporal scale; ~~understanding the; which are~~ Birch effect is important for gaining mechanistic understanding of CO₂ efflux production (Högberg and Read, 2006; Vargas et al., 2011). ~~The models of Z~~ Zhang et al. (2014) ~~developed a total five models, including~~ are based on an existing four-carbon pool model, ~~but have -and four new models with-~~ additional carbon pools and/or explicit representations of soil moisture controls on carbon degradation and microbial uptake rates. The models ~~Zhang et al. (2014)~~ were calibrated, and Bayesian model selection was used to select the best model (Zhang et al., 2014). However, this effort was based on a single data model. It is unknown whether the best model still remains the best (in terms of reproducing the both calibration data and the cross-validation data) if a different data model is used. In addition, since predictive performance of the models was not evaluated in Zhang et al. (2014), it is unknown whether the best model will give the best predictions. These two questions are addressed in this study by considering eight data models and by evaluating predictive performance in a manner of cross-validation. The top two models (also the two most high fidelity models) ranked by Zhang et al. (2014) are considered in this study, and the worst model (also the low fidelity model) is also considered in this study for comparison. We use the terms model fidelity and model discrepancy interchangeably. Model fidelity refers to the degree of realism of representing our scientific knowledge with respect to the real world system. That is a high fidelity model has less discrepancy. Conducting Bayesian inference and Evaluating

predictive performance for the three models with different degrees of fidelity provides more insights than ~~for~~ a single model.

Bayesian inference in general uses the Bayes' theorem to update the prior distributions of model parameters to posterior parameter distributions given a likelihood function of data. The mathematical formulation of the (formal and informal) likelihood function requires a probabilistic data model that however is intrinsically unknown due to unknown errors in all model components such as ~~observation data~~, model structures, parameters, and driving forces. Bayesian inference of soil respiration models often adopts the assumption of independent, normally distributed and homoscedastic residuals (e.g. Ahrens et al., 2014; Bagnara et al., 2015, 2018; Barr et al., 2013; Barron-gafford et al., 2014; Braakhekke et al., 2014; Braswell et al., 2015; Correia et al., 2012; Du et al., 2015, 2017; Hararuk et al., 2014; Hashimoto et al., 2011; He et al., 2018; Klemedtsson et al., 2008; Menichetti et al., 2016; Raich et al., 2002; Ren et al., 2013; Richardson and Hollinger, 2005; Steinacher and Joos, 2016; Tucker et al., 2014; Tuomi et al., 2008; Xu et al., 2006; Yeluripati et al., 2009; Yuan et al., 2012, 2016; Zhang et al., 2014; Zhou et al., 2010). These assumptions are conveniently adopted ~~to satisfy~~ ~~sinee~~ the requirement of using an unknown probability model in Bayesian statistics, which is called “a basic dilemma” by (Box and Tiao, 1992). ~~Box and Tiao (1992).~~

Postulating the data models is always based on assumptions about residual statistics, and the most widely used assumptions are paired as follows: (i) independent vs. correlated residuals, (ii) homoscedastic vs. heteroscedastic residuals, and (iii) Gaussian vs. non-Gaussian residuals. For soil respiration modeling few studies have relaxed the non-correlation assumption (e.g. Cable et al., 2008, 2011; Li et al., 2016b), the homoscedasticity assumption (e.g. Berryman et al., 2018; Elshall et al., 2018; Ogle et al., 2016; Tucker et al., 2013), and the non-Gaussian and homoscedasticity

assumptions (e.g. Elshall et al., 2018; Ishikura et al., 2017; Kim et al., 2014). ~~The~~A recent study of (Scholz et al., (2018) relaxed these three assumptions using the generalized likelihood function developed by (Schoups and Vrugt, (2010). However, few studies have focused on investigating appropriateness and impact of these assumptions for soil respiration modeling, by relaxing the independent residuals assumption (Ricciuto et al., 2011) and the Gaussian residuals assumption (Ricciuto et al., 2011; van Wijk et al., 2008). By relaxing these three assumptions stepwise resulting in eight data models, to our knowledge this is the first study that systematically evaluates the impact of data model selection on Bayesian inference and predictive performance of soil respiration modeling. In addition, to our knowledge this is the first soil respiration modeling study that investigates the impact of data models in relation to model fidelity.

Relaxing these three assumption results in eight data models, which are shown in details in Section 2. For example, combining the assumptions of independent, homoscedastic, and Gaussian residuals leads to the standard least squares data model. This model is the simplest one among the eight data models, since it requires only one parameter, i.e., the constant variance of the Gaussian distribution. Note that there is a difference between the soil respiration model parameters and the data model parameters. They technically can be jointly estimated~~together~~, but one arises from assumptions about soil respiration processes, and the other from assumptions about the residuals. Relaxing the homoscedastic assumption to heteroscedastic gives the weighted least squares data model. It is more complex because it has extra parameters to account for multiple variances for multiple data. Whenever one or combinations of the three assumptions (independence, homoscedasticity, and normality) are relaxed, the resulting data models become more complex and require more parameters. Such systematic evaluation of data models (McInerney et al., 2017;

Smith et al. 2010b, 2015) is necessary to evaluate appropriateness of residuals assumptions and their impacts on Bayesian inference.

The assumptions of heteroscedastic, correlated, and non-Gaussian residuals are accounted for by using the method of Schoups and Vrugt (2010) in the following procedure: (i) the correlation is removed from the residuals by using an autoregressive model; (ii) the resulting residuals are normalized by a linear model of variance; and (iii) the normalized residuals are characterized by using the skew exponential power distribution. The data model parameters (i.e., coefficients of the autoregressive model, the linear variance model, and the skew exponential power distribution) are not specified by users, but estimated together with soil respiration model parameters during the Bayesian inference. The skew exponential power distribution is general in that by adjusting the values of its kurtosis and skewness parameters the distribution can produce other distributions such as the Laplace distribution (van Wijk et al., 2008; Ricciuto et al., 2011) and other distributions through using an exponential model with different kurtosis parameters (Tang and Zhuang, 2009). It is worth pointing out that there exist other methods to account for the three assumptions. Evin et al. (2013) suggested accounting for residual heteroscedasticity before accounting for residual autocorrelation. Lu et al. (2013) developed an iterative two-stage procedure to separately estimate physical model parameters and data model parameters. Evin et al. (2014) developed a similar procedure to first estimate model parameters and then estimate heteroscedasticity and autocorrelation parameters. While this study uses the method of Schoups and Vrugt (2010), exploring other methods is warranted in future studies.

After investigating the impacts of the data models on Bayesian inference, this study evaluates the impacts of the data models on predictive performance of the three soil respiration models. Using random samples generated during the Bayesian inference, a prediction ensemble is produced

for each soil respiration model. The ensemble is used to evaluate predictive performance of the models in a stochastic sense by estimating to what extent the models can predict future events. The evaluation in this study is done in a cross-validation manner by splitting the dataset of CO₂ efflux into two parts for Bayesian inference and cross-validation, respectively. The evaluation of predictive performance is important because different data models may give different parameter distributions and accordingly different predictive performance. For example, the study of van Wijk et al. (2008) concluded that the choice of the residual function is crucial to achieve accurate model prediction and parameter estimation. Shi et al. (2014) showed that the posterior parameter distributions and predictive performance given by two data models (weighted least square and skew exponential power distribution after removing heteroscedasticity and autocorrelation) are dramatically different, and a definitive conclusion was drawn that one data model is better than the other. The evaluation of predictive analysis is conducted for the following two cases: (1) the prediction ensemble is generated by random samples of the soil respiration models only (i.e. credible interval), and (2) the prediction ensemble is generated by random samples of not only the soil respiration models but also the data models (i.e. predictive interval). The two cases lead to different conclusions about the predictive performance. It is expected that the evaluation of predictive performance conducted in this study can help select the most appropriate data model to achieve optimal model predictions.

The remainder of the paper is organized as follows. Section 2 starts with a description of the evolving data models and their corresponding likelihood functions used in Bayesian inference, followed by a brief summary of the three soil respiration models. The results of Bayesian inference are discussed in Section 3 and Section 4, addressing the data model implications on parameter

estimation and predictive performance, respectively. Section 5 summarizes the key findings and limitations of this study, and provides recommendations for approaching data model selection.

2 Methodology

This section starts with a description of the eight data models that account for the three pairs of assumptions about residuals in a stepwise manner in Section 2.1. The data models are used to build the likelihood functions used in Section 2.2 for Bayesian inference. The three soil respiration models and observations of CO₂ efflux are described in Sections 2.3 and 2.4, respectively. Metrics for evaluating predictive performance are presented in Section 2.5.

2.1 Data models

This study considers eight evolving data models starting from a data model that assumes independent, homoscedastic, and Gaussian residuals to a data model that relaxes all the three assumptions. The eight data models are based on the generic normalized residual,

$$a_t = \frac{\varepsilon_t}{\sigma_t} \quad a_t \sim X, \quad (1)$$

where $\varepsilon_t = d_t - Y_t$ is the residual (the difference between data d_t and its corresponding model simulation Y_t) at time or location t ; σ_t is the standard deviation of the residual; and X is the probability density function (PDF) of a_t . The eight data models are formulated with different forms of ε_t , σ_t , and X . The standard least square (SLS) data model is

$$a_t = \frac{\varepsilon_t}{\sigma_0} \quad a_t \sim N(0,1), \quad (2)$$

where $\sigma_t = \sigma_0$ is a constant for all the data (i.e., homoscedasticity), and X is the standard normal distribution, $N(0,1)$. The unknown parameter σ_0 is estimated jointly with unknown physical

251 model parameters. If σ_t is not a constant (i.e., heteroscedastic), SLS becomes the weighted least
 252 squared (WLS) data model. While heteroscedasticity can be accounted for through residuals
 253 transformation (e.g. Thiemann et al., 200; Smith et al., 2010b) or other similar approaches (Gagne
 254 et al., 2015) a linear heteroscedastic model $\sigma_t = \sigma_0 + \sigma_1 Y_t$ is assumed here by following ~~other~~the
 255 studies of (Thyer et al., (2009)), Schoups and Vrugt, (2010), and Evin et al., (2013, 2014). With
 256 the linear model, there is no need to estimate σ_t for each data. Instead, σ_t is calculated by
 257 estimating only two parameters, σ_0 and σ_1 . The WSL data model is written as

$$258 \quad a_t = \frac{\varepsilon_t}{\sigma_0 + \sigma_1 Y_t} \quad a_t \sim N(0,1). \quad (3)$$

259 The two unknown parameters σ_0 and σ_1 are estimated jointly with unknown physical model
 260 parameters. The linear model assigns smaller weight to the data with larger simulation, Y_t . If the
 261 simulation is small and $\sigma_0 \gg \sigma_1 Y_t$, the weight becomes constant for all data. Both SLS and WLS
 262 assume that a_t is independently and identically distributed.

263 It is not uncommon that residuals are correlated in space and time, due to propagation of
 264 measurement errors (Tiedeman and Green, 2013) and model structure errors (Evin et al., 2014;
 265 Kavetski et al., 2013; Lu et al., 2013). The temporal correlation that occurs in the numerical
 266 example of this study can be accounted for by using a p -order autoregressive model. This leads to
 267 the data model of standard least square with autocorrelation (SLS-AC),

$$268 \quad a_t = \frac{\varepsilon_t - \sum_{i=1}^p \phi_i \varepsilon_{t-i}}{\sigma_0} \quad a_t \sim N(0,1) \quad (4)$$

269 where p is the order of autocorrelation, and ϕ_i is an autocorrelation coefficient. The unknown ϕ_i
 270 and σ_0 are estimated together with unknown model parameters. By extending the concept of
 271 correlated residuals to WLS leads to the weight least square with autocorrelation (WLS-AC),

$$272 \quad a_t = \frac{\varepsilon_t - \sum_{i=1}^p \phi_i \varepsilon_{t-i}}{\sigma_0 + \sigma_1 Y_t} \quad a_t \sim N(0,1) \quad (5)$$

273 The unknown parameters of σ_0 , σ_1 , and ϕ_i are estimated jointly with physical model
 274 parameters. Equations (2) – (5) assume that the residuals are Gaussian.

275 The next four data models are similar to the previous four models except that the standard
 276 normal distribution of a_t is replaced by the skew exponential power distribution, $SEP(0,1,\xi,\beta)$,
 277 with zero mean and unit standard deviation (Schoups and Vrugt, 2010)

$$278 \quad p(a_t | \xi, \beta) = \frac{2\sigma_\xi}{\xi + \xi^{-1}} \omega_\beta \exp \left[-c_\beta |a_{\xi,t}|^{2/(1+\beta)} \right], \quad (6)$$

279 where ~~zero is mean, one is standard deviation,~~ ξ is skewness, β is kurtosis,

$$280 \quad a_{\xi,t} = (\mu_\xi + \sigma_\xi a_t) / \xi^{\text{sign}(\mu_\xi + \sigma_\xi a_t)}, \quad \mu_\xi = M(\xi - \xi^{-1}), \quad \omega_\beta = \frac{\Gamma^{1/2}[3(1+\beta)/2]}{(1+\beta)\Gamma^{3/2}[(1+\beta)/2]},$$

$$281 \quad \sigma_\xi = \sqrt{(1 - M^2)(\xi^2 + \xi^{-2}) + 2M^2 - 1}, \quad M = \frac{\Gamma[1+\beta]}{\Gamma^{1/2}[3(1+\beta)/2]\Gamma^{1/2}[(1+\beta)/2]}, \quad \text{and}$$

$$282 \quad c_\beta = \left(\frac{\Gamma[3(1+\beta)/2]}{\Gamma[(1+\beta)/2]} \right)^{1/(1+\beta)} \text{ are derived variables of } \beta \text{ and } \xi, \text{ and } \Gamma[\cdot] \text{ is the gamma function. The}$$

283 kurtosis parameter $\{\beta \in \mathbb{R} : -1 \leq \beta \leq 1\}$ determines the peakness of the pdf such that the β values
 284 of -1, 0, and 1 give uniform, Gaussian and Laplace distributions, respectively. The skewness
 285 parameter $\{\xi \in \mathbb{R} : 0.1 \leq \xi \leq 10\}$ determines the skewness of the pdf such that the ξ values of 0.1,
 286 1, and 10 give positively skewed, symmetric, and negatively skewed distributions, respectively.

287 Setting $\beta=0$ and $\xi=1$ leads to $\mu_\xi=0$, $\sigma_\xi=1$, $\omega_\beta=1/\sqrt{2\pi}$, $c_\beta=1/2$ and $a_{\xi,t}=a_t$, and the
 288 skew exponential power distribution $SEP(0,1,\xi=1,\beta=0)$ becomes the standard normal
 289 distribution,

$$290 \quad p(a_t | \xi=1, \beta=0) = \frac{1}{\sqrt{2\pi}} \exp\left[-\frac{1}{2}(a_t)^2\right]. \quad (7)$$

291 which is the data model of SLS in equation (2).

292 Replacing $a_t \sim N(0,1)$ with $a_t \sim SEP(0,1,\xi,\beta)$ in equations (2) – (5) leads to the data models
 293 SEP, WSEP, SEP-AC, and WSEP-AC as follows,

$$294 \quad a_t = \frac{\varepsilon_t}{\sigma_0} \quad a_t \sim SEP(0,1,\xi,\beta) \quad (8)$$

$$295 \quad a_t = \frac{\varepsilon_t}{\sigma_0 + \sigma_1 Y_t} \quad a_t \sim SEP(0,1,\xi,\beta). \quad (9)$$

$$296 \quad a_t = \frac{\varepsilon_t - \sum_{i=1}^p \phi_i \varepsilon_{t-i}}{\sigma_0} \quad a_t \sim SEP(0,1,\xi,\beta) \quad (10)$$

$$297 \quad a_t = \frac{\varepsilon_t - \sum_{i=1}^p \phi_i \varepsilon_{t-i}}{\sigma_0 + \sigma_1 Y_t} \quad a_t \sim SEP(0,1,\xi,\beta) \quad (11)$$

298 In comparison with the Gaussian data models, the SEP-based data models have two more
 299 parameters (ξ and β) to be estimated jointly with physical model parameters. ~~WSEP-AC-d~~Data
 300 model WSEP-AC, which is known as the generalized likelihood function, is the most commonly
 301 used SEP-based data model (e.g. Vrugt and Ter Braak, 2011; Hublart et al., 2016; Scholz et al.,
 302 2018). A summary table of the eight data models with corresponding parameters is provided in the
 303 supplementary materials.

2.2 Bayesian inference and likelihood functions

Consider a Bayesian inference problem for a nonlinear model, f , used to simulate state variables (e.g., CO₂ efflux), $\mathbf{d} = \mathbf{Y}(\boldsymbol{\theta}) + \mathbf{e}\boldsymbol{\xi}$, where \mathbf{d} is a vector of data, $\boldsymbol{\theta}$ is a vector of model parameters, and $\mathbf{e}\boldsymbol{\xi}$ is a vector of residuals that may include errors in data, model parameters, and model structures. The goal of Bayesian inference is to estimate the posterior distributions, $p(\boldsymbol{\theta}|\mathbf{d})$, of model parameters, $\boldsymbol{\theta}$, given data, \mathbf{d} , using Bayes' theorem (Box and Tiao, 1992)

$$p(\boldsymbol{\theta}|\mathbf{d}) = \frac{p(\mathbf{d}|\boldsymbol{\theta})p(\boldsymbol{\theta})}{\int p(\mathbf{d}|\boldsymbol{\theta})p(\boldsymbol{\theta})d\boldsymbol{\theta}} \quad (12)$$

where $p(\boldsymbol{\theta})$ is the prior distribution, and $p(\mathbf{d}|\boldsymbol{\theta})$ is the likelihood function to measure goodness-of-fit between model simulations, $\mathbf{Y}(\boldsymbol{\theta})$, and data, \mathbf{d} . The prior distribution can be obtained from data of previous studies (e.g. Elshall and Tsai, 2014) or expert judgment. When prior information is lacking, a common practice is to assume uniform distributions with relatively large parameter ranges so that the prior distributions do not affect the estimation of posterior distributions.

The data models above can be used to construct the likelihood functions. For the Gaussian data models given in equations (2) – (5), the corresponding Gaussian likelihood functions are straightforward, and an example is equation (7). For the SEP data models, the corresponding likelihood that is called generalized likelihood function is (Schoups and Vrugt, 2010)

$$p(\mathbf{d}|\boldsymbol{\theta}) = p(\boldsymbol{\varepsilon}_t|\boldsymbol{\theta}) = \prod_{t=1}^n \sigma_t^{-1} \frac{2\sigma_\xi}{\xi + \xi^{-1}} \omega_\beta \exp\left(-c_\beta |a_{\xi,t}|^{2/(1+\beta)}\right). \quad (13)$$

where n is the dimension of \mathbf{d} . The Gaussian likelihood functions are special case of the generalized likelihood functions. For example, by setting $\beta=0$, $\xi=1$, $\phi_t=0$, $\sigma_t=\sigma_0$, $\sigma_\xi=1$, $\mu_\xi=0$,

$\omega_\beta=1/\sqrt{2\pi}$, $c_\beta=1/2$, and $a_{\xi,t}=a_t$, equation (13) becomes the likelihood function corresponding

to the SLS data model. Replacing $\sigma_t = \sigma_0$ by $\sigma_t = \sigma_0 + \sigma_1 E_t$, equation (13) becomes the likelihood function of the WLS data model.

In this study, the posterior distributions of the data model parameters are jointly estimated with the soil respiration model parameters using the MT-DREAM_(ZS) code (Laloy and Vrugt, 2012). MT-DREAM_(ZS) implements a Markov chain Monte Carlo (MCMC) algorithm by running multiple Markov chains in parallel with adaptive proposal distribution, multiple-try sampling, and sampling from an archive of past states. These state-of-the-art features assist in overcoming common challenges in the sampling [landscape-space](#) such as multimodality, ill-conditioning, and high dimensionality, and thus allow for accurate exploration of the targeted distributions.

2.3 Soil respiration models

Zhang et al. (2014) studied the Birch effect (the peak soil microbial respiration pulses in response to episodic rainfall pulses), and developed five models, evolving from an existing four-carbon pool model to models with additional carbon pools and/or explicit representations of soil moisture controls on carbon degradation and microbial uptake rates. Three of the five models are used in this study, and they are denoted as 4C, 5C, and 6C. Note that model 4C is model 4C_NOSM of Zhang et al. (2014), not [their](#) model 4C. Figure 1 is the diagram of model 6C, the most complex one among the five models. The simplest one, model 4C, has four carbon pools, i.e., soil organic carbon (SOC), dissolved organic carbon (DOC), microbial biomass (MIC), and enzymes (ENZ), and does not consider the soil moisture control on carbon degradation and microbial uptake rates. Models 5C and 6C [have](#) an explicit representation of soil moisture controls on the rates. Based on the dual Arrhenius and Michaelis–Menten kinetics model, the original SOC degradation rate, V_{decom} , is (Davidson et al., 2011; Davidson and Janssens, 2006)

$$V_{decom} = V_{\max} C_{ENZ} \frac{C_{SOC}}{K_m + C_{SOC}} \quad (14)$$

where V_{\max} [s^{-1}] is the maximum SOC degradation rate per unit enzyme when the substrates is not limiting, C_{ENZ} [gCm^{-3}] is enzyme pool size, C_{SOC} [gCm^{-3}] is SOC pool size, and K_m is the half-saturation for SOC. The original microbial uptake rate, V_{uptake} , is (Davidson et al., 2011; Davidson and Janssens, 2006)

$$V_{uptake} = V_{\max_up} C_{MIC} \frac{C_{DOC}}{K_{m_up} + C_{DOC}} \frac{C_{O_2}}{K_{m_upO_2} + C_{O_2}}, \quad (15)$$

where V_{\max_up} [s^{-1}] is the maximum DOC uptake rate when the substrates is not limiting, C_{MIC} [gCm^{-3}] is the **microbial biomass**MIC pool size, C_{DOC} [gCm^{-3}] is the DOC pool size, C_{O_2} [m^3m^{-3}] is the gas concentration of O_2 in the soil pore, and K_{m_up} [gCm^{-3}] and $K_{m_upO_2}$ [m^3m^{-3}] are the corresponding half-saturation constants for DOC and O_2 , respectively. With the explicit representation of soil moisture control, the two rates become (Zhang et al., 2014)

$$V_{decom} = V_{\max} C_{ENZ} \frac{C_{SOC}}{K_m + C_{SOC}} \left(\frac{\theta}{\theta_s} \right) \quad (16)$$

$$V_{uptake} = V_{\max_up} C_{MIC} \frac{C_{DOC}}{K_{m_up} + C_{DOC}} \frac{C_{O_2}}{K_{m_upO_2} + C_{O_2}} \left(\frac{\theta}{\theta_s} \right) \quad (17)$$

where θ [-] is the volumetric soil moisture, and θ_s [-] is the porosity.

In addition to using the new rate equations, models 5C and 6C have more carbon pools. In model 5C, DOC is split into two sub-pools for wet zone and dry zone of soil pores, and only the wet DOC is used by MIC, as shown in Figure 1. The moisture-controlled microbial uptake rate becomes

$$V_{uptake} = V_{\max_up} C_{MIC} \frac{C_{DOC_w}}{K_{m_up} + C_{DOC_w}} \frac{C_{O2}}{K_{m_upO2} + C_{O2}} \left(\frac{\theta}{\theta_s} \right). \quad (18)$$

where C_{DOC_w} [gCm⁻³] is the DOC pool size in the wet soil pores. Model 6C is more complex in that ENZ is further split into two sub-pools for wet and dry pores, and both the wet and dry ENZ are subject to degradation, as shown in Figure 1. The moisture-controlled SOC degradation rate becomes

$$V_{decom} = V_{\max} C_{ENZ_W} \frac{C_{SOC}}{K_m + C_{SOC}} \left(\frac{\theta}{\theta_s} \right) \quad (19)$$

for the wet ENZ and

$$V_{decom} = V_{\max} C_{ENZ_D} \frac{C_{SOC}}{K_m + C_{SOC}} \left(1 - \frac{\theta}{\theta_s} \right) \varepsilon_D \quad (20)$$

for the dry ENZ, where C_{ENZ_W} [gCm⁻³] is the wet soil pores enzyme pool size, C_{ENZ_D} [gCm⁻³] is the enzyme pool size in the dry soil pores, and ε_D is the catalysis efficiency of the dry zone enzyme.

Due to considering the moisture control and adding more soil pools, model 5C is expected to be significantly better than model 4C for simulating the Birch effect. Since the accumulated ENZ in dry soil is secondary, model 6C is expected to be slightly better than model 5C. In terms of model structural error, model 4C has the largest model structure error, model 5C has significantly less model structure error, and model 6C has the smallest model structural error. [In other words, model 6C has the highest model fidelity \(i.e. lowest model discrepancy\) among the three models.](#)

As shown below, the degree of model structural error is reflected in the process of Bayesian inference and verified by the cross-validation.

2.4 Observations and parameter estimation

Figure 2 plots the time series of 17,016 observations of soil ~~moister~~moisture and CO₂ efflux used in this study. The observations were obtained during the entire year of 2007, covering a long period of dry season prior to monsoon and episodic rainfall events during monsoon. The first two third of this dataset is used for the Bayesian inference, and the last one third is used for cross-validation. The inference and cross-validation periods have both dry and wet periods, as shown in Figure 2. The observation site is located within the Santa Rita Experimental Range (SRER, 31.8214°N, 110.8661°W, elevation 1,116 m) outside of Tucson, Arizona (Barron-Gafford et al., 2011; Scott et al., 2009). This savanna site was covered by 22% of perennial grass, forbs and subshrubs and 35% of mesquite. The soils are uniformly Comoro loamy sand (77.6% sand, 11.0% clay, and 11.4% silt). The half-hourly atmospheric forcing data were collected from measurements through an eddy covariance tower (Scott et al., 2009). This includes downward shortwave, longwave, precipitation, wind, air temperature, humidity, and pressure. Volumetric CO₂ concentration was measured at a half-hourly interval through compact probes. The CO₂ efflux was estimated from the gradient of CO₂ concentration measured at two depths of 2 cm and 10 cm through Fick's first law of diffusion, and the estimates were validated against measurements from a portable CO₂ gas analyzer.

The parameters estimated in this study include the parameters of the soil respiration models (4C – 6C) and the parameters of the data models described in Section 2.1. The estimated parameters of models 4C and 5C include the microbial carbon use efficiency (CUE) [g/g], enzyme production rate, k_e [g/m³s], microbial turnover rate, τ_m [1/s], and enzyme turnover rate τ_e [1/s]. Uniform distributions are used as the prior in the Bayesian inference, and the ranges of the four parameters are 0.2 – 1.00, 1×10^{-12} – 1×10^{-7} , 1×10^{-12} – 1×10^{-5} and 1×10^{-11} – 1×10^{-6} , respectively.

The values of other parameters are fixed at the values used in Allison et al. (2010). Model 6C has two more parameters, and they are the catalysis efficiency ε_D [-] and the turnover rate of the dry-zone enzymes τ_{en} [1/s]. The prior of the two parameters are uniform distributions with the ranges of 0.2 – 0.8 and $1 \times 10^{-12} - 1 \times 10^{-8}$, respectively.

The DREAM-based MCMC simulation is conducted for a total of 24 cases, the combinations of eight data models and three soil respiration models. For each case, the parameter distributions are obtained after drawing a total of 5×10^5 samples using five Markov chains. The Gelman and Rubin (1992) R-statistic is used for convergence diagnostic, and it approaches one in less than ~~40,000~~ $\times 10^4$ samples. The initial 50% of the samples are discarded during the burn-in period.

2.5 Metrics for evaluating predictive performance

Three criteria are used to evaluate the predictive performance of the soil respiration models and data models, and they are central mean tendency, dispersion, and reliability. Each criterion is measured by a single metric. In addition, a newly defined metric [by \(Elshall et al., 2018\)](#) is also used for simultaneously measuring the three criteria.

The central mean tendency is measured in this study using the Nash-Sutcliffe model efficiency (NSME) coefficient (Nash and Sutcliffe, 1970),

$$NSME = 1 - \frac{\sum_{i=1}^n (d_i - \bar{Y}_i)^2}{\sum_{i=1}^n (d_i - \bar{\mathbf{d}})^2}, \quad (21)$$

where n is the number of cross-validation data, d_i is the i -th data, $\bar{\mathbf{d}}$ is the mean of the data, and \bar{Y}_i is the mean of the prediction ensemble, Y_i , for d_i . NSME ranges from $-\infty$ to 1, with $NSME = 1$ corresponding to a perfect match between data and mean prediction, i.e., the ensemble is centered on the data. $NSME = 0$ indicates that the model predictions are as only accurate as the mean of the

data, while an efficiency $NSME < 1$ indicates that the mean of data is a better prediction than the mean prediction.

In addition to the central mean tendency, it is also desirable that the ensemble is precise with small dispersion and reliable to cover all the data. This study uses a nonparametric metric for dispersion, and it is the sharpness of a prediction interval (e.g. Smith et al., 2010a)

$$Sharpness = 1/n \sum_{i=1}^n [Max(\mathbf{Y}_i) - Min(\mathbf{Y}_i)] \quad (22)$$

where \mathbf{Y}_i is the prediction ensemble within the 95% prediction interval, (the Bayesian credible interval, not the confidence interval used in nonlinear regression (Lu et al., 2013). Smaller values of sharpness indicate better prediction precision. Reliability is measured using predictive coverage. (e.g. Hoeting et al., 1999), which is the percentages of data contained in the prediction interval. Larger predictive coverage values are preferred.

To account for the trade-off between the three metrics, (Elshall et al., (2018b)) -defined relative model score (RMS) that simultaneously measure all the three criteria. Scoring rules are commonly used in hydrology to assess predictive performance (e.g., Weijs et al., 2010; Westerberg et al., 2011). RMS is used in this study to measure the relative predictive performance of the combinations of soil respiration models and data models. For combination M_j , RMS is defined as

$$RMS(M_j) = \frac{\sum_{i=1}^n p(d_i | \mathbf{Y}_{ij}, M_j)}{\sum_{j=1}^m p(d_i | \mathbf{Y}_{ij}, M_j)} \times 100 \quad (23)$$

where m is the number of combinations; ~~and~~ the ensemble prediction \mathbf{Y}_{ij} is similar to \mathbf{Y}_i above ~~where is with~~ i index i over time and index j specific to the j -th combination. The density function $p(d_i | \mathbf{Y}_{ij})$ can be evaluated by first obtaining the density function $p(\mathbf{Y}_{ij})$ of the ensemble prediction \mathbf{Y}_{ij} (e.g., by using the kernel density function) and then evaluating $p(d_i | \mathbf{Y}_{ij})$ using interpolation

methods based on the intersection of Y_{ij} and d_i . [More details about evaluating RMS can be found in Elshall et al. \(2018\).](#) This evaluation is based purely on the model predictions, and does not involve any assumptions on the models, their parameters, and likelihood functions. Larger RMS values indicate better overall predictive performance. A figure of our workflow scheme is presented in the supplementary materials.

3 Results of Bayesian Inverse Modeling

This section analyzes the residuals of the best realization (with the highest likelihood value) of the MCMC simulation to understand whether the assumptions of the eight data models hold. The impacts of the data models on the posterior parameter distributions are also analyzed.

3.1 Residual characterization

Figure 3 shows residual plots for model 6C based on data models SLS and WSEP-AC. SLS is the simplest [data model](#) with the assumptions of homoscedastic, independent, and Gaussian residuals, and the WSEP-AC is the most complex one without the assumptions. Model 6C is the most complex model and also the best one as ranked by Zhang et al. (2014) using Bayesian model selection. The variable a_t plotted in Figures 3a-3c and Figures 3d-3f is defined in equations (2) and (11), respectively. Figures 3a – 3c show that [all](#) the three residual assumptions are violated when SLS is used, because (i) the residual variance is not constant, but increases as a function of the simulated CO₂ efflux (Figure 3a); (ii) the autocorrelation function at most lags is beyond the 95% confidence interval (Figure 3b); (iii) ~~and~~ the standard normal density function cannot adequately characterize the residuals (Figure 3c). Figures 3d-f show that, after relaxing the three assumptions, the processed residuals, a_t , can be well characterized by WSEP-AC. Figure 3d shows that, after normalizing ε_t with the linear variance ($\sigma_t = 0.034 + 0.099 E_t$), the variation of the variance of a_t becomes significantly smaller, although the variance is still not constant. Figure 3e shows that,

after removing a first-order autoregressive model from ε_t , a_t becomes less correlated, although the correlation is not fully removed. The two coefficients of the autoregressive model are $\phi_1 = 0.989$ and $\phi_2 = 4.5 \times 10^{-6}$; the small value of ϕ_2 indicates that there is no need to attempt an autoregressive model of higher order. Figure 3f shows that a_t follows the SEP distribution with the estimated skewness coefficient of $\xi = 0.933$ and kurtosis coefficient of $\beta = 0.998$. As a summary, Figure 3 shows that it is important to examine the residuals and to determine whether the selected data model is adequate for charactering the residuals. Although WSEP-AC still cannot perfectly characterize ε_t , it is significantly better than SLS.

Although the Gaussian assumption used in SLS is violated for model 64C (Figure 3c), this is not generally the case for other data models and soil respiration models. This is shown in Figure 4, which presents the quantile-quantile (Q-Q) plot for the eight data models and the three soil respiration models. For SLS, WLS, SLS-AC, and WLS-AC, the theoretical quantiles are based on the standard normal distribution, $N(0,1)$; for SEP, WSEP, SEP-AC, and WSEP-AC, the theoretical quantiles are based on the standard skew exponential power distribution, $SEP(0,1,1,0)$. If the residuals follow the assumed standard distributions, the Q-Q plots fall on the 1:1 lines, which is marked as the theoretical lines in Figure 4. If the residuals are Gaussian or SEP but not standard, the Q-Q plots fall on a straight line but not the 1:1 line. Figures 4a and 4e show that, for all the soil respiration models, the Q-Q plots of SLS and SEP deviate significantly from the theoretical lines and exhibit fat-tail behaviors, which is an indication of outliers (Thyer et al., 2009). The deviation is reduced after accounting for autocorrelation in SLS-AC and SEP-AC, as shown in Figures 4c and 4g. It is interesting to observe from the two figures that the Q-Q plots of the three models are almost visually identical. The deviation is almost fully removed after accounting for heteroscedasticity in WLS and WSEP in that their corresponding Q-Q plots fall on the 1:1 lines,

especially for models 5C and 6C, as shown in Figures 4b and 4f. However, the Q-Q plots start deviating from the 1:1 lines as shown in Figures 4d and 4h, after accounting for both heteroscedasticity and autocorrelation in WLS-AC and WSEP-AC. As a summary, Figure 4 shows that, for the numerical example of this study, either the Gaussian or the SEP distribution is valid if heteroscedasticity is accounted for in the data models. However, accounting for autocorrelation in the data models does not help improve the characterization of the residual distributions.

3.2 Posterior parameter distributions

While Figures 3 and 4 help understand validity of the three assumptions used in the data models, the impacts of the data models on estimating model parameter distributions must be evaluated separately. This section discusses the impact of the data model selection on parameter estimation with the objective of understanding whether incorrect specification of the data model, will necessarily lead to biased parameter estimates. Such assessment is not a trivial task for two main reasons. First, microbial soil respiration models aggregate complex natural processes and spatial details into simpler conceptual representations. As a result several model parameters are effective values of several complex natural processes that cannot be actually measured in the field as discussed by Vrugt et al. (2013). Second In addition, even for model parameter that can be measured in the field, since the model structure is imperfect, calibrated it can be the case that parameter values are sometimes can be accepted beyond their physically reasonable range, as discussed by Pappenberg and Beven (2006). This is often undesirable, if we seek to make the models more mechanistically descriptive.

We focus our discussion on carbon use efficiency (CUE) for microbial growth due to two reasons: (1) since CUE is a fundamental parameter in microbial soil respiration models, and (2) a physically reasonable physical range for CUE can be estimated. The concept of microbial

CUE(Allison et al., 2010; Bradford et al., 2008; Manzoni et al., 2012; Wieder et al., 2013) has been used to present fundamental microbial processes in recent microbial enzyme models (Allison et al., 2010; German et al., 2011; Schimel and Weintraub, 2003; Wang et al., 2013). The microbial CUE, which is marked between MIC and CO₂ in Figure 1, controls microbial growth, enzyme production and microbial respiration. A physically reasonable range of CUE can be estimated from the physical viewpoint (Tang and Riley, 2014). Sinsabaugh et al. (2013) ~~study show~~eds that the thermodynamic calculations support a maximum CUE of 0.60 and that ~~methods used to previous~~ studies that estimate CUE in terrestrial systems report a mean value of 0.55. Theoretically, there ~~is~~ no lower limit for CUE as it can approach zero, and CUE < 0.1 ~~has been~~are reported for terrestrial ecosystems (e.g., Fernández-Martínez et al., 2014) and used in modeling studies (Li et al., 2014). Note that, for inverse modeling with MCMC sampling, we did not assume CUE maximum value of 0.6. In other words, for parameter estimation and predictive performance we did not impose the constraint that CUE is less than 0.6. We merely use this CUE maximum value of 0.6 to evaluate whether the posterior CUE parameter samples obtained using different data models and different soil respiration models are within the physically reasonable range of 0 ~ 0.6.

Figure 5 plots the CUE posterior marginal density of the three soil respiration models obtained using the eight data models. The physical range between zero and 0.6 is marked in yellow. Figure 5 shows that the CUE posterior parameter distribution ~~offer~~ Model 6C obtained using ~~for all the~~ data models likelihood functions that does not account for autocorrelation are within the physically reasonable ~~physical~~ range. For models 4C and 5C, the posterior parameter samples are outside the ~~physical~~ range for six data models. For model 4C, the posterior parameters are within the physical range only for data models SEP and WSEP; for model 5C, the two data models are WLS and WSEP. It is not surprising to find the posterior parameter distribution of models 4C and

5C, which have a certain degree of model structure error, to be out of the physically plausible ~~physical~~ range. This can be attributed to two reasons. First, the model solution can be biased toward the missing processes in the model structure such as the additional carbon pool in both 4C and 5C or missing the explicit accounting for soil ~~moister~~moisture in 4C. Second, biased parameter estimation can compensate for model structure inadequacy and other sources of discrepancy in both the physical models and the ~~statistical-data~~ models.

In addition, it is important to understand how accounting for autocorrelation, heteroscedasticity and non-Gaussian residuals can affect the parameter estimation. First, it is observed in Figure 5e-h that we obtained biased parameter estimates ~~that is are~~ outside the ~~reasonable~~ physically reasonable range when autocorrelation is explicitly accounted for ~~as shown in Figure 5e-h~~. This may suggest again that accounting for heteroscedasticity is desirable but accounting for autocorrelation is not. A possible reason is that filtering autocorrelation may reduce the residual space such that the transformed residual space cannot correspond to the parameter space of the models. In other words, parameter information may be lost due to filtering out autocorrelation. However, it is not fully understood why this does not occur for the model 6C under data model SLS-AC (Figure 5e), and more research is warranted. Second, unlike accounting for autocorrelation, accounting only for heteroscedasticity (i.e., WLS and WSEP) ~~since this will~~ only ~~amplifies~~ or ~~reduces~~ the variance without affecting the structure of the residual space. Figures 5c-d show ~~s~~ that account for heteroscedasticity (i.e. WLS and WSEP) tends to improve the parameter estimation in comparison with homoscedastic data models (i.e., SLS and SEP) shown in Figure 5a-b. Finally, with respect to non-Gaussian residuals, Schoups and Vrugt (2010) ~~proposes~~ suggested that, compared to Gaussian pdf, the peaked pdf of the SEP with ~~heavier-a longer~~ tails ~~compared to Gaussian pdf~~ is useful for making parameter inference robust against outliers. To a

certain degree, this can be substantiated by the results in Figure 5a-d, ~~insuch~~ that SEP and WSEP provide more favorable parameter estimates than SLS and WLS.

Finally, ~~from Figure 5a we can also notice shows~~ that the posterior parameter distributions of SLS ~~(Figure 5a) isare~~ very narrow for the three soil respiration models. These narrow ~~posterior parameter distributions of SLS compared to other likelihood functions~~ can be attributed to several reasons. Since SEP distribution can have heavier-longer tails than Gaussian distribution, this can further increase the samples acceptance ratio from tails resulting in wider distribution (Figure 5b). In addition, accounting for heteroscedasticity will result in wider ~~the~~-posterior parameter distribution (Figure 5c) due to accepting higher variances at peak effluxes. Moreover, filtering correlation (Figure 5e-h) increases the entropy, and leads to wider distributions.

4. Results of Predictive Performance

Based on the last one third of the CO₂ efflux observations, a cross-validation test was conducted for ~~all the 24 models~~, the combinations of three soil respiration models and eight data models. ~~Given For~~ the cross-validation period, the predictive performance is examined using the four statistical metrics that are defined in Section 2.5. The metrics are also calculated for the calibration period. This is not to perform Bayesian model selection given the calibration data, but to better understand the impact of data models on predictive performance of the three soil respiration models. For each calibration and each cross-validation data, a prediction ensemble is generated from the two perspectives of parametric uncertainty only and total uncertainty, as presented in Section 4.1 and 4.2, respectively.

4.1 Predictive performance with parametric uncertainty of soil respiration model

In this section the ensemble is generated by running the soil respiration models with the posterior samples (obtained from the Bayesian inference) of the physical model parameters. In other words, the ensemble addresses parametric uncertainty of the soil respiration models only. Considering the relative contribution of parametric uncertainty only will provide insights for modeling approaches that attempt to segregate various sources of uncertainty (e.g., Thyer et al., 2009 ; Tsai and Elshall, 2013).

The four statistics above (i.e. NSME, sharpness, coverage, and RMS) are calculated for the three soil respiration models and the eight data models. Taking data models SLS and WSEP-AC as an example, Figure 6 plots the data (for the calibration and cross-validation periods separately) along with the mean and 95% credible intervals of the prediction ensemble for the three models.

Figure 6 shows that the data models affect model simulations for all the models. The statistics, especially RMS, indicate that WSEP-AC has better predictive performance than SLS. This is most visually obvious for model 6C during the cross-validation period after 330 days, as the prediction ensemble of SLS (Figure 6k) cannot cover the observations, ~~whereas~~~~unlike~~ the prediction ensemble of WSEP-AC can (Figure 6l). This conclusion that WSEP-AC outperforms SLS agrees with that drawn from Figures 3 and 4.

Figure 7 plots the four statistics for all the soil respiration models and data models. Figures 7a and 7b show the predictive performance with respect to the central mean tendency measured by ~~using~~ NSME for both the calibration and cross-validation periods respectively. The results indicate that, under all data models, the low fidelity model 4C ~~under all data models will~~ over-fits the data ~~and~~ results~~ing~~ in biased predictions, ~~such in~~ that the NSME values become significantly worse (e.g., from 0.6 to -0.6) from the calibration to the cross-validation period. This is confirmed by the visual inspection of Figures 6a ~~and, 6b,~~ 6g for data model SLS; and of Figures 6b and 6h

for data models ~~s-SLS-and~~ WSEP-AC. For models 5C and 6C, their NSME values vary with the data models; and the central mean accuracy is the worst for SLS-AC that considers only autocorrelation (Figure 6b).

With respect to parametric uncertainty estimation, Figures 7c and 7d show that sharpness generally increases when the three assumptions in the data models are gradually relaxed from SLS to WSEP-AC. This is even more obvious during the validation period. Given that the prediction ensemble does not center on the data, the increasing sharpness is desirable as it improves reliability. This is confirmed by the reliability plots in Figures 7e and 7f. The exceptions are again for SLS-AC and SEP-AC that generally have the lowest coverage.

With respect to the overall predictive performance measured by RMS, the same variation pattern and exception are also observed in the RMS plots in Figures 7g and 7h. This is not surprising because RMS is the metric that can be used to measure all the three criteria (central mean tendency, sharpness, and reliability). Since the prediction ensemble is not centered on the data, the sharpness and reliability are the decisive factors for evaluating the predictive performance.

As a summary, while it is necessary to account for heteroscedasticity in a data model, caution is needed when accounting for autocorrelation in the manner described in Section 2.1. In addition, after comparing the RMS values of the residuals using the Gaussian and SEP distributions, the conclusion is that the SEP distribution outperforms the Gaussian distribution with respect to predictive performance. Finally, uncertainty underestimation ~~is as evidenced~~ is not unexpected by the very small predictive coverage. The underestimation of uncertainty for all the physical models with all ~~likelihood-functionsthe data model makes sense~~ is not unexpected because only parametric

uncertainty is considered in this study. Considering the overall predictive uncertainty is the subject of the next section.

4.2 Predictive performance with total uncertainty

The simulated output $Y(\theta_p)$ ~~is will~~ generally not ~~be~~ equal to the observed output \mathbf{d} , and we have a residual term $\underline{\varepsilon}$ ~~e~~ due to measurement, input and model structure errors such that $\mathbf{d} = Y(\theta_p) + \varepsilon$. Accounting for the error term ~~e~~ ε can be through separating various error terms. For example, in section 4.1 we obtained uncertainty due to the physical model parameters. Accounting for other sources of uncertainty can be done using a single model approach (e.g. Thyer et al., 2009) or a multi-model approach (e.g. Tsai and Elshall, 2013). Alternatively, we can quantify the uncertainty based on total residuals that separates out parametric uncertainty, so the residual error includes s errors ins measurements, model inputs, and model structures s-uncertainty (e.g. Thyer et al., 2009; Schoups and Vrugt, 2010). This lumped approach is based on sampling the residuals model $\varepsilon(\theta_\varepsilon)$ with parameters θ_ε . SLS has one fixed parameter that is the constant variance, and other data models have two to six parameters. Thus in this section the prediction ensemble addresses parametric uncertainty of not only the soil respiration models but also the data models. When generating the prediction ensemble in the procedure described by Schoups and Vrugt (2010), an ensemble of residuals is first generated by running the data models with posterior samples of the data model parameters for the positive carbon efflux domain; the residual ensemble is then added to the prediction ensemble generated in Section 4.1.

We start by ~~a~~the visual assessment of the predictive performance. Figure 8 is similar to Figure 6 with the exception that Figure 8 considers the overall ~~all~~ predictive uncertainty (i.e. parametric and output uncertainty), while Figure 6 considers the parametric uncertainty only. Figure 8 reveals a practical observation about accounting for the overall uncertainty through the lumped approach

of sampling the ~~data~~residuals models. ~~For example,~~ Figure 8b shows that, despite the wide prediction interval of model 4C, ~~which has the model with~~ significant model structure error, ~~it could~~ cannot capture the birch pulse around day 180. ~~This clearly~~It indicates that proper using a data model for model residuals cannot compensate modeling of the residuals will not make up for of significant model structure error.

Figure 9 plots the four statistics (NSME, sharpness, predictive coverage, and RMS) of the three soil respiration models under the eight data models to assess the predictive performance. ~~First~~ ~~w~~With respect to central mean tendency, ~~t~~The NSME values in Figures 9a-9b are visually the same as those in Figures 7a-7b, indicating that the central mean accuracy under parametric uncertainty is the same as that under predictive uncertainty.

With respect to uncertainty, the values of sharpness and predictive coverage increase substantially (Figures 9c – 9f). In particular, Figures 9e and 9f show that, except for SLS and SEP, the predictive coverage of the rest of the six data models are close to 100% for all the three soil respiration models, indicating that the prediction intervals cover almost all the data. This is demonstrated in Figures 6 for WSEP-AC. Similar to Figures 7c and 7d, Figures 9c and 9d also show a general pattern that the sharpness increases when the three assumptions in the data models are gradually relaxed from SLS to WSEP-AC. The data models that account for autocorrelation are still the exceptions.

With respect to the overall predictive performance, the RMS values are largely determined by the mean accuracy and sharpness as the predictive coverage is similar for different data models. Figures 9g and 9h of RMS show that the predictive performance of the four data models that account for autocorrelation is worse than that of the other four data models. This suggests again that one needs to be cautious when building autocorrelation into a data model. This is consistent

with the finding of Evin et al. (2013, 2014) that accounting for autocorrelation before accounting for heteroscedasticity or jointly accounting for autocorrelation and heteroscedasticity can result in poor predictive performance. In summary, Figures 9g and 9h show for both the calibration and prediction periods that accounting for heteroscedasticity (~~i.e. in~~ WLS and WSEP) ~~will~~ gives the best overall predictive performance, and accounting for autocorrelation without heteroscedasticity (~~i.e. in~~ SLS-AC and SEP-AC) ~~will~~ gives the worst overall predictive performance. Finally, for the three soil respiration models, RMS shows that model 4C has the worst predictive performance for both the calibration and cross-validation data. Generally speaking, the high fidelity model 6C outperforms model 5C for both the calibration and cross-validation data, which justifies the complexity of model 6C.

To demonstrate the impacts of the data models on predictive performance of the soil respiration models, Figure 10 plots the model simulations and predictions given by model 6C during the calibration and cross-validation periods using all the eight data models.

~~In Figure 10 is used to investigate we try to understand the~~ predictive performance characteristics of the different data models. ~~By looking at examining~~ the predictive performance of model 6C, ~~specific~~ predictive performance patterns can be identified. Figures 10-a ~~– 10d~~ show that SLS and SEP have similar predictive performance with SEP generally having better predictive performance especially during the validation period. ~~Accounting for heteroscedasticity using WLS as shown in Figures 10e and 10h~~ Not accounting for heteroscedasticity will underestimate the predication uncertainty (Figure 10b and Figure 10d). This is mainly because the variance of the efflux residuals increases with the magnitude of the carbon effluxes (Figure 3a), and thus assuming constant variance is not representative. Accordingly, accounting for heteroscedasticity using WLS (Figure 10e) or WSEP (Figure 10h) will make the predictions more

sensitive to peak carbon effluxes. This and will generally improve the predictive coverage on the expense of sharpness and the central mean tendency. While WLS and WSEP have similar predictive performance, ~~However,~~ WSEP ~~maintains~~ has ~~slightly~~ better central mean tendency and overall predictive performance than WLS. ~~Figures 10i – 10l show that Accounting for autocorrelation using SLS-AC and SEP-AC as shown in Figures 10i and 10l reduces the information content of the residuals, and thus resultings~~ in wider uncertainty bands and insensitivity to peak carbon effluxes as compared to SLS and SEP (Figures 10a-d), which may be due to reduction of information content of the residuals. This ~~resulteds~~ in deteriorating the sharpness, the central mean tendency and the capturing of peak carbon fluxes, especially during the validation period. ~~Figures 10m – 10p show that Accounting for both heteroscedasticity and autocorrelation using WLS-AC and WSEP-AC will makes~~ the inference robust against ~~peak~~ carbon effluxes. ~~However,~~ ~~yet~~ due to the loss of information content, the uncertainty bands are still wider, and uncertainty becomes overestimated especially during validation period as compared to WLS and WSEP (Figures 10e – 10h). The results of Models 4C and 5C, which are not shown here, also show the same prediction patterns with respect to non-Gaussian residuals, heteroscedasticity, and autocorrelation.

Finally, we observe in Figure 10 that the data models that have good overall predictive performance as measured by RMS during the calibration period will maintain this good predictive performance during the validation period. For model 6C, RMS values for the calibration and validation periods are very well correlated with a correlation coefficient of 0.92. However, we note that for models 4C and 5C the overall predictive performances during the calibration and validation periods are not that well correlated as 6C, with correlation coefficients of 0.52 for model 4C and

0.61 for model 5C. This suggests that model 6C is more robust than 4C and 5C for forecasting and hindcasting.

4.3 Discussion on handling residual correlation

Accounting for autocorrelation can lead to biased parameter estimation (Figure 5) and poor predictive performance (Figure 10). Auto-correlated residuals may be attributed to model discrepancy, as shown in Lu et al. (2013). The most obvious solution to handle the autocorrelation is to reduce the autocorrelation by improving the soil respiration model. If model improvement is difficult for practical reasons, we can improve the data model to better characterize the autocorrelation. Addressing autocorrelation in a data model is challenging since it involves several interlinked factors as follows:

- (1) Non-stationarity due to wet-dry periods could be a reason for this problem. By drawing on similarity from surface hydrology, the study of Ammann et al. (2018) suggests that auto-correlated residuals might be attributed to non-stationarity due to wet-dry periods with half-hourly data. Accounting for non-stationarity could better address the problem of auto-correlated residuals (Ammann et al., 2018; Smith et al., 2010b).
- (2) The way of implementing autocorrelation could have an impact. Autocorrelation could be applied to raw residuals directly (e.g., Li et al., 2015), to transformed residuals based on covariance matrix of residuals $L(\mathbf{e})$ (e.g., Lu et al., 2013), or to normalized residuals $L(\mathbf{a})$ (e.g., Schoups and Vrugt, 2010; Evin et al., 2013). Note that \mathbf{e} is a vector of transformed residuals, while \mathbf{a} denotes a vector of independent and identically distributed random errors with zero mean and unit standard deviation. The $L(\mathbf{e})$ approach based on covariance matrix of residuals is generally limited to Gaussian data models (e.g. Lu et al., 2013), while the $L(\mathbf{a})$ approach for normalized residuals can be readily adopted for non-Gaussian data models.

(3) The autocorrelation model could have an impact. Using an autoregressive model is a popular technique to account for auto-correlated residuals. However, using an autoregressive model with either joint inversion approach (e.g., this study and Schoups and Vrugt, 2010) or sequential approaches (e.g., Evin et al., 2013, 2014; Lu et al., 2013) removes correlation errors through a filter approach, which can lead to a loss of information content. As this may cause overcorrection of prediction especially at surge events, Li et al. (2015) developed a restricted autoregressive model to overcome this adverse effect. Other autocorrelation models include moving average model and mixed autoregressive-moving averaging model (Chatfield, 2004).

(4) Joint versus sequential inversion for autocorrelation could have an impact. Sequential inversion approaches include two-step procedures (e.g. Evin et al., 2013, 2014; Lu et al., 2013) or the multi-step procedure (Li et al., 2016a). These sequential approach estimates the autoregressive parameters sequentially in a later step after estimating the physical model parameters and other data model parameters. Evin et al. (2013, 2014) used a sequential approach to avoid the interaction between the parameters of the heteroscedasticity model and the autocorrelation model. In addition, the autoregressive model parameters can be deterministically calculated as an internal variables of the data model similar to Lu et al. (2013), and not as calibration parameters (e.g. Schoups and Vrugt; Evin et al. 2013; 2014). While the first step in the sequential approach would avoid the biased parameter estimation (Figure 10a-d), the second step can still lead a poor predicative performance since we are essentially using a filter approach to remove residual correlation. To address this problem, Li et al. (2016) multi-step procedure that is based on Gaussian data model uses restricted autoregressive model. Generally, Ammann et al. (2018) states that the joint inversion is still preferred, and

understanding the conditions where accounting for auto-correlation can be achieved remain poorly understood.

5. Conclusions

In parameter estimation and prediction of soil carbon fluxes to the atmosphere, ~~we~~one often assumes that residuals, which include errors in observations, model inputs, parameter estimates, and model structures~~errors~~, are normally distributed, homoscedastic and uncorrelated. We ~~studied~~ed these assumptions by calibrating three ~~microbial-enzymes~~soil inspiration models, which have varying degrees of model structure errors. We further explore ~~tested~~ eight data models that characterize the residuals statistically by starting with the standard least squares (SLS) and skew exponential power (SEP) data models that assume homoscedastic and non-correlated residuals. ~~Given-For~~ these two distributions, we evaluated ~~d~~ six other data models that account for heteroscedasticity (WLS and WSEP), autocorrelation (SLS-AC and SEP-AC), and joint inversion of heteroscedasticity and autocorrelation (WLS-AC and WSEP-AC). To our knowledge this is the first study that provides such detailed analysis for soil respiration inverse modeling. We also used ~~d~~ three ~~solid~~ed respiration models with different degrees of model fidelity (i.e., model ~~realism~~discrepancy) and model complexity (i.e. number of model parameters); to understand the impact of model discrepancy on the calibration results under different data models. We analyzed ~~d~~ the ~~calibration~~ results with respect to (1) residual characterization, (2) parameter estimation, (3) predictive performance, and (4) impacts of model discrepancy. The main findings of this study ~~can be~~are summarized as follows:

(1) With respect to residual characterization, residual analysis results suggest that the common assumption of not accounting for heteroscedasticity and residual autocorrelation ~~of residuals~~ (i.e. in the data models SLS and SEP) results in poor characterization of residuals. Explicit

accounting for heteroscedasticity ~~in(i.e. WLS and WSEP) can~~ results in ~~good~~significantly improved characterization of the residuals, and the improvement is larger than that obtained by, and is followed by joint the inversion of accounting for both heteroscedasticity and autocorrelation ~~(i.e.in WLS-AC and WSEP-AC)~~. Accounting for autocorrelation only ~~(i.e.in SLS-AC and SEP-AC)~~ ~~may not~~does not significantly improve ~~much~~ the characterization of the residuals.

(2ii) With respect to parameter estimation, the impacts of the data models are evaluated by we focused ~~ed~~ing on carbon use efficiency (CUE), which is a central parameter in soil respiration modeling. ~~We found the~~Using SLS yields with relatively reasonable posterior parameter distributions for CUE, yet very narrow posterior. ~~The D~~data models ~~consider autocorrelation (i.e. SLS-AC, SEP-AC, WLS-AC and WSEP-AC that consider autocorrelation)~~ tend to ~~generally~~ yield CUE estimates that are physically ~~non-un~~reasonable. We speculate that filtering residual correlation can affect the mapping of the model physics (as implicitly included in the residuals) into the ~~likelihood-parameter~~ space, which might result in biased parameter estimates that are physically unreasonable.

(3iii) With respect to predictive performance, it is measured by four statistical criteria: we assessed ~~the~~ central mean tendency, sharpness, coverage, and relative model score~~uncertainty bands and the overall predictive performance~~ for both the calibration and the cross-validation periods. Results show that accounting for autocorrelation ~~(i.e.in SLS-AC, SEP-AC, WLS-AC, and WSEP-AC)~~ deteriorates the predictive performance, such that the predictive performance is inferior to that of SLS in terms of the central mean tendency and overall predictive performance (measured by the relative model score), especially during the cross-validation period. Results also indicates that using ~~thea~~ SEP distribution can potentially improve the predictive

performance. The same is true for accounting for heteroscedasticity. Using SEP distribution and accounting for heteroscedasticity (i.e. WSEP) can potentially improve the predictive performance.

(4iv) With respect to the impact of model discrepancy, the high fidelity ~~complex~~ model (6C) gives the best results with respect to parameter estimation and predictive performance. Model 6C generally maintain~~ed~~s its superior performance under different data models. This justifies the complexity of model 6C relative to model 5C that has one less carbon pool. Model 4C ~~that has~~ with the lowest fidelity maintains its poor performance for different data models, because the model model-with has only four carbon pools and lacks the explicit representation of soil moisture control, ~~maintains its poor performance for different data models.~~

Based on ~~From~~ the empirical findings ~~above of this research,~~ we conclude the following:

(i1) Not accounting for heteroscedasticity and autocorrelation using a Gaussian or non-Gaussian data model might not necessarily result in biased parameter estimates or biased predictions with respect to central mean tendency, but will definitely underestimate uncertainty resulting in lower overall predictive performance.

(2ii) Using a non-Gaussian ~~residual error data~~ model can improve ~~the~~ parameter estimation~~s~~, and ~~the~~and predictive performance with respect to central mean tendency and uncertainty ~~estimation~~quantification.

(3iii) Accounting for heteroscedasticity ~~will definitely~~ improve~~s~~ the uncertainty estimation with respect to reliability at the cost of having a wider predictive interval.

(4iv) This study confirms ~~the other~~ empirical findings and theoretical analysis~~es~~ (Evin et al., 2013; 2014; Li et al., 2015, Ammann et al. 2018) that separately accounting for autocorrelation or jointly accounting for autocorrelation and heteroscedasticity can be problematic. While the

reasons remain poorly understood (Ammann et al., 2018), it might be attributed to non-stationarity due to wet-dry periods with half-hourly data (Ammann et al., 2018) or to the method of handling autocorrelation (e.g., Schoups and Vrugt, 2010, Evin et al., 2013; 2014; Lu et al., 2013; Li et al., 2015, 2016a; Ammann et al. 2018). Further investigation to address autocorrelation in soil respiration modeling is warranted in a future study. Accounting for non-stationarity (Smith et al., 2010b, Ammann et al. 2018) could address this problem. Relatively poor performance with respect to autocorrelation can be also attributed to the implementation scheme. The inference scheme such as joint inference as in this study, post-processing inference approach for autocorrelation (Evin et al., 2013; 2014), residuals transformation approach (e.g. Lu et al., 2013) or other strategies (Li et al., 2015, 2016a) could have an impact. Yet Ammann et al., (2018) study states that the joint inversion is still preferred, and understanding the conditions where accounting for auto-correlation can be achieved remain poorly understood. Further investigation of this point is warranted in a future study.

The above conclusions are subject to several limitations. First, the conclusions are specific to the soil respiration models developed and validated for semi-arid savannah. Performance variations across different soil respiration models with different levels of complexities is possible. Second, the conclusions are conditioned on the data that were obtained at the half-hour interval over a one-year period. Different conclusions are possible if the data are thinned to daily or weekly scales or data of longer observation periods are used. Third, ~~our~~the study investigates effects of the residual assumptions of formal likelihood functions through direct conditioning of the residuals model parameters, yet this can also be done through other approaches such as residuals transformation (Thiemann et al., 2001), autorgressive bias model (Del Giudice et al., 2013), approximate Bayesian computation (Sadegh and Vrugt, 2013), ~~and~~ data assimilation (Spaaks and

Bouten, 2013). Comparing different methods for accounting the residual assumptions are beyond the scope of this work. Fourth, this study focuses on formal Bayesian computation using formal likelihood functions, and comparison with other inference functions such as informal likelihood functions or approximate Bayesian computation is warranted in a future study.

—Based on the aforesaid conclusions and limitations, we recommend to start calibrating soil respiration models with simple SLS or SEP likelihood function. If the residuals characterization is adequate (e.g., Scharnagl et al., 2011), then the underlying assumptions are met. Otherwise, increase complexity of the data model until satisfactory results are obtained in terms of residuals characterization, posterior parameter estimation, and predictive performance. This is similar to the procedure given in Smith et al. (2015). Although the empirical findings of this study provide general guidelines for data model selection ~~of microbial~~for soil respiration modelings, more comparative studies are needed to validate and refute the findings of this study.

Acronyms

4C	Four carbon pool model
5C	Five carbon pool model
6C	Six carbon pool model
CUE	Microbial carbon use efficiency
DOC	Dissolved organic carbon
ENZ	Enzymes
MCMC	Markov chain Monte Carlo
MIC	Microbial biomass
NSME	Nash-Sutcliffe model efficiency
PDF	Probability density function
RMS	Relative model score
SEP	Skew exponential power distribution
SEP-AC	Skew exponential power distribution with autocorrelation
SLS	Standard least square
SLS-AC	Standard least square with autocorrelation
SOC	Soil organic carbon
WLS	Weighted least squared
WLS-AC	Weight least square with autocorrelation
WSEP	Weighted skew exponential power distribution
WSEP-AC	Weighted skew exponential power distribution with autocorrelation

Code and data availability

The data and codes and models used to produce this paper are available on contact of the corresponding author at mye@fsu.edu. We cannot publicly share the workflow because MT-DREAM_(ZS) code (Laloy and Vrugt, 2012) , which is a main component in the workflow, is in the process of becoming a commercial code.

Author contributions

ASE developed and implemented the code for the eight data models for soil respiration modeling, and prepared the manuscript with contribution of all co-authors. MY developed the research idea and outline, and supervised the research implementation when ASE was a post-doc at Florida State University. GN developed the soil respiration models. GAB collected and processed the eddy-covariance data used for model calibration.

Competing interests

The authors declare that they have no conflict of interest.

Acknowledgement

~~The first two authors~~^{is work-was-were} supported by the U.S. Department of Energy Early-Career grantAward, DE-SC0008272~~and~~. The first author was also partly supported by the -U.S. National Science Foundation Award# OIA-1557349. The second author was also partly supported by U.S. Department of Energy grant DE-SC0019438 and U.S. National Science Foundation grant EAR-1552329. We thank two anonymous reviewers for providing comments that helped to improve the manuscript.

References

Ahrens, B., Reichstein, M., Borken, W., Muhr, J., Trumbore, S. E. and Wutzler, T.: Bayesian calibration of a soil organic carbon model using ΔC measurements of soil organic carbon

914 and heterotrophic respiration as joint constraints, *Biogeosciences*, 11(8), 2147–2168,
 915 doi:10.5194/bg-11-2147-2014, 2014.

916 Allison, S. D., Wallenstein, M. D. and Bradford, M. A.: Soil-carbon response to warming
 917 dependent on microbial physiology, *Nat. Geosci.*, 3, 336 [online] Available from:
 918 <http://dx.doi.org/10.1038/ngeo846>, 2010.

919 Ammann, L., Reichert, P. and Fenicia, F.: A framework for likelihood functions of deterministic
 920 hydrological models, *Hydrol. Earth Syst. Sci.*, (August), 2018.

921 Bagnara, M., Sottocornola, M., Cescatti, A., Minerbi, S., Montagnani, L., Gianelle, D. and
 922 Magnani, F.: Bayesian optimization of a light use efficiency model for the estimation of
 923 daily gross primary productivity in a range of Italian forest ecosystems, *Ecol. Modell.*, 306,
 924 57–66, doi:10.1016/j.ecolmodel.2014.09.021, 2015.

925 Bagnara, M., Oijen, M. Van, Cameron, D., Gianelle, D., Magnani, F. and Sottocornola, M.:
 926 Bayesian calibration of simple forest models with multiplicative mathematical structure :
 927 A case study with two Light Use Efficiency models in an alpine forest, *Ecol. Modell.*,
 928 371(January), 90–100, doi:10.1016/j.ecolmodel.2018.01.014, 2018.

929 Barr, J. G., Engel, V., Fuentes, J. D., Fuller, D. O. and Kwon, H.: Modeling light use efficiency in
 930 a subtropical mangrove forest equipped with CO₂ eddy covariance, *Biogeosciences*, 10(3),
 931 2145–2158, doi:10.5194/bg-10-2145-2013, 2013.

932 Barron-gafford, G. A., Cable, J. M., Bentley, L. P., Scott, R. L., Huxman, T. E., Jenerette, G. D.
 933 and Ogle, K.: Quantifying the timescales over which exogenous and endogenous
 934 conditions affect soil respiration, *New Phytol.*, 2014.

935 Barron-Gafford, G. A., Scott, R. L., Jenerette, G. D. and Huxman, T. E.: The relative controls of
 936 temperature, soil moisture, and plant functional group on soil CO₂ efflux at

937 diel, seasonal, and annual scales, *J. Geophys. Res. Biogeosciences*, 116(1), 1–16,
 938 doi:10.1029/2010JG001442, 2011.

939 Berryman, E. M., Frank, J. M., Massman, W. J. and Ryan, M. G.: Agricultural and Forest
 940 Meteorology Using a Bayesian framework to account for advection in seven years of
 941 snowpack CO₂ fluxes in a mortality-impacted subalpine forest, *Agric. For. Meteorol.*,
 942 249(April 2017), 420–433, doi:10.1016/j.agrformet.2017.11.004, 2018.

943 Box, G. E. P. and Tiao, G. C.: *Bayesian inference in statistical analysis*, Wiley., 1992.

944 Braakhekke, M. C., Beer, C., Schrumpf, M., Ekici, A., Ahrens, B., Hoosbeek, M. R., Kruijt, B.,
 945 Kabat, P. and Reichstein, M.: The use of radiocarbon to constrain current and future soil
 946 organic matter turnover and transport in a temperate forest, *J. Geophys. Res.*
 947 *Biogeosciences*, 372–391, doi:10.1002/2013JG002420.Received, 2014.

948 Bradford, M. A., Davies, C. A., Frey, S. D., Maddox, T. R., Melillo, J. M., Mohan, J. E., Reynolds,
 949 J. F., Treseder, K. K. and Wallenstein, M. D.: Thermal adaptation of soil microbial
 950 respiration to elevated temperature, *Ecol. Lett.*, 11(12), 1316–1327, doi:10.1111/j.1461-
 951 0248.2008.01251.x, 2008.

952 Braswell, B. H., Sacks, W. J., Linder, E. and Schimel, D. S.: Estimating diurnal to annual
 953 ecosystem parameters by synthesis of a carbon flux model with eddy covariance net
 954 ecosystem exchange observations, *Glob. Chang. Biol.*, 335–355, doi:10.1111/j.1365-
 955 2486.2005.00897.x, 2015.

956 Cable, J. M., Ogle, K., Williams, D. G., Weltzin, J. F. and Huxman, T. E.: Soil Texture Drives
 957 Responses of Soil Respiration to Precipitation Pulses in the Sonoran Desert : Implications
 958 for Climate Change, *Ecosystems*, 961–979, doi:10.1007/s10021-008-9172-x, 2008.

959 Cable, J. M., Ogle, K., Lucas, R. W., Huxman, T. E., Loik, M. E., Smith, S. D., Tissue, D. T.,

- Ewers, B. E., Pendall, E., Welker, J. M., Charlet, T. N., Cleary, M., Griffith, A., Nowak, R. S., Rogers, M., Steltzer, H., Sullivan, P. F. and Gestel, N. C. Van: The temperature responses of soil respiration in deserts : a seven desert synthesis, *Biogeochemistry*, 71–90, doi:10.1007/s10533-010-9448-z, 2011.
- Chatfield, C.: The analysis of time series : an introduction, Chapman & Hall/CRC. [online] Available from: <https://www.crcpress.com/The-Analysis-of-Time-Series-An-Introduction-Sixth-Edition/Chatfield/p/book/9781584883173> (Accessed 9 April 2019), 2004.
- Chevallier, F. and O'Dell, C. W.: Error statistics of Bayesian CO₂ flux inversion schemes as seen from GOSAT, *Geophys. Res. Lett.*, 40(6), 1252–1256, doi:10.1002/grl.50228, 2013.
- Correia, A. C., Minunno, F., Caldeira, M. C., Banza, J., Mateus, J., Carneiro, M., Wingate, L., Shvaleva, A., Ramos, A., Jongen, M., Bugalho, M. N., Nogueira, C., Lecomte, X. and Pereira, J. S.: Agriculture , Ecosystems and Environment Soil water availability strongly modulates soil CO₂ efflux in different Mediterranean ecosystems : Model calibration using the Bayesian approach, *Agric. Ecosyst. Environ.*, 161, 88–100, doi:10.1016/j.agee.2012.07.025, 2012.
- Davidson, E. A. and Janssens, I. A.: Temperature sensitivity of soil carbon decomposition and feedbacks to climate change, *Nature*, 440, 165 [online] Available from: <http://dx.doi.org/10.1038/nature04514>, 2006.
- Davidson, E. A., Samanta, S., Caramori, S. S. and Savage, K.: The Dual Arrhenius and Michaelis–Menten kinetics model for decomposition of soil organic matter at hourly to seasonal time scales, *Glob. Chang. Biol.*, 18(1), 371–384, doi:10.1111/j.1365-2486.2011.02546.x, 2011.
- Du, Z., Nie, Y., He, Y., Yu, G. and Wang, H.: *Tellus B : Chemical and Physical Meteorology*

983 Complementarity of flux- and biometric-based data to constrain parameters in a terrestrial
984 carbon model Complementarity of flux- and biometric-based data to constrain parameters
985 in a terrestrial carbon model, *Tellus B Chem. Phys. Meteorol.*, 0889,
986 doi:10.3402/tellusb.v67.24102, 2015.

987 Du, Z., Zhou, X., Shao, J., Yu, G., Wang, H., Zhai, D., Xai, J. and Luo, Y.: *Journal of Advances*
988 in Modeling Earth Systems, *J. Adv. Model. Earth Syst.*, 548–565,
989 doi:10.1002/2016MS000687.Received, 2017.

990 Elshall, A. S. and Tsai, F. T.-C.: Constructive epistemic modeling of groundwater flow with
991 geological structure and boundary condition uncertainty under the Bayesian paradigm, *J.*
992 *Hydrol.*, 517, doi:10.1016/j.jhydrol.2014.05.027, 2014.

993 Elshall, A. S., Ye, M., Pei, Y., Zhang, F., Niu, G.-Y. and Barron-Gafford, G. A.: Relative model
994 score: a scoring rule for evaluating ensemble simulations with application to microbial soil
995 respiration modeling, *Stoch. Environ. Res. Risk Assess.*, doi:10.1007/s00477-018-1592-3,
996 2018a.

997 Elshall, A. S., Ye, M., Pei, Y., Zhang, F., Niu, G.-Y. and Barron-Gafford, G. A.: Relative model
998 score: a scoring rule for evaluating ensemble simulations with application to microbial soil
999 respiration modeling, *Stoch. Environ. Res. Risk Assess.*, 32(10), 2809–2819,
1000 doi:10.1007/s00477-018-1592-3, 2018b.

1001 Evin, G., Kavetski, D., Thyer, M. and Kuczera, G.: Pitfalls and improvements in the joint inference
1002 of heteroscedasticity and autocorrelation in hydrological model calibration, *Water Resour.*
1003 *Res.*, 49(7), 4518–4524, doi:10.1002/wrcr.20284, 2013.

1004 Evin, G., Thyer, M., Kavetski, D., McInerney, D. and Kuczera, G.: Comparison of joint versus
1005 postprocessor approaches for hydrological uncertainty estimation accounting for error

1006 autocorrelation and heteroscedasticity, *Water Resour. Res.*, 50(3), 2350–2375,
 1007 doi:10.1002/2013WR014185, 2014.

1008 Fernández-Martínez, M., Vicca, S., Janssens, I. A., Sardans, J., Luyssaert, S., Campioli, M.,
 1009 Chapin III, F. S., Ciais, P., Malhi, Y., Obersteiner, M., Papale, D., Piao, S. L., Reichstein,
 1010 M., Rodà, F. and Peñuelas, J.: Nutrient availability as the key regulator of global forest
 1011 carbon balance, *Nat. Clim. Chang.*, 4, 471 [online] Available from:
 1012 <http://dx.doi.org/10.1038/nclimate2177>, 2014.

1013 Gelman, A. and Rubin, D. B.: Inference from Iterative Simulation Using Multiple Sequences, *Stat.*
 1014 *Sci.*, 7(4), 457–472, doi:10.1214/ss/1177011136, 1992.

1015 German, D. P., Marcelo, K. R. B., Stone, M. M. and Allison, S. D.: The Michaelis–Menten kinetics
 1016 of soil extracellular enzymes in response to temperature: a cross-latitudinal study, *Glob.*
 1017 *Chang. Biol.*, 18(4), 1468–1479, doi:10.1111/j.1365-2486.2011.02615.x, 2011.

1018 Del Giudice, D., Honti, M., Scheidegger, A., Albert, C., Reichert, P. and Rieckermann, J.:
 1019 Improving uncertainty estimation in urban hydrological modeling by statistically
 1020 describing bias, *Hydrol. Earth Syst. Sci.*, 17(10), 4209–4225, doi:10.5194/hess-17-4209-
 1021 2013, 2013.

1022 Gragne, A. S., Sharma, A., Mehrotra, R. and Alfredsen, K.: Improving real-time inflow forecasting
 1023 into hydropower reservoirs through a complementary modelling framework, *Hydrol. Earth*
 1024 *Syst. Sci.*, 19(8), 3695–3714, doi:10.5194/hess-19-3695-2015, 2015.

1025 Hararuk, O., Xia, J. and Luo, Y.: Evaluation and improvement of a global land model against soil
 1026 carbon data using a Bayesian Markov chain Monte Carlo method, *J. Geophys. Res.*
 1027 *Biogeosciences*, 119(3), 403–417, doi:10.1002/2013JG002535, 2014.

1028 Hashimoto, S., Morishita, T., Sakata, T., Ishizuka, S., Kaneko, S. and Takahashi, M.: Simple

1029 models for soil CO₂, CH₄, and N₂O fluxes calibrated using a Bayesian approach and
 1030 multi-site data, *Ecol. Modell.*, 222(7), 1283–1292, doi:10.1016/j.ecolmodel.2011.01.013,
 1031 2011.

1032 He, H., Meyer, A., Jansson, P., Svensson, M., Rütting, T. and Klemmedtsson, L.: Simulating
 1033 ectomycorrhiza in boreal forests : implementing ectomycorrhizal fungi model MYCOFON
 1034 in CoupModel (v5), *Geosci. Model Dev.*, 725–751, 2018.

1035 Hilton, T. W., Davis, K. J. and Keller, K.: Evaluating terrestrial CO₂flux diagnoses and
 1036 uncertainties from a simple land surface model and its residuals, *Biogeosciences*, 11(2),
 1037 217–235, doi:10.5194/bg-11-217-2014, 2014.

1038 Hoeting, J. A., Madigan, D., Raftery, A. E. and Volinsky, C. T.: Bayesian model averaging: a
 1039 tutorial (with comments by M. Clyde, David Draper and E. I. George, and a rejoinder by
 1040 the authors, *Stat. Sci.*, 14(4), 382–417, doi:10.1214/ss/1009212519, 1999.

1041 Höglberg, P. and Read, D. J.: Towards a more plant physiological perspective on soil ecology,
 1042 *Trends Ecol. Evol.*, 21(10), 548–554, doi:10.1016/j.tree.2006.06.004, 2006.

1043 Hublart, P., Ruelland, D., De Cortázar-Atauri, I. G., Gascoin, S., Lhermitte, S. and Ibacache, A.:
 1044 Reliability of lumped hydrological modeling in a semi-arid mountainous catchment facing
 1045 water-use changes, *Hydrol. Earth Syst. Sci.*, 20(9), 3691–3717, doi:10.5194/hess-20-3691-
 1046 2016, 2016.

1047 Ishikura, K., Yamada, H., Toma, Y., Takakai, F., Darung, U., Limin, A. and Limin, S. H.: Soil
 1048 Science and Plant Nutrition Effect of groundwater level fluctuation on soil respiration rate
 1049 of tropical peatland in Central Kalimantan , Indonesia, *Soil Sci. Plant Nutr.*, 63(1), 1–13,
 1050 doi:10.1080/00380768.2016.1244652, 2017.

1051 Janssens, I. A., Freibauer, A., Ciais, P., Smith, P., Nabuurs, G.-J., Folberth, G., Schlamadinger, B.,

1052 Hutjes, R. W. A., Ceulemans, R., Schulze, E.-D., Valentini, R. and Dolman, A. J.: Europe's
 1053 terrestrial biosphere absorbs 7 to 12% of European anthropogenic CO₂ emissions.,
 1054 Science, 300(5625), 1538–42, doi:10.1126/science.1083592, 2003.

1055 Katz, R. W., Craigmile, P. F., Guttorp, P., Haran, M., Sansó, B. and Stein, M. L.: Uncertainty
 1056 analysis in climate change assessments, Nat. Clim. Chang., 3, 769 [online] Available from:
 1057 <http://dx.doi.org/10.1038/nclimate1980>, 2013.

1058 Kavetski, D., Franks, S. W. and Kuczera, G.: Confronting Input Uncertainty in Environmental
 1059 Modelling, Calibration Watershed Model., doi:doi:10.1029/WS006p0049, 2013.

1060 Keenan, T. F., Davidson, E., Moffat, A. M., Munger, W. and Richardson, A. D.: Using model-data
 1061 fusion to interpret past trends, and quantify uncertainties in future projections, of terrestrial
 1062 ecosystem carbon cycling, Glob. Chang. Biol., 18(8), 2555–2569, doi:10.1111/j.1365-
 1063 2486.2012.02684.x, 2012.

1064 Kim, Y., Nishina, K., Chae, N., Park, S. J., Yoon, Y. J. and Lee, B. Y.: Constraint of soil moisture
 1065 on CO₂ efflux from tundra lichen, moss, and tussock in Council, Alaska, using a
 1066 hierarchical Bayesian model, Biogeosciences, 5567–5579, doi:10.5194/bg-11-5567-2014,
 1067 2014.

1068 Klemetsson, L., Jansson, P. E., Gustafsson, D., Karlberg, L., Weslien, P., Von Arnold, K.,
 1069 Ernfors, M., Langvall, O. and Lindroth, A.: Bayesian calibration method used to elucidate
 1070 carbon turnover in forest on drained organic soil, Biogeochemistry, 89(1), 61–79,
 1071 doi:10.1007/s10533-007-9169-0, 2008.

1072 Laloy, E. and Vrugt, J. A.: High-dimensional posterior exploration of hydrologic models using
 1073 multiple-try DREAM(ZS) and high-performance computing, Water Resour. Res., 48(1),
 1074 doi:10.1029/2011WR010608, 2012.

1075 Li, J., Wang, G., Allison, S. D., Mayes, M. A. and Luo, Y.: Soil carbon sensitivity to temperature
 1076 and carbon use efficiency compared across microbial-ecosystem models of varying
 1077 complexity, *Biogeochemistry*, 119, 67–84 [online] Available from:
 1078 <http://www.jstor.org/stable/24716883>, 2014.

1079 Li, M., Wang, Q. J., Bennett, J. C. and Robertson, D. E.: A strategy to overcome adverse effects
 1080 of autoregressive updating of streamflow forecasts, *Hydrol. Earth Syst. Sci.*, 19(1), 1–15,
 1081 doi:10.5194/hess-19-1-2015, 2015.

1082 Li, M., Wang, Q. J., Bennett, J. C. and Robertson, D. E.: Error reduction and representation in
 1083 stages (ERRIS) in hydrological modelling for ensemble streamflow forecasting, *Hydrol.*
 1084 *Earth Syst. Sci.*, 20(9), 3561–3579, doi:10.5194/hess-20-3561-2016, 2016a.

1085 Li, Q., Xia, J., Shi, Z., Huang, K., Du, Z. and Lin, G.: Variation of parameters in a Flux-Based
 1086 Ecosystem Model across 12 sites of terrestrial ecosystems in the conterminous USA, *Ecol.*
 1087 *Modell.*, 336, 57–69, doi:10.1016/j.ecolmodel.2016.05.016, 2016b.

1088 Lu, D., Ye, M., Meyer, P. D., Curtis, G. P., Shi, X., Niu, X.-F. and Yabusaki, S. B.: Effects of error
 1089 covariance structure on estimation of model averaging weights and predictive performance,
 1090 *Water Resour. Res.*, 49(9), 6029–6047, doi:10.1002/wrcr.20441, 2013.

1091 Luo, Y., Ogle, K., Tucker, C., Fei, S., Gao, C., LaDeau, S., Clark, J. S. and Schimel, D. S.:
 1092 Ecological forecasting and data assimilation in a data-rich era, *Ecol. Appl.*, 21(5), 1429–
 1093 1442, doi:10.1890/09-1275.1, 2011.

1094 Luo, Y., Keenan, T. F. and Smith, M.: Predictability of the terrestrial carbon cycle, *Glob. Chang.*
 1095 *Biol.*, 21(5), 1737–1751, doi:10.1111/gcb.12766, 2014.

1096 Manzoni, S., Taylor, P., Richter, A., Porporato, A. and Ågren, G. I.: Environmental and
 1097 stoichiometric controls on microbial carbon-use efficiency in soils, *New Phytol.*, 196(1),

1098 79–91, doi:10.1111/j.1469-8137.2012.04225.x, 2012.

1099 McInerney, D., Thyer, M., Kavetski, D., Lerat, J. and Kuczera, G.: Improving probabilistic
 1100 prediction of daily streamflow by identifying Pareto optimal approaches for modeling
 1101 heteroscedastic residual errors, *Water Resour. Res.*, 53, 2199–2239,
 1102 doi:10.1002/2016WR019168.Received, 2017.

1103 Menichetti, L., Kätterer, T. and Leifeld, J.: Parametrization consequences of constraining soil
 1104 organic matter models by total carbon and radiocarbon using long-term field data,
 1105 *Biogeosciences*, 3003–3019, doi:10.5194/bg-13-3003-2016, 2016.

1106 Nash, J. E. and Sutcliffe, J. V: River flow forecasting through conceptual models part I — A
 1107 discussion of principles, *J. Hydrol.*, 10(3), 282–290, doi:https://doi.org/10.1016/0022-
 1108 1694(70)90255-6, 1970.

1109 Ogle, K., Ryan, E., Dijkstra, F. A. and Pendall, E.: *Journal of Geophysical Research:*
 1110 *Biogeosciences*, *J. Geophys. Res. Biogeoscoences*, 1–14, doi:10.1002/2016JG003385,
 1111 2016.

1112 Pappenberger, F. and Beven, K. J.: Ignorance is bliss: Or seven reasons not to use uncertainty
 1113 analysis, *Water Resour. Res.*, 42(5), doi:10.1029/2005WR004820, 2006.

1114 Peters, W., Jacobson, A. R., Sweeney, C., Andrews, A. E., Conway, T. J., Masarie, K., Miller, J.
 1115 B., Bruhwiler, L. M. P., Pétron, G., Hirsch, A. I., Worthy, D. E. J., van der Werf, G. R.,
 1116 Randerson, J. T., Wennberg, P. O., Krol, M. C. and Tans, P. P.: An atmospheric perspective
 1117 on North American carbon dioxide exchange: CarbonTracker., *Proc. Natl. Acad. Sci. U. S.*
 1118 *A.*, 104(48), 18925–30, doi:10.1073/pnas.0708986104, 2007.

1119 Le Quéré, C., Peters, G. P., Andres, R. J., Andrew, R. M., Boden, T. A., Ciais, P., Friedlingstein,
 1120 P., Houghton, R. A., Marland, G., Moriarty, R., Sitch, S., Tans, P., Arneth, A., Arvanitis,

1121 A., Bakker, D. C. E., Bopp, L., Canadell, J. G., Chini, L. P., Doney, S. C., Harper, A.,
 1122 Harris, I., House, J. I., Jain, A. K., Jones, S. D., Kato, E., Keeling, R. F., Klein Goldewijk,
 1123 K., Körtzinger, A., Koven, C., Lefèvre, N., Maignan, F., Omar, A., Ono, T., Park, G.-H.,
 1124 Pfeil, B., Poulter, B., Raupach, M. R., Regnier, P., Rödenbeck, C., Saito, S., Schwinger, J.,
 1125 Segschneider, J., Stocker, B. D., Takahashi, T., Tilbrook, B., van Heuven, S., Viovy, N.,
 1126 Wanninkhof, R., Wiltshire, A. and Zaehle, S.: Global carbon budget 2013, *Earth Syst. Sci.*
 1127 *Data*, 6(1), 235–263, doi:10.5194/essd-6-235-2014, 2014.

1128 Raich, J. W. J. J. W., Potter, C. S. C. and Bhagawati, D.: Interannual variability in global soil
 1129 respiration, 1980-94, *Glob. Chang. Biol.*, 8, 800–812, doi:10.1046/j.1365-
 1130 2486.2002.00511.x, 2002.

1131 Ren, X., He, H., Moore, D. J. P., Zhang, L., Liu, M., Li, F., Yu, G. and Wang, H.: Uncertainty
 1132 analysis of modeled carbon and water fluxes in a subtropical coniferous plantation, *J.*
 1133 *Geophys. Res. Biogeosciences*, 118(4), 1674–1688, doi:10.1002/2013JG002402, 2013.

1134 Ricciuto, D. M., King, A. W., Dragoni, D. and Post, W. M.: Parameter and prediction uncertainty
 1135 in an optimized terrestrial carbon cycle model: Effects of constraining variables and data
 1136 record length, *J. Geophys. Res. Biogeosciences*, 116(1), 1–17,
 1137 doi:10.1029/2010JG001400, 2011.

1138 Richardson, A. D. and Hollinger, D. Y.: Statistical modeling of ecosystem respiration using eddy
 1139 covariance data: Maximum likelihood parameter estimation, and Monte Carlo simulation
 1140 of model and parameter uncertainty, applied to three simple models, *Agric. For. Meteorol.*,
 1141 131(3–4), 191–208, doi:10.1016/j.agrformet.2005.05.008, 2005.

1142 Sadegh, M. and Vrugt, J. A.: Bridging the gap between GLUE and formal statistical approaches:
 1143 Approximate Bayesian computation, *Hydrol. Earth Syst. Sci.*, 17(12), 4831–4850,

doi:10.5194/hess-17-4831-2013, 2013.

Scharnagl, B., Vrugt, J. A., Vereecken, H. and Herbst, M.: Inverse modelling of in situ soil water dynamics: Investigating the effect of different prior distributions of the soil hydraulic parameters, *Hydrol. Earth Syst. Sci.*, 15(10), 3043–3059, doi:10.5194/hess-15-3043-2011, 2011.

Schimel, J. P. and Weintraub, M. N.: The implications of exoenzyme activity on microbial carbon and nitrogen limitation in soil: a theoretical model, *Soil Biol. Biochem.*, 35(4), 549–563, doi:10.1016/S0038-0717(03)00015-4, 2003.

Schmidt, M. W. I., Torn, M. S., Abiven, S., Dittmar, T., Guggenberger, G., Janssens, I. A., Kleber, M., Kögel-Knabner, I., Lehmann, J., Manning, D. A. C., Nannipieri, P., Rasse, D. P., Weiner, S. and Trumbore, S. E.: Persistence of soil organic matter as an ecosystem property, *Nature*, 478(7367), 49–56, doi:10.1038/nature10386, 2011.

Scholz, K., Hammerle, A., Hiltbrunner, E. and Wohlfahrt, G.: Analyzing the Effects of Growing Season Length on the Net Ecosystem Production of an Alpine Grassland Using Model – Data Fusion, *Ecosystems*, 21(5), 982–999, doi:10.1007/s10021-017-0201-5, 2018.

Schoups, G. and Vrugt, J. A.: A formal likelihood function for parameter and predictive inference of hydrologic models with correlated, heteroscedastic, and non-Gaussian errors, *Water Resour. Res.*, 46(10), 1–17, doi:10.1029/2009WR008933, 2010.

Scott, R. L., Jenerette, G. D., Potts, D. L. and Huxman, T. E.: Effects of seasonal drought on net carbon dioxide exchange from a woody-plant-encroached semiarid grassland, *J. Geophys. Res. Biogeosciences*, 114(4), doi:10.1029/2008JG000900, 2009.

Shi, X., Ye, M., Curtis, G. P., Miller, G. L., Meyer, P. D., Kohler, M., Yabusaki, S. and Wu, J.: Assessment of parametric uncertainty for groundwater reactive transport modeling, *Water*

1167 Resour. Res., 50(5), 4416–4439, doi:10.1002/2013WR013755, 2014.

1168 Sinsabaugh, R. L., Manzoni, S., Moorhead, D. L. and Richter, A.: Carbon use efficiency of
1169 microbial communities: stoichiometry, methodology and modelling, *Ecol. Lett.*, 16(7),
1170 930–939, doi:10.1111/ele.12113, 2013.

1171 Smith, M. W., Bracken, L. J. and Cox, N. J.: Toward a dynamic representation of hydrological
1172 connectivity at the hillslope scale in semiarid areas, *Water Resour. Res.*, 46(12),
1173 doi:10.1029/2009WR008496, 2010a.

1174 Smith, T., Sharma, A., Marshall, L., Mehrotra, R. and Sisson, S.: Development of a formal
1175 likelihood function for improved Bayesian inference of ephemeral catchments, *Water*
1176 *Resour. Res.*, 46(12), 1–11, doi:10.1029/2010WR009514, 2010b.

1177 Smith, T., Marshall, L. and Sharma, A.: Modeling residual hydrologic errors with Bayesian
1178 inference, *J. Hydrol.*, 528, 29–37, doi:10.1016/j.jhydrol.2015.05.051, 2015.

1179 Spaaks, J. H. and Bouten, W.: Resolving structural errors in a spatially distributed hydrologic
1180 model using ensemble Kalman filter state updates, *Hydrol. Earth Syst. Sci.*, 17(9), 3455–
1181 3472, doi:10.5194/hess-17-3455-2013, 2013.

1182 Steinacher, M. and Joos, F.: Transient Earth system responses to cumulative carbon dioxide
1183 emissions: Linearities, uncertainties, and probabilities in an observation-constrained model
1184 ensemble, *Biogeosciences*, 13(4), 1071–1103, doi:10.5194/bg-13-1071-2016, 2016.

1185 Tang, J. and Riley, W. J.: Weaker soil carbon–climate feedbacks resulting from microbial and
1186 abiotic interactions, *Nat. Clim. Chang.*, 5, 56 [online] Available from:
1187 <http://dx.doi.org/10.1038/nclimate2438>, 2014.

1188 Tang, J. and Zhuang, Q.: A global sensitivity analysis and Bayesian inference framework for
1189 improving the parameter estimation and prediction of a process-based Terrestrial

1190 Ecosystem Model, *J. Geophys. Res. Atmos.*, 114(D15), doi:10.1029/2009JD011724, 2009.

1191 Thiemann, M., Trosset, M., Gupta, H. and Sorooshian, S.: Bayesian recursive parameter estimation
 1192 for hydrologic models, *Water Resour. Res.*, 37(10), 2521–2535,
 1193 doi:10.1029/2000WR900405, 2001.

1194 Thyer, M., Renard, B., Kavetski, D., Kuczera, G., Franks, S. W. and Srikanthan, S.: Critical
 1195 evaluation of parameter consistency and predictive uncertainty in hydrological modeling:
 1196 A case study using Bayesian total error analysis, *Water Resour. Res.*, 45(12), 1–22,
 1197 doi:10.1029/2008WR006825, 2009.

1198 Tiedeman, C. R. and Green, C. T.: Effect of correlated observation error on parameters,
 1199 predictions, and uncertainty, *Water Resour. Res.*, 49(10), 6339–6355,
 1200 doi:10.1002/wrcr.20499, 2013.

1201 Tsai, F. T.-C. and Elshall, A. S.: Hierarchical Bayesian model averaging for hydrostratigraphic
 1202 modeling: Uncertainty segregation and comparative evaluation, *Water Resour. Res.*, 49(9),
 1203 doi:10.1002/wrcr.20428, 2013.

1204 Tucker, C. L., Bell, J., Pendall, E. and Ogle, K.: Does declining carbon-use efficiency explain
 1205 thermal acclimation of soil respiration with warming?, *Glob. Chang. Biol.*, 252–263,
 1206 doi:10.1111/gcb.12036, 2013.

1207 Tucker, C. L., Young, J. M., Williams, D. G. and Ogle, K.: Process-based isotope partitioning of
 1208 winter soil respiration in a subalpine ecosystem reveals importance of rhizospheric
 1209 respiration, *Biogeochemistry*, 121, 389–408 [online] Available from:
 1210 <http://www.jstor.org/stable/24717586>, 2014.

1211 Tuomi, M., Vanhala, P., Karhu, K., Fritze, H. and Liski, J.: Heterotrophic soil respiration-
 1212 Comparison of different models describing its temperature dependence, *Ecol. Modell.*,

1213 211(1–2), 182–190, doi:10.1016/j.ecolmodel.2007.09.003, 2008.

1214 Vargas, R., Carbone, M. S., Reichstein, M. and Baldocchi, D. D.: Frontiers and challenges in soil
 1215 respiration research: from measurements to model-data integration, *Biogeochemistry*,
 1216 102(1), 1–13, doi:10.1007/s10533-010-9462-1, 2011.

1217 Vrugt, J. A. and Ter Braak, C. J. F.: DREAM(D): An adaptive Markov Chain Monte Carlo
 1218 simulation algorithm to solve discrete, noncontinuous, and combinatorial posterior
 1219 parameter estimation problems, *Hydrol. Earth Syst. Sci.*, 15(12), 3701–3713,
 1220 doi:10.5194/hess-15-3701-2011, 2011.

1221 Vrugt, J. A., ter Braak, C. J. F., Diks, C. G. H. and Schoups, G.: Hydrologic data assimilation using
 1222 particle Markov chain Monte Carlo simulation: Theory, concepts and applications, *Adv.*
 1223 *Water Resour.*, 51, 457–478, doi:10.1016/j.advwatres.2012.04.002, 2013.

1224 Wang, G., Post, W. M. and Mayes, M. A.: Development of microbial-enzyme-mediated
 1225 decomposition model parameters through steady-state and dynamic analyses, *Ecol. Appl.*,
 1226 23(1), 255–272, doi:10.1890/12-0681.1, 2013.

1227 Weijs, S. V., Schoups, G. and Van De Giesen, N.: Why hydrological predictions should be
 1228 evaluated using information theory, *Hydrol. Earth Syst. Sci.*, 14(12), 2545–2558,
 1229 doi:10.5194/hess-14-2545-2010, 2010.

1230 Westerberg, I. K., Guerrero, J. L., Younger, P. M., Beven, K. J., Seibert, J., Halldin, S., Freer, J.
 1231 E. and Xu, C. Y.: Calibration of hydrological models using flow-duration curves, *Hydrol.*
 1232 *Earth Syst. Sci.*, 15(7), 2205–2227, doi:10.5194/hess-15-2205-2011, 2011.

1233 Wieder, W. R., Bonan, G. B. and Allison, S. D.: Global soil carbon projections are improved by
 1234 modelling microbial processes, *Nat. Clim. Chang.*, 3, 909 [online] Available from:
 1235 <http://dx.doi.org/10.1038/nclimate1951>, 2013.

1236 Wieder, W. R., Allison, S. D., Davidson, E. A., Georgiou, K., Hararuk, O., He, Y., Hopkins, F.,
 1237 Luo, Y., Smith, M. J., Sulman, B., Todd-Brown, K., Wang, Y.-P., Xia, J. and Xu, X.:
 1238 Explicitly representing soil microbial processes in Earth system models, *Global*
 1239 *Biogeochem. Cycles*, 29(10), 1782–1800, doi:10.1002/2015GB005188, 2015.
 1240 Van Wijk, M. T., Van Putten, B., Hollinger, D. Y. and Richardson, A. D.: Comparison of different
 1241 objective functions for parameterization of simple respiration models, *J. Geophys. Res.*
 1242 *Biogeosciences*, 113(3), 1–11, doi:10.1029/2007JG000643, 2008.
 1243 Xu, T., White, L., Hui, D. and Luo, Y.: Probabilistic inversion of a terrestrial ecosystem model:
 1244 Analysis of uncertainty in parameter estimation and model prediction, *Global Biogeochem.*
 1245 *Cycles*, 20(2), 1–15, doi:10.1029/2005GB002468, 2006.
 1246 Xu, X., Schimel, J. P., Thornton, P. E., Song, X., Yuan, F. and Goswami, S.: Substrate and
 1247 environmental controls on microbial assimilation of soil organic carbon: a framework for
 1248 Earth system models, *Ecol. Lett.*, 17(5), 547–555, doi:10.1111/ele.12254, 2014.
 1249 Yeluripati, J. B., van Oijen, M., Wattenbach, M., Neftel, A., Ammann, A., Parton, W. J. and Smith,
 1250 P.: Bayesian calibration as a tool for initialising the carbon pools of dynamic soil models,
 1251 *Soil Biol. Biochem.*, 41(12), 2579–2583, doi:10.1016/j.soilbio.2009.08.021, 2009.
 1252 Yuan, W., Liang, S., Liu, S., Weng, E., Luo, Y. and Hollinger, D.: Improving model parameter
 1253 estimation using coupling relationships between vegetation production and ecosystem
 1254 respiration, *Ecol. Modell.*, 240, 29–40, doi:10.1016/j.ecolmodel.2012.04.027, 2012.
 1255 Yuan, W., Xu, W., Ma, M., Chen, S. and Liu, W.: Agricultural and Forest Meteorology Improved
 1256 snow cover model in terrestrial ecosystem models over the Qinghai – Tibetan Plateau,
 1257 *Agric. For. Meteorol.*, 218–219, 161–170, doi:10.1016/j.agrformet.2015.12.004, 2016.
 1258 Zhang, X., Niu, G.-Y., Elshall, A. S., Ye, M., Barron-Gafford, G. A. and Pavao-Zuckerman, M.:

1259 Assessing five evolving microbial enzyme models against field measurements from a
1260 semiarid savannah - What are the mechanisms of soil respiration pulses?, *Geophys. Res.*
1261 *Lett.*, 41(18), doi:10.1002/2014GL061399, 2014.

1262 Zhou, X., Luo, Y., Gao, C., Verburg, P. S. J., Arnone, J. A., Darrouzet-Nardi, A. and Schimel, D.
1263 S.: Concurrent and lagged impacts of an anomalously warm year on autotrophic and
1264 heterotrophic components of soil respiration: A deconvolution analysis, *New Phytol.*,
1265 187(1), 184–198, doi:10.1111/j.1469-8137.2010.03256.x, 2010.

1266

Figure 1. Diagram of model 6C representing the processes of (1) degradation of soil organic carbon (SOC) to dissolved organic carbon (DOC) through catalysis of enzymes (ENZ) produced by microbes (MIC), (2) MIC uptake of DOC, and (3) microbial (MIC) respiration to produce CO₂ (CUE is the carbon use efficiency). SOC degradation and microbial uptake rates are controlled by water saturation (θ / θ_s). The DOC and ENZ pools are split into two subpools, one for the wet zone and the other for the dry zone of the soil pore space. Microbial uptake of DOC occurs only in the wet zone, and the uptake rate is linearly related to θ / θ_s . Catalysis through ENZ in the wet zone is proportional to θ / θ_s , while that in the dry zone is proportional to $1 - \theta / \theta_s$. V_{max} (s⁻¹) is the maximum rate, and K_m is the half-saturation concentration.

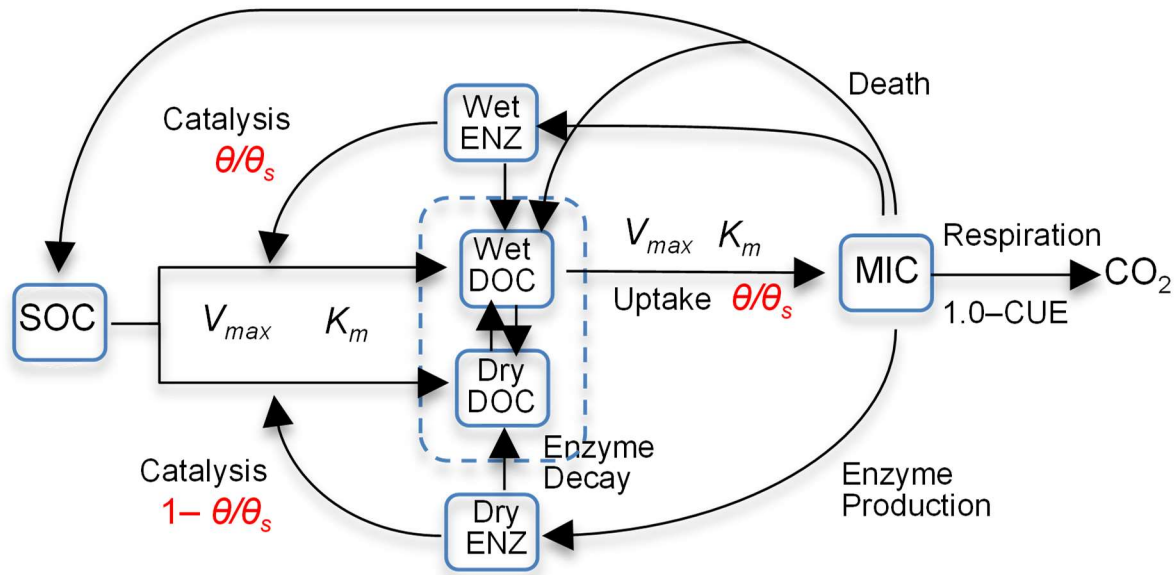


Figure 2. Time series of soil moisture and efflux observations. The dashed line marks the divide of the dataset into calibration and validation periods.

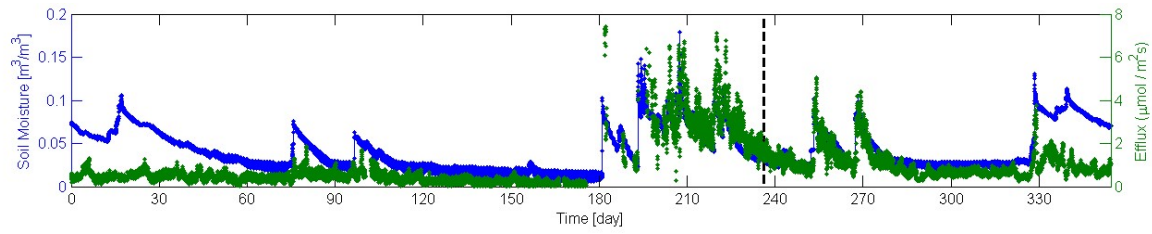


Figure 3. Residual analysis of the best realization (among multiple MCMC realizations) for model 6C using data models (a-c) SLS and (d-f) WSEP-AC.

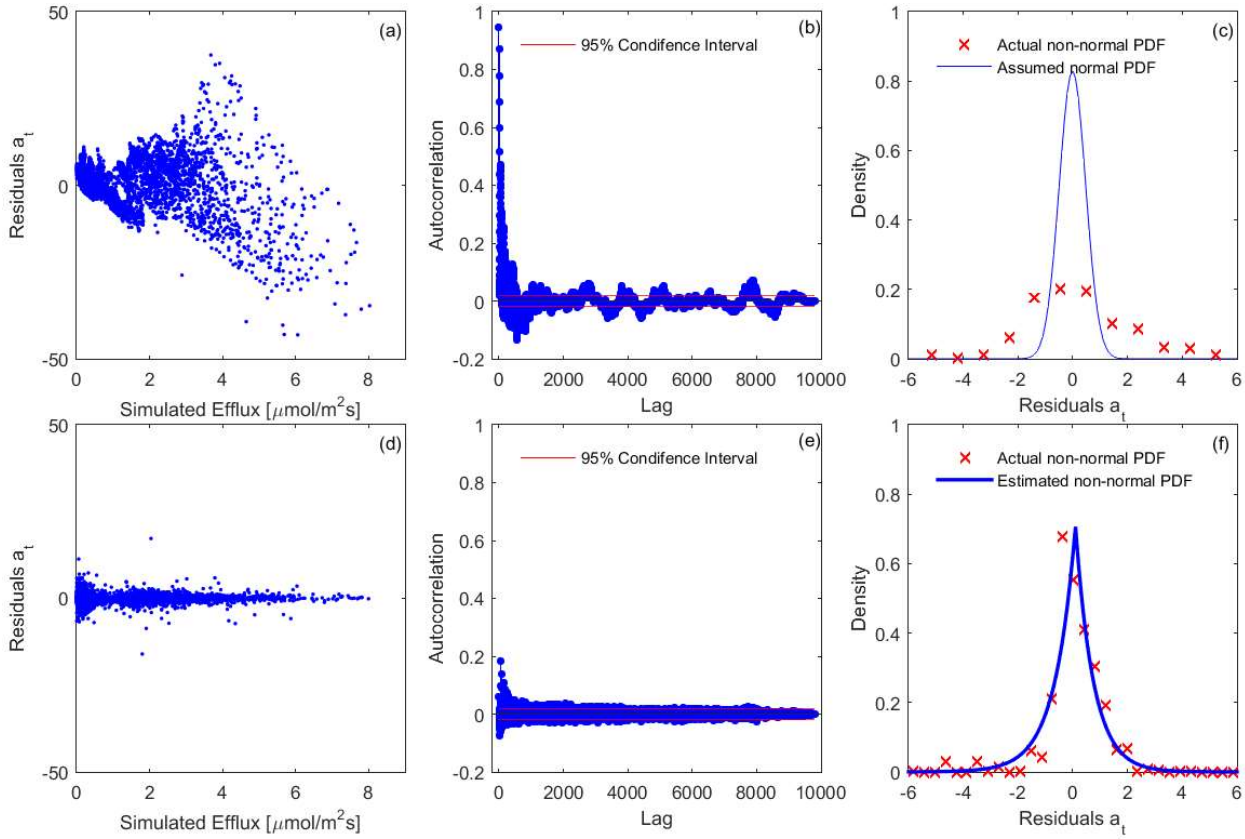


Figure 4. Residual quantile-quantile (Q-Q) plots of the best realization (among multiple MCMC realizations) for the three soil respiration models and eight data models.

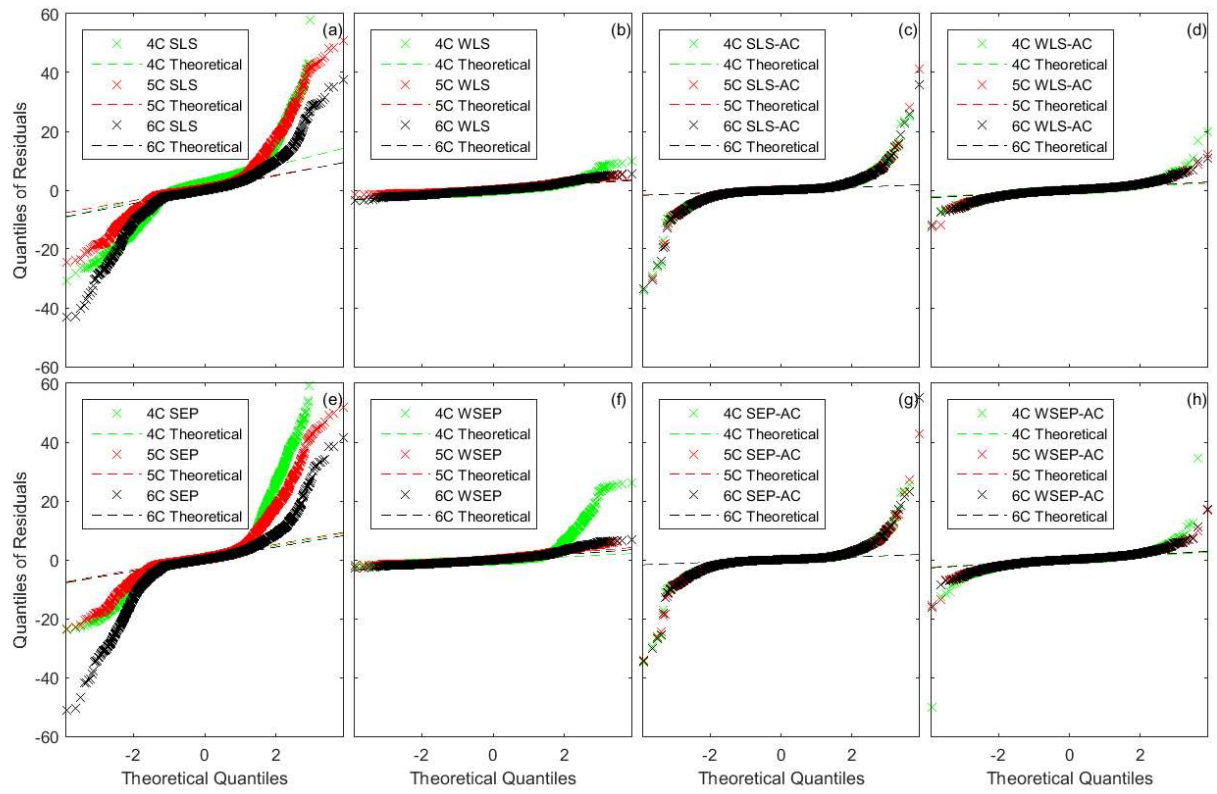


Figure 5. Marginal posterior parameter density of carbon use efficiency (CUE) for the three soil respiration models and eight data models.

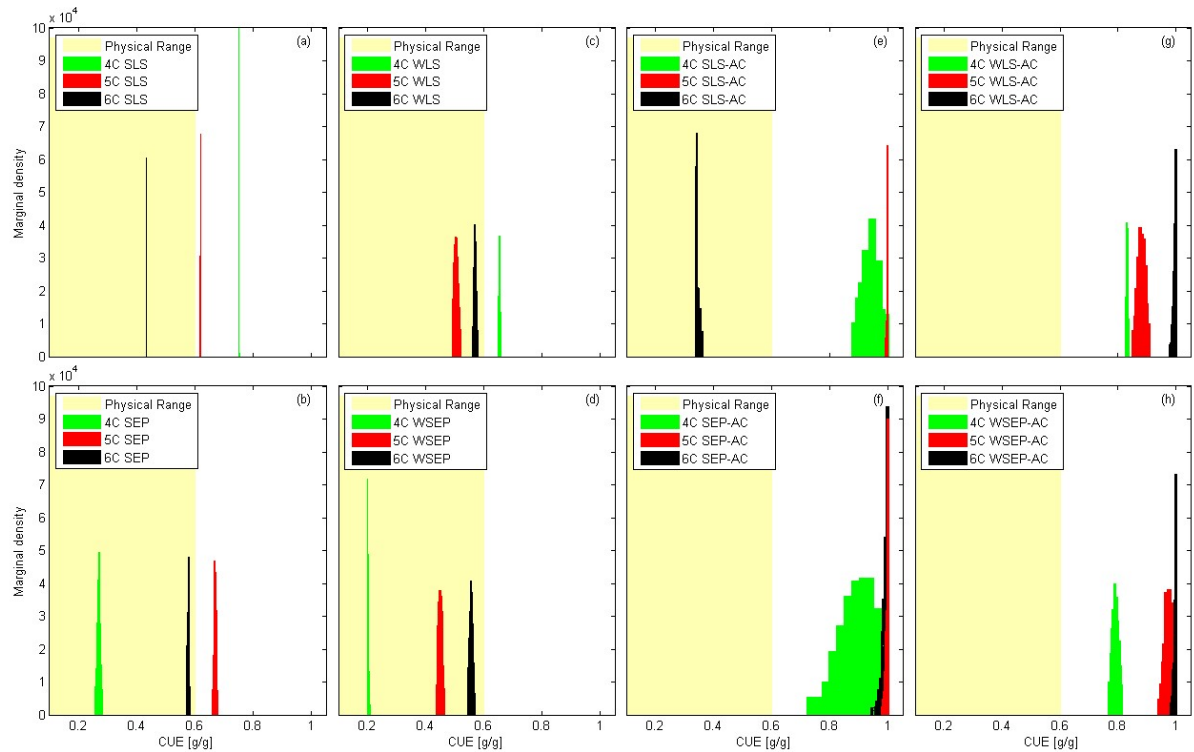


Figure 6. Observation data (blue dots) and mean prediction (green line) and 95% credible intervals (red line) of prediction ensembles for (a)-(f) the calibration period and (g)-(l) the validation period. The plots are for the three soil respiration models using data models SLS and WSEP-AC. *The prediction ensembles are generated to consider parametric uncertainty of the soil respiration models only.*

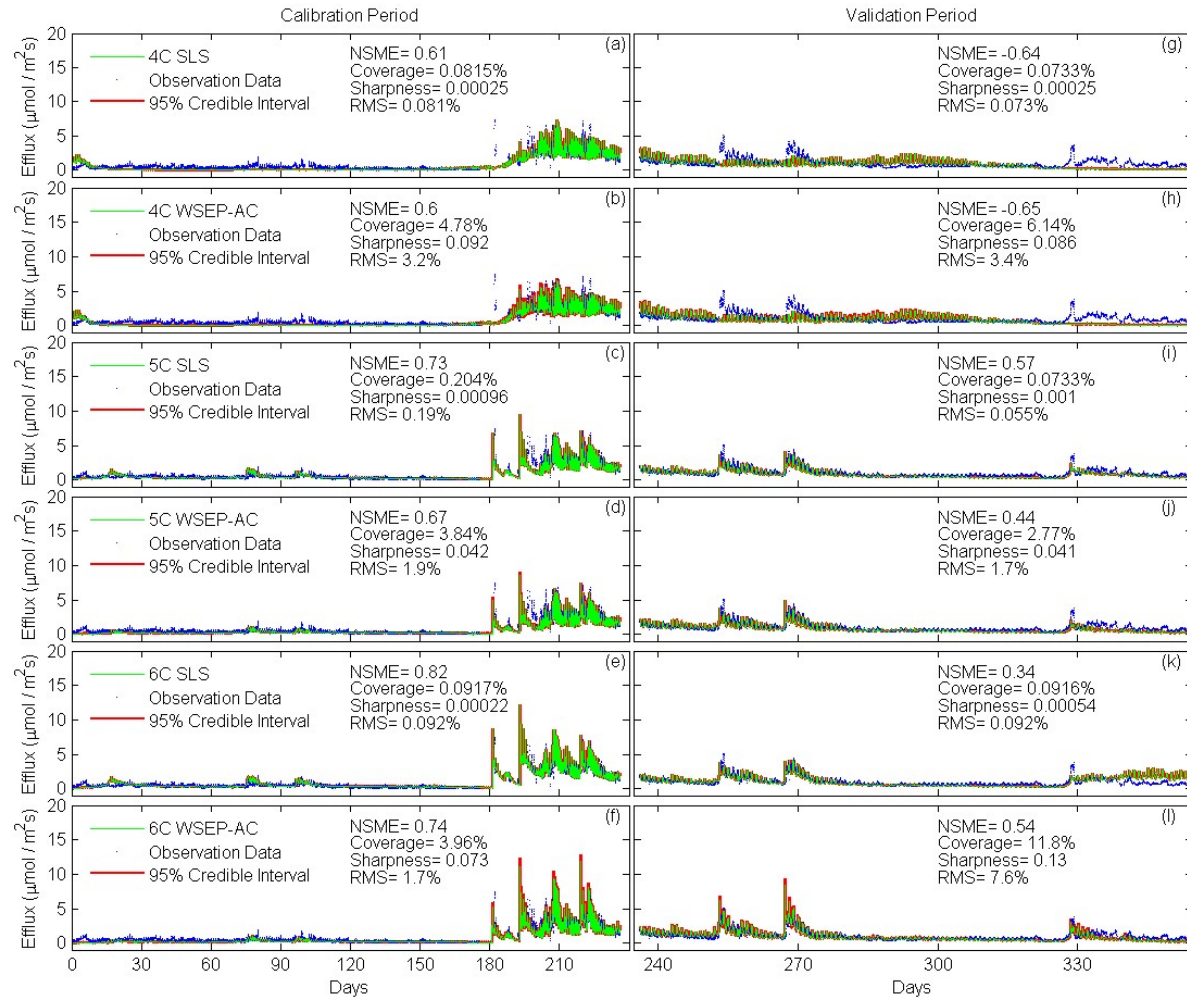


Figure 7. (a-b) Nash-Sutcliffe model efficiency (NSME), (c)-(d) sharpness, (e)-(f) predictive coverage, and (g)-(h) relative model score for measuring predictive performance of the three soil respiration models and the eight data models during the calibration and cross-validation periods. *The statistics are evaluated from the prediction ensembles generated to consider parametric uncertainty of the soil respiration models only.*

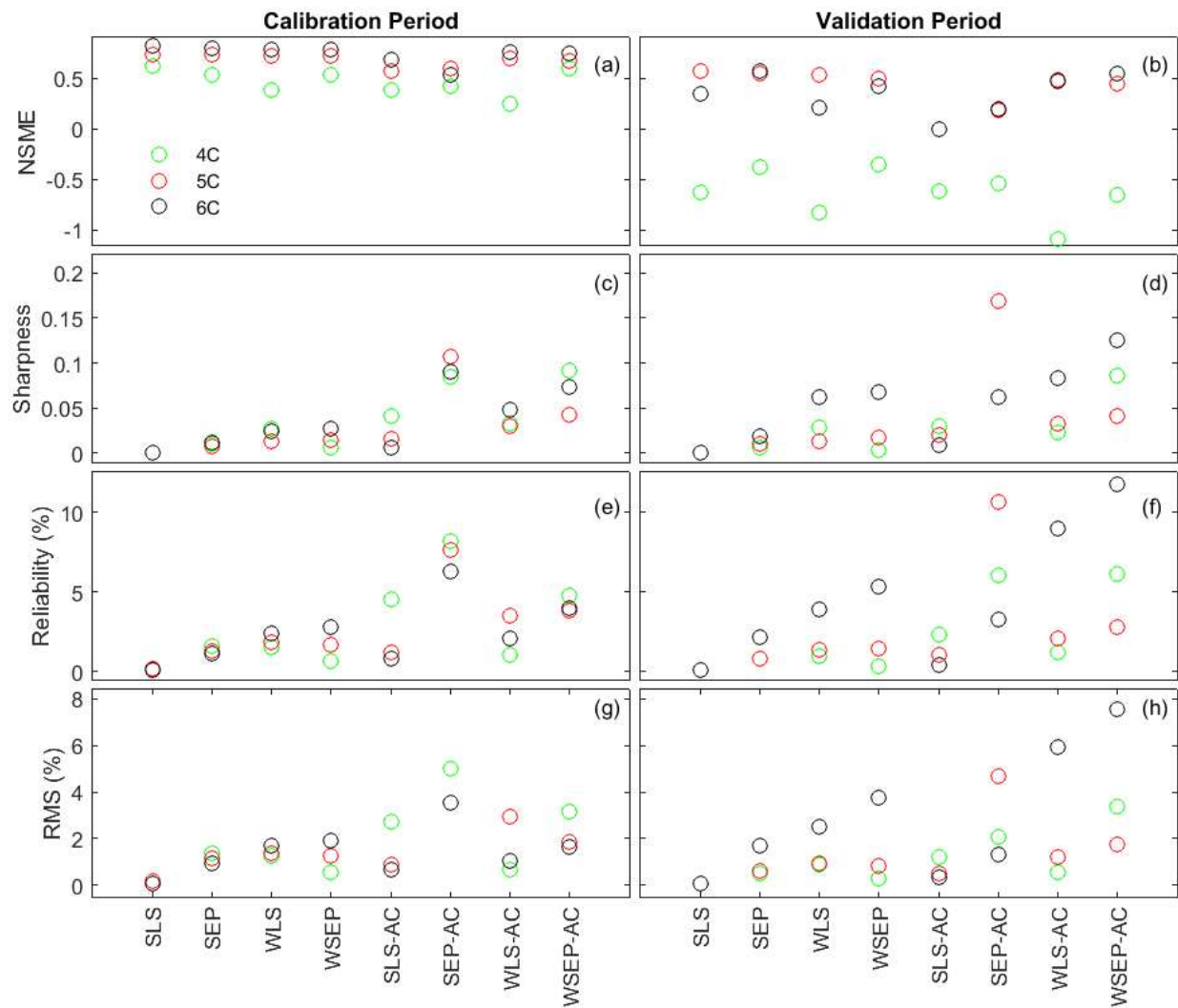


Figure 8. Observation data (blue dots) and mean prediction (green line) and 95% credible intervals (red line) of prediction ensembles for (a)-(f) the calibration period and (g)-(l) the validation period. The plots are for the three soil respiration models using data models SLS and WSEP-AC. *The prediction ensembles are generated to consider parametric uncertainty of not only the soil respiration models but also the data models.*

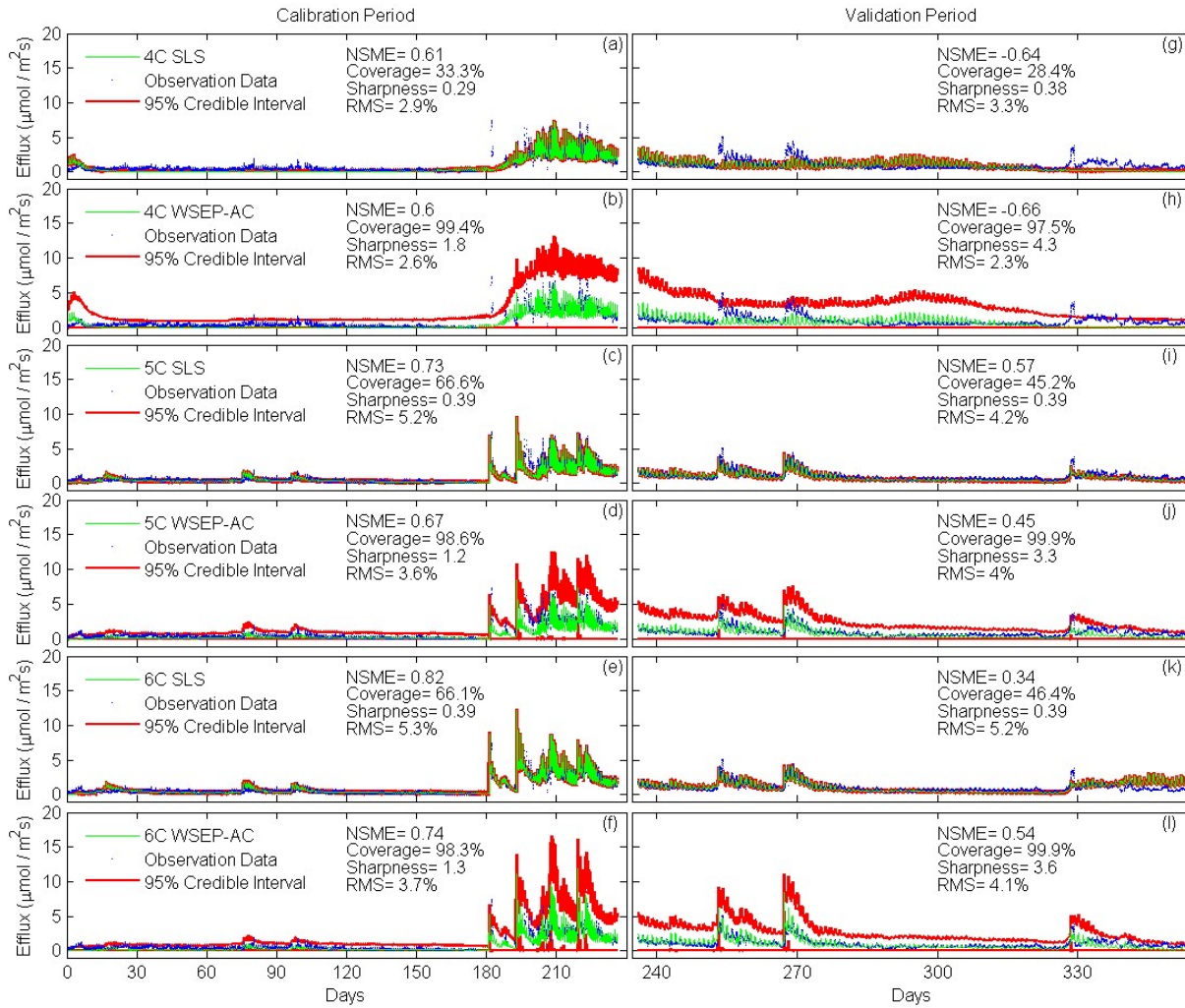


Figure 9. (a-b) Nash-Sutcliffe model efficiency (NSME), (c)-(d) sharpness, (e)-(f) predictive coverage, and (g)-(h) relative model score for measuring predictive performance of the three soil respiration models and the eight data models during the calibration and cross-validation periods. *The statistics are evaluated from the prediction ensembles generated to consider parametric uncertainty of not only the soil respiration models but also the data models.*

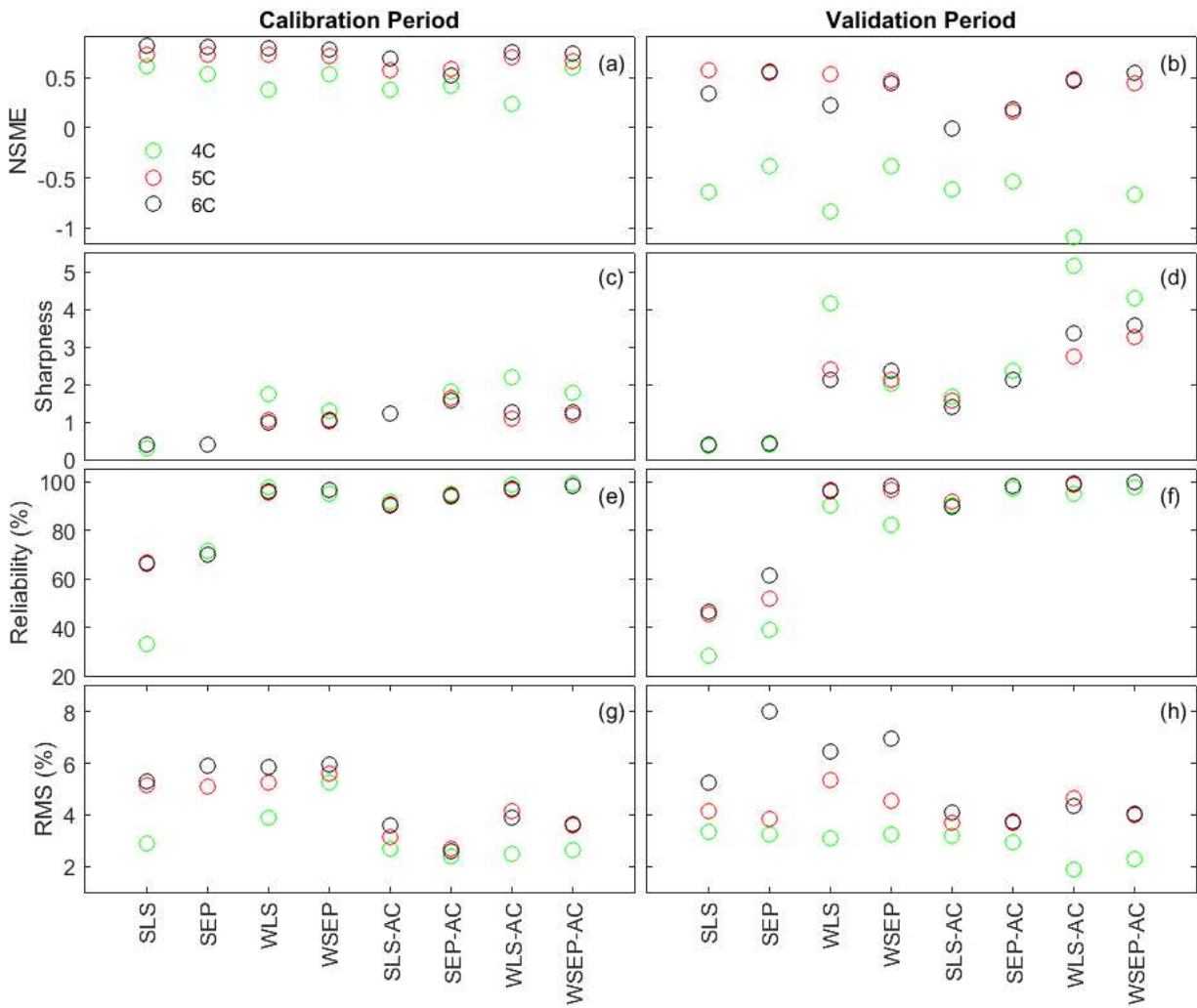
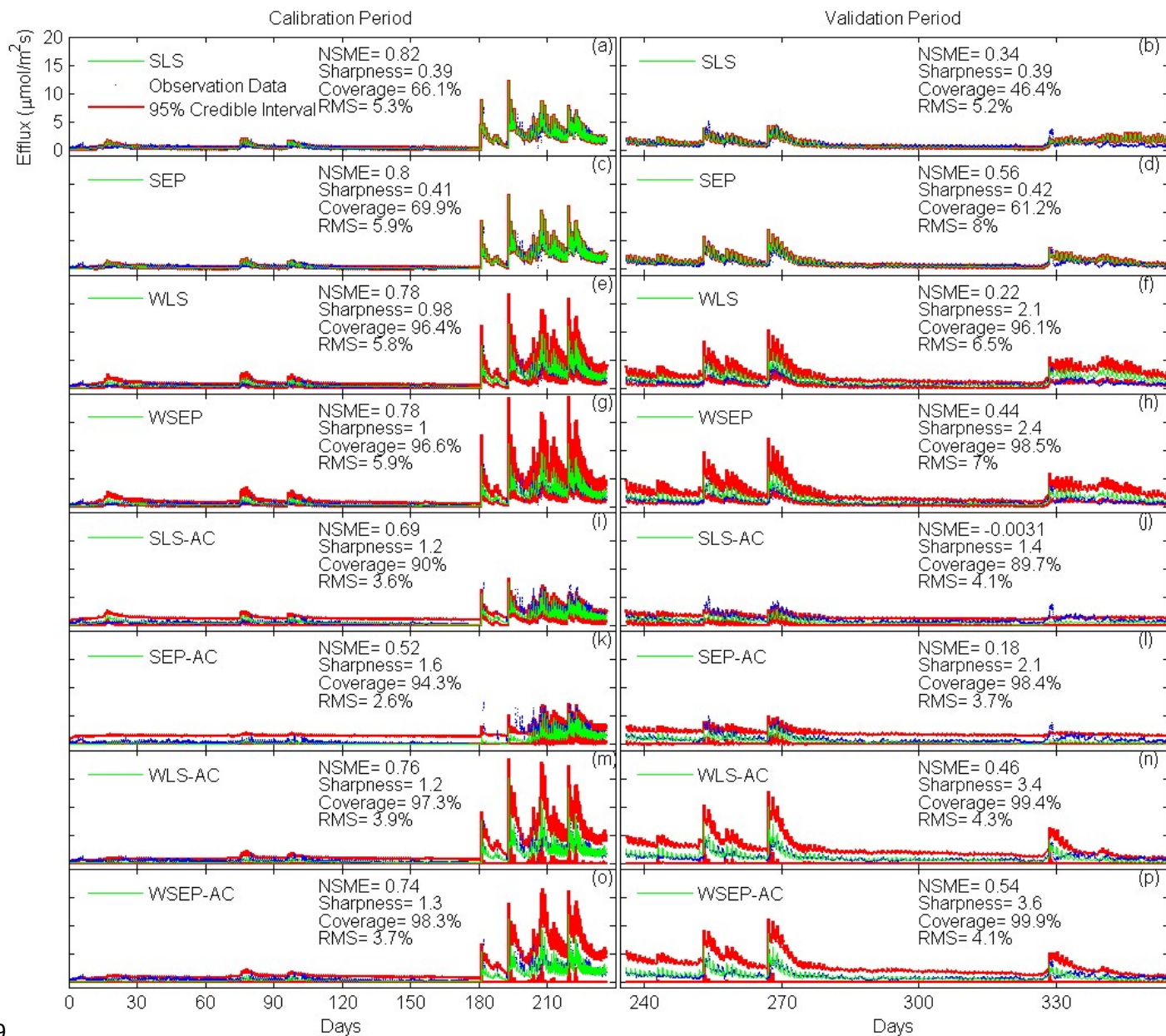


Figure 10. Observation data (blue dots) and mean prediction (green line) and 95% credible intervals (red line) for 6C for the eight likelihood functions during the calibration period (a)-(h) and the validation period (i)-(p). *The prediction ensembles are generated to consider parametric uncertainty of not only the soil respiration models but also the data models.*



Supplementary Table 1. Summary of the data models, their parameters, and corresponding likelihood functions

Residual Assumptions	Likelihood function	Data model	Residuals	Variance	<u>Likelihood function</u> <u>Data model</u> parameters
	<u>Generic data model</u>	Generic data model $a_t = \frac{\varepsilon_t}{\sigma_t} \quad a_t \sim X$	ε_t	σ_t	
Independent, normally distributed, and homoscedastic	Standard least square (SLS)	$a_t = \frac{\varepsilon_t}{\sigma_0} \quad a_t \sim N(0,1)$	$\varepsilon_t = d_t - Y_t$	$\sigma_t = \sigma_0$	Constant σ_0
Independent, and homoscedastic	Skew exponential power (SEP)	$a_t = \frac{\varepsilon_t}{\sigma_0} \quad a_t \sim SEP(0,1,\xi,\beta)$	$\varepsilon_t = d_t - Y_t$	$\sigma_t = \sigma_0$	Constant σ_0 Skewness ξ , Kurtosis β
Independent and normally distributed	Weighted least square (WLS)	$a_t = \frac{\varepsilon_t}{\sigma_0 + \sigma_1 Y_t} \quad a_t \sim N(0,1)$	$\varepsilon_t = d_t - Y_t$	$\sigma_t = \sigma_0 + \sigma_1 Y_t$	Heteroscedasticity model parameters σ_0, σ_1
Independent	Weighted skew exponential power (WSEP)	$a_t = \frac{\varepsilon_t}{\sigma_0 + \sigma_1 Y_t} \quad a_t \sim SEP(0,1,\xi,\beta)$	$\varepsilon_t = d_t - Y_t$	$\sigma_t = \sigma_0 + \sigma_1 Y_t$	Heteroscedasticity model parameters σ_0, σ_1 Skewness ξ , Kurtosis β
Normally distributed, and homoscedastic	Standard least square with auto-correlation (SLS-AC)	$a_t = \frac{\varepsilon_t - \sum_{i=1}^p \phi_i \varepsilon_{t-i}}{\sigma_0} \quad a_t \sim N(0,1)$	$\varepsilon_t - \sum_{i=1}^p \phi_i \varepsilon_{t-i}$	$\sigma_t = \sigma_0$	Constant σ_0 , Autoregressive model parameters ϕ_i
Homoscedastic	Skew exponential power with auto-correlation (SEP-AC)	$a_t = \frac{\varepsilon_t - \sum_{i=1}^p \phi_i \varepsilon_{t-i}}{\sigma_0} \quad a_t \sim SEP(0,1,\xi,\beta)$	$\varepsilon_t - \sum_{i=1}^p \phi_i \varepsilon_{t-i}$	$\sigma_t = \sigma_0$	Constant σ_0 , Autoregressive model parameters ϕ_i Skewness ξ , Kurtosis β
Normally distributed	Weighted least square with auto-correlation (WLS-AC)	$a_t = \frac{\varepsilon_t - \sum_{i=1}^p \phi_i \varepsilon_{t-i}}{\sigma_0 + \sigma_1 Y_t} \quad a_t \sim N(0,1)$	$\varepsilon_t - \sum_{i=1}^p \phi_i \varepsilon_{t-i}$	$\sigma_t = \sigma_0 + \sigma_1 Y_t$	Heteroscedasticity model parameters σ_0, σ_1 Autoregressive model parameters ϕ_i
	Generalized likelihood function (WSEP-AC)	$a_t = \frac{\varepsilon_t - \sum_{i=1}^p \phi_i \varepsilon_{t-i}}{\sigma_0 + \sigma_1 Y_t} \quad a_t \sim SEP(0,1,\xi,\beta)$	$\varepsilon_t - \sum_{i=1}^p \phi_i \varepsilon_{t-i}$	$\sigma_t = \sigma_0 + \sigma_1 Y_t$	Heteroscedasticity model parameters σ_0, σ_1 Autoregressive model parameters ϕ_i , Skewness ξ , Kurtosis β

Supplementary Figure 1. Workflow scheme

