

Black font: Reviewer comments

Black font: Author response

Blue font: Verbatim copy and paste from the revised manuscript

Anonymous Referee #1

The paper evaluates the impacts of statistical data assumptions in soil microbial respiration modeling on estimated model parameters and on model predictions. Inference is done using various soil respiration models and various likelihood functions, using half hourly CO₂ flux data from a field site. It's an interesting study, but I suggest additional effort to clarify and increase contribution of the work.

We are very thankful for the reviewer for talking the time to evaluate the manuscript, and for providing constructive comments.

1. Contribution: the authors should more clearly spell out the explicit contributions of the paper. On the one hand, the methodology is not new and has been developed and applied in hydrological studies. On the other hand, the application to CO₂ modeling may also not be entirely new since the likelihood approach used here has already been applied to ecological modeling (including carbon flux modeling); a recent example is Scholz, K., Hammerle, A., Hiltbrunner, E. et al. *Ecosystems* (2018) 21: 982. <https://doi.org/10.1007/s10021-017-0201-5>.

Response: We explicitly spelled out the novel contribution of this paper, which is the systematic evaluation of the impact of data model selection on Bayesian inference and predictive performance of soil respiration modeling with different degrees of model fidelity. We did a systematic review of Bayesian inference for soil respiration modeling. Most studies assume independent, Gaussian, and homoscedastic residuals. Few studies have relaxed these assumptions. However, only very few studies have focused on investigating the impacts 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 models on Bayesian inference and predictive performance of soil respiration modeling. The revised manuscript reads: "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 since the requirement of using an unknown probability model in Bayesian statistics is called "a basic dilemma" by 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 independent residuals 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). A recent study (Scholz et al., 2018) relaxed these three assumptions using the generalized likelihood function (Schoups and Vrugt, 2010). However, few studies have focused on investigating appropriateness and impact of these assumptions for soil

respiration modeling. This was performed 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.” In the first paragraph of the introduction we also stated “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.”

2. The authors find some problems with the estimation of autocorrelation and suggest an alternative approach (Evin et al.). Why not test this approach as well? I’m not sure this would warrant a separate publication. Including it here would enhance novelty of the paper in my opinion. Note also that the high temporal resolution (half hourly) of the data used by the authors may be a complicating factor; see the following paper that discusses this: <https://www.hydrol-earth-sci-discuss.net/hess-2018-406/>.

Thank you very much for bring our attention to this recent article of Ammann et al. (2018).

This manuscript provides a systematic evaluation of the impact of data model selection on Bayesian inference and predictive performance of soil respiration modeling. Figure 10 for example shows specific trends that would occur when relaxing the three assumptions of non-correlation, normality, and homoscedasticity using joint inversion approach, which has never been reported before in literature.

Autocorrelation is a complicated problem that we are currently working on. Joint inversion of heteroscedasticity and autocorrelation parameters can lead to poor predictive performance (Evin et al., 2013, 2014; Ammann et al. 2018; and this study). To address this problem a two-step procedure (e.g. Lu et al., 2013; Evin et al., 2013, 2014) is proposed. Our preliminary results show that using the sequential approach of Evin et al. (2013; 2014) by estimating the autoregressive parameters sequentially (after estimating the soil respiration model parameters and data-model parameters) did not solve this problem. Ammann et al. (2018) even states that the joint inversion is still preferred, and understanding the conditions where accounting for auto-correlation can be achieved remain poorly understood.

The problem of autocorrelation has several interlinked aspects that we would like to address in another manuscript. Auto-correlated errors might be attributed to a systematic error in the soil respiration model. The most obvious solution is to improve the soil respiration model. Otherwise, we can improve our data model. Our hypothesis that we would like to test is that omitting autocorrelation error through a filter approach (e.g. Schoups and Vrugt, 2010; Evin et al., 2013; 2014; this study) could be tricky as this leads to a loss of information content. Thus, joint approach may lead to biased parameter estimation (Figure 5) and poor predictive performance (Figure 10). While sequential approach would avoid the biased parameter estimation, but would still lead a poor predicative performance. Our current understanding is that this problem could emerge from several interlinked factors:

- Non-stationarity due to wet-dry periods as proposed by Ammann et al. (2018) could be a reason for this problem and thus accounting for non-stationarity (Smith et al., 2010b, Ammann et al. 2018) could alleviate this problem.
- The method for accounting for autocorrelation could have an impact. Autocorrelation could be addressed using a likelihood function based on covariance matrix of residuals $L(\mathbf{e})$ (e.g. Lu et al., 2013) with transformed residuals, and likelihood function of normalized residuals $L(\mathbf{a})$ (e.g. Schoups and Vrugt, 2010; Evin et al., 2013; 2014; this study) with autoregressive model that filter out autocorrelation. Note

that \mathbf{e} is a vector of transformed residuals, while \mathbf{a} is a vector independent and identically distributed random errors with zero mean and unit standard deviation. We would like to study these two methods.

- Joint versus sequential inversion for autocorrelation could also have an impact. Ammann et al. (2018) suggests that the joint inversion is still preferred over sequential inversion. This will be investigated under both L(e) and L(a) approaches. In addition, we would introduce a novel joint inversion procedure based on L(a) approach as follows. First, the parameters of the linear heteroscedastic model will be estimated similar to Schoups and Vrugt (2010) to remove heteroscedasticity. For each MCMC sample, after applying the linear heteroscedasticity model, the auto-correlation parameters can be deterministically calculated as internal variables of the data model similar to Lu et al. (2013) and not as calibration parameters as in Schoups and Vrugt (2010). This is mainly to avoid interaction between heteroscedasticity and autocorrelation parameters. The auto-correlation parameters can be calculated following Lu et al. (2013). We have revised the manuscript to further clarify these issues. The revised manuscript reads [“This study confirms the empirical findings and theoretical analysis \(Evin et al., 2013; 2014; Ammann et al. 2018\) that separate accounting for autocorrelation or joint inversion of correlation and heteroscedasticity can be problematic. By drawing on similarity from surface hydrology, the study of Ammann et al. \(2018\) suggests that this might be attributed to non-stationarity due to wet-dry periods with half-hourly data. 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.”](#)

3. The paper should be checked for various grammatical errors and typos. One example is "heteroscedasticity", which is spelled in multiple creative ways throughout the paper.

Response: Thank you very for pointing this out and we have corrected "heteroscedasticity" at eight different locations throughout the manuscript. We corrected several other grammatical errors and typos.

4. Description of the various evaluation metrics seems better placed in the methods than results section.

We moved the description of the various evaluation metrics from the results to the methods section.

5. Terminology: the distinction between model fidelity and discrepancy is not clear

We clarified these two terms as follows: [“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.”](#)

6. Line 305, "discrete proposal distribution": I don't think the proposal is discrete, it is a proposal distribution over a continuous parameter space.

Response: We revised "discrete proposal distribution" to [“adaptive proposal distribution.”](#)

7. Line 477: please rephrase; I don't think it's "expected" that accounting for autocorrelation leads to biased parameter values. I would expect the opposite, since autocorrelation provides a (simple) way to account for model errors.

Response: We rephrased this sentence to "First, we obtained biased parameter estimates that is out the reasonable physical range."

8. Eq. 23: is index i an index over time or is it an ensemble index? Please clarify.

Thank you very much for point this out. We clarified that this is an ensemble prediction Y_{ij} where i is index over time, and revised other parts of the manuscript accordingly. The new sentence read "the ensemble prediction Y_{ij} is similar to Y_i above where i is index over time and specific to the j -th combination."

9. Line 598: approaches that use "total residual error" typically still separate out parametric uncertainty, so the residual error includes measurement, model input, and model structure uncertainty, but not parameter uncertainty.

That is true. We rephrased that sentence to "total residuals that separates out parametric uncertainty, so the residual error includes measurement, model input, and model structure uncertainty."

Thank you very much for your constructive comments.

Anonymous Referee #2

The manuscript submitted by Elshall et al. is an interesting study dealing with the complexity of soil C model parameterization. In recent decades, the complexity of those model as well as the different tools to parameterize has increased substantially leading to potential misuses of powerful but complex mathematical approaches. The goal of Elshall et al is therefore to evaluate the impact on process-based model predictions of neglecting a couple of assumptions of the Bayesian framework as it is often done by soil modelers to avoid complexity.

We thank the reviewer very much evaluating the manuscript and for providing constructive feedback and suggestions.

The present study might not be super novel for the entire modeling communities in geoscience as mentioned by the other referee. Nevertheless, it underlines a flaw of several carbon soil modeling studies and might be considered as novel in this context. It is a pity that the author may not freely communicate their models and scripts it would have definitely increased the impact of the paper.

We feel sorry for this too, and we would love to share the code and the soil respiration models upon request.

Even though the objectives of the paper are important and deserve to be published, in my opinion, the manuscript in its present form is sometimes too hard to read and needs some simplifications. A first recommendation might be to have a table summarizing all the acronyms and try to reduce them when not necessary.

We added a list of acronyms as follows:

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

Secondly, a workflow scheme might also be useful to understand the logic of the authors, which is not always super clear.

We added a summary table of the data models and corresponding likelihood functions. The revised manuscripts states [“A summary table of the eight data models with corresponding parameters is provided in the supplementary materials.”](#) We added a workflow scheme as a supplementary figure. The revised manuscript reads [“ Our workflow scheme is presented in the supplementary materials.”](#) The new table and figure are presented below and in the attached supplementary file.

Finally, I missed some definition to be sure I fully understood the text. In particular, it is not crystal clear to me what the author means by 'data model'. From my understanding, a data model is based on data but the observed data are presented quite fare from the data model.

In the revised manuscript we clarified that [“A data model that is also known as a residuals model or an error model is used to characterize residuals \(i.e., the difference between data and corresponding model simulations\).”](#) In addition, please see our response to the previous comment.

Another point is that I still do not fully understood how the authors link their data model with their process-based model. I understood that the data models are used for posterior parameter estimation but sometimes the text makes me doubt.

The parameters of the data model are jointly estimated with the parameters of the soil respiration model using MCMC. We clarified this in the revised manuscript [“the posterior distributions of the data model parameters are jointly estimated with the soil respiration model parameters using the MT-DREAM\(ZS\) code \(Laloy and Vrugt, 2012\).”](#) In addition, a summary of the data model parameters is presented in the supplementary materials as we clarified in a previous response.

I don't understand why the author fixed the upper limit of the physical range of CUE to 0.6 (the mean over terrestrial systems) whereas in the paper they cited several observations are above 0.6

The thermodynamic maximum limit of CUE is 0.6 and the empirical observations show that CUE over a wide range of field conditions converges to ~ 0.30 with a mean value of 0.55 for terrestrial ecosystems (Sinsabaugh et al., 2013). We used this upper limit for analysis only. We did not fix this limit for Bayesian inverse modeling to understand the impact of data model on parameter estimation.

Some typo: l121 'and' not necessary L176 please correct the parenthesis L611: despite instead of desp8ite

Thank you very much for pointing out these typos and we corrected them.

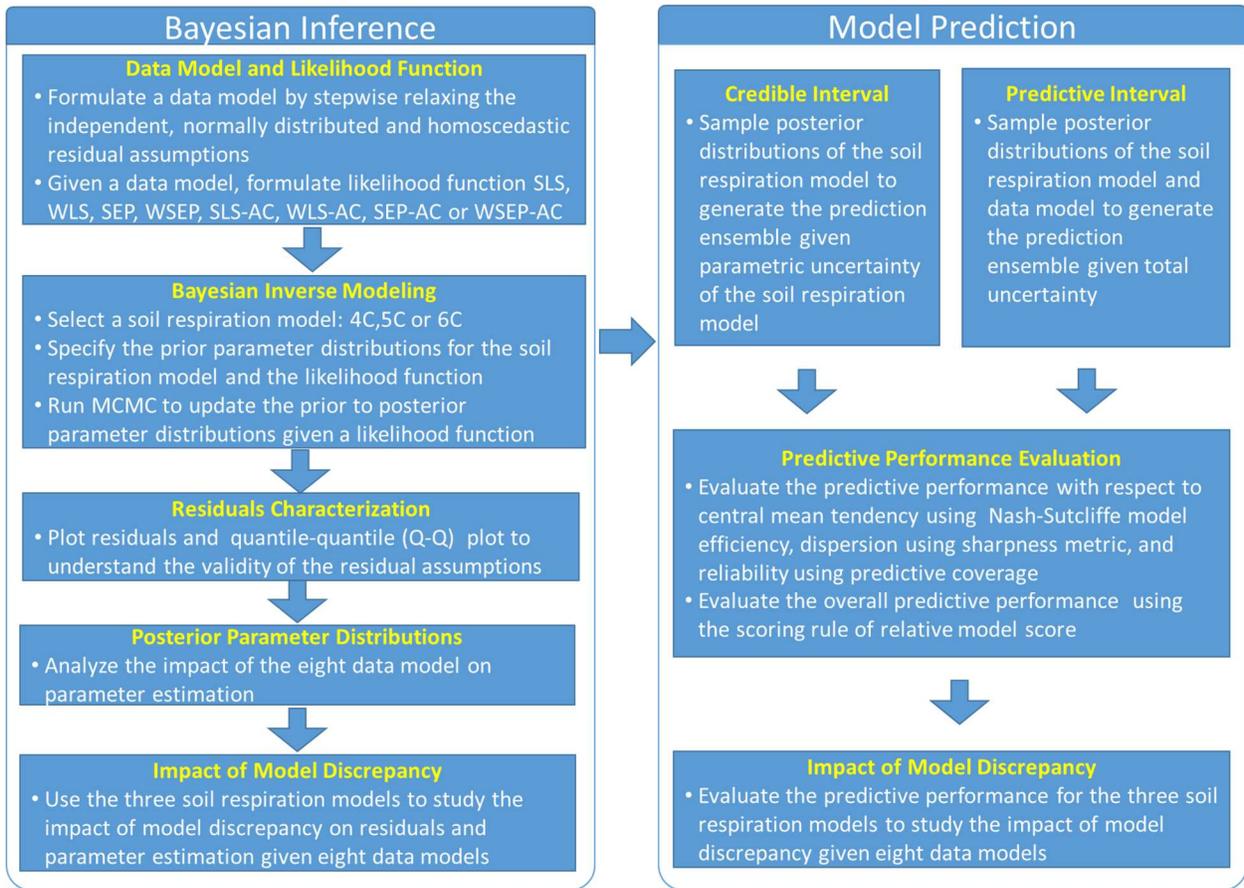
I, therefore, think that this manuscript deserves publication after a deep rewriting to clarify the methods used

Addressing the review comments helped us to rewrite and clarify several parts of the manuscript. Thank you very much.

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

Residual Assumptions	Likelihood function	Data model	Residuals	Variance	Likelihood function parameters
		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



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

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

5
6 ¹ Department of Earth Sciences, University of Hawai‘i Manoa, Honolulu, Hawaii, USA

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

8 ³ Department of Scientific Computing, Florida State University, Tallahassee, Florida

9 ⁴ Biosphere 2, University of Arizona, Tucson, Arizona

10 ⁵ Department of Hydrology and Water Resources, University of Arizona, Tucson, Arizona

11 ⁶ School of Geography and Development, University of Arizona, Tucson, Arizona

12
13
14 *Corresponding Author: Ming Ye, Telephone: (850) 644-4587, Email: mye@fsu.edu

15
16
17 Submitted for publication in Geoscientific Model Development

18
19 ~~October~~ March, 2019~~8~~

47 **Key Points**

48

49 (1) Bayesian inference and prediction are useful to evaluate multiple soil respiration models
50 with different levels of model complexity.

51 (2) Data models used in Bayesian inference have substantial impacts on model parameter
52 distributions and subsequently model predictions.

53 (3) Using exponential power distribution and considering heteroscedasticity in data models
54 improves Bayesian inference and prediction.

55

56

57

58

59

60

61

62

63

64

65

66

67 Keywords: Soil respiration, Bayesian, likelihood function, data model, autocorrelation,
68 heteroscedasticity, skew exponential power distribution, cross-validation, relative model score

69

70

71 **Abstract**

72 Bayesian inference of microbial soil respiration models is often based on the assumptions that the
73 residuals are independent (i.e. no temporal or spatial correlation), identically distributed (i.e.
74 Gaussian noise) and with constant variance (i.e. homoscedastic). In the presence of model
75 discrepancy, since no model is perfect, this study shows that these assumptions are generally
76 invalid in soil respiration modeling such that residuals have high temporal correlation, an
77 increasing variance with increasing magnitude of CO₂ efflux, and non-Gaussian distribution.
78 Relaxing these three assumptions stepwise results in eight data models. Data models are the basis
79 of formulating likelihood functions of Bayesian inference. This study presents a systematic and
80 comprehensive investigation of the impacts data model selection on Bayesian inference and
81 predictive performance. We use three mechanistic soil respiration models with different levels of
82 model fidelity (i.e. model discrepancy) with respect to number of carbon pools and explicit
83 representations of soil moisture controls on carbon degradation, and accordingly have different
84 levels of model complexity with respect to the number of model parameters. The study shows data
85 models have substantial impacts on Bayesian inference and predictive performance of the soil
86 respiration models such that: (i) the level of complexity of the best model is generally justified by
87 the cross-validation results for different data models; (ii) not accounting for heteroscedasticity and
88 autocorrelation might not necessarily result in biased parameter estimates or predictions, but will
89 definitely underestimate uncertainty; (iii) using a non-Gaussian data model improves the parameter
90 estimates and the predictive performance; and (iv) separate accounting for autocorrelation or joint
91 inversion of correlation and heteroscedasticity can be problematic and requires special treatment.
92 Although the conclusions of this study are empirical, the analysis may provide insights for
93 selecting appropriate data models for soil respiration modelings.

94 1 Introduction

95 Developing accurate soil respiration models is important for realistic projection of global
96 carbon [C] cycle, as global soils store 2,300Pg carbon, an amount more than 3 times that of the
97 atmosphere (Schmidt et al., 2011) and release 60–75 Pg C/yr, about 7 times more CO₂ to the
98 atmosphere than all human-caused emissions (Le Quéré et al., 2014). The major work on soil
99 respiration modeling has been focused on advancing knowledge about model inputs and
100 calibration data (e.g. Janssens et al., 2003; Peters et al., 2007; Scott et al., 2009; Barron-Gafford et
101 al., 2011; Hilton et al., 2014) and on developing more advanced models for better representing
102 soil microbial processes (e.g. Schimel and Weintraub, 2003; Allison et al., 2010; Davidson et al.,
103 2011; Wieder et al., 2013, 2015; Xu et al., 2014; Zhang et al., 2014) . Integration of data and
104 models is indispensable for improving predictability of the terrestrial carbon cycle, and statistical
105 modeling is a vital tool for the model-data integration (Luo et al., 2011, 2014; Wieder et al., 2015).
106 In addition, use of state-of-the-art statistical methods is necessary to accurately quantify
107 uncertainty in parameters and structures of soil respiration models for improvement and practical
108 uses of the models (Katz et al., 2013). ~~Statistical modeling always requires adequately~~A data model
109 that is also known as a residuals model or an error model-characterizing is used to characterize
110 residuals residuals(,i.e., the difference between data and corresponding model simulations). While
111 a large number of data models have been used (e.g. (Elshall et al., 2018a; Scholz et al., 2018); to
112 our knowledge, comprehensive and systematic evaluation of data models for soil respiration
113 ~~models-modeling~~ has not been reported in literature.

114 The goal of this study is to evaluate the impacts of data models on Bayesian inference and
115 predictive performance of three mechanistic soil respiration models, and use these findings to
116 make broader recommendations. The three models were developed by Zhang et al. (2014) to

117 simulate the Birch effect (the peak soil microbial respiration pulses in response to episodic rainfall
118 pulses) at a site scale and a short temporal scale, which are important for gaining mechanistic
119 understanding of CO₂ efflux production (Högberg and Read, 2006; Vargas et al., 2011). Zhang et
120 al. (2014) developed a total five models, including an existing four-carbon pool model and four
121 new models with additional carbon pools and/or explicit representations of soil moisture controls
122 on carbon degradation and microbial uptake rates. The models Zhang et al. (2014) were calibrated,
123 and Bayesian model selection was used to select ~~and~~ the best model. However, this effort was
124 based on a single data model. It is unknown whether the best model still remains the best (in terms
125 of reproducing the both calibration data and the cross-validation data) if a different data model is
126 used. In addition, since predictive performance of the models was not evaluated in Zhang et al.
127 (2014), it is unknown whether the best model will give the best predictions. These two questions
128 are addressed in this study by considering eight data models and by evaluating predictive
129 performance in a manner of cross-validation. The top two models (also the two most high fidelity
130 models) ranked by Zhang et al. (2014) are considered in this study, and the worst model (also the
131 low fidelity model) is also considered in this study for comparison. We use the terms model fidelity
132 and model discrepancy interchangeably. Model fidelity refers to the degree of realism of
133 representing our scientific knowledge with respect to the real world system. That is - a ~~That is~~ high
134 fidelity model has~~model with~~ less discrepancy. Conducting Bayesian inference and evaluating
135 predictive performance for the three models with different degrees of fidelity provides more
136 insights than for a single model.

137 Bayesian inference in general uses the Bayes' theorem to update the distributions of model
138 parameters to posterior parameter distributions given a likelihood function. The mathematical
139 formulation of the (formal and informal) likelihood function requires a probabilistic data model

140 that however is intrinsically unknown due to unknown errors in all model components such as
141 observation data, model structures, parameters, and driving forces. Bayesian inference of soil
142 respiration models often adopts the assumption of independent, normally distributed and
143 homoscedastic residuals (e.g. Ahrens et al., 2014; Bagnara et al., 2015, 2018; Barr et al., 2013;
144 Barron-gafford et al., 2014; Braakhekke et al., 2014; Braswell et al., 2015; Correia et al., 2012; Du
145 et al., 2015, 2017; Hararuk et al., 2014; Hashimoto et al., 2011; He et al., 2018; Klemedtsson et
146 al., 2008; Menichetti et al., 2016; Raich et al., 2002; Ren et al., 2013; Richardson and Hollinger,
147 2005; Steinacher and Joos, 2016; Tucker et al., 2014; Tuomi et al., 2008; Xu et al., 2006; Yeluripati
148 et al., 2009; Yuan et al., 2012, 2016; Zhang et al., 2014; Zhou et al., 2010). These assumptions are
149 conveniently adopted since the requirement of using an unknown probability model in Bayesian
150 statistics is called “a basic dilemma” by Box and Tiao (1992). Postulating the data models is always
151 based on assumptions about residual statistics, and the most widely used assumptions are paired
152 as follows: (i) independent vs. correlated residuals, (ii) homoscedastic vs. heteroscedastic
153 residuals, and (iii) Gaussian vs. non-Gaussian residuals. For soil respiration modeling few studies
154 have relaxed the non-correlation assumption(e.g. Cable et al., 2008, 2011; Li et al., 2016b), the
155 homoscedasticity assumption(e.g. Berryman et al., 2018; Elshall et al., 2018; Ogle et al., 2016;
156 Tucker et al., 2013), and the non-Gaussian and homoscedasticity assumptions (e.g. Elshall et al.,
157 2018; Ishikura et al., 2017; Kim et al., 2014). A recent study (Scholz et al., 2018) relaxed these
158 three assumptions using the generalized likelihood function (Schoups and Vrugt, 2010). There are
159 ~~many diagnostics available to assess these choices (a number of them is used in this paper).~~
160 However, few studies have focused on investigating appropriateness and impact of these
161 assumptions for soil respiration modeling, by relaxing the independent residuals assumption (
162 Ricciuto et al., 2011) and the Gaussian residuals assumption (Ricciuto et al., 2011; van Wijk et al.,

163 2008). By relaxing these three assumptions stepwise resulting in eight data models, to our
164 knowledge this is the first study that systematically evaluates the impact of data model selection
165 on Bayesian inference and predictive performance of soil respiration modeling. In addition, to our
166 knowledge this is the first soil respiration modeling study that investigates the impact of data
167 models in relation to model fidelity.

168 Relaxing these three assumption results in eight data models, which are shown in details in
169 Section 2. For example, combining the assumptions of independent, homoscedastic, and Gaussian
170 residuals leads to the standard least squares data model. This model is the simplest one among the
171 eight data models, since it requires only one parameter, i.e., the constant variance of the Gaussian
172 distribution. Note that there is a difference between the physical-soil respiration model parameters
173 and data model parameters. They technically can be estimated together, but one arises from
174 assumptions about soil respiration processes, and the other assumptions about the residuals data
175 models. Relaxing the homoscedastic assumption to heteroscedastic gives the weighted least
176 squares data model. It is more complex, because it has extra parameters to account for-it requires
177 multiple variances for multiple data. Whenever one or combinations of the three assumptions
178 (independence, homoscedasticity, and normality) are relaxed, the resulting data models become
179 more complex and require more parameters. Such This sSsystematic evaluation-way-of of
180 formulating data models ((McInerney et al., 2017; McInerney et al., 201; is similar to that of Smith
181 et al. (2010b, 2015),-and it is necessary to evaluate appropriateness of residualsthe three basic
182 assumptions and their impacts on Bayesian inference.

183 The assumptions of heteroscedastic, correlated, and non-Gaussian residuals are accounted for
184 using the method of Schoups and Vrugt (2010) in the following procedure: (i) the correlation is
185 removed from the residuals by using an autoregressive model; (ii) the resulting residuals are

186 normalized by a linear model of variance; and (iii) the normalized residuals are characterized by
187 using the skew exponential power distribution. The data model parameters (i.e., coefficients of the
188 autoregressive model, the linear variance model, and the skew exponential power distribution) are
189 not specified by users, but estimated together with ~~physical-soil respiration~~ model parameters
190 during the Bayesian inference. The skew exponential power distribution is general in that by
191 adjusting the values of its kurtosis and skewness parameters the distribution can produce other
192 distributions such as the Laplace distribution ~~used by~~ (van Wijk et al., 2008; Ricciuto et al.,
193 ~~2011~~) and ~~(Ricciuto et al., 2011)~~, and other distributions ~~through given by~~ using an different
194 ~~kurtosis parameters of an~~ exponential model with different kurtosis parameters (Tang and Zhuang,
195 2009). It is worth pointing out that there exist other methods to account for the three assumptions.
196 Evin et al. (2013) suggested accounting for residual heteroscedasticity before accounting for
197 residual autocorrelation. Lu et al. (2013) developed an iterative two-stage procedure to separately
198 estimate physical model parameters and data model parameters. Evin et al. (2014) developed a
199 similar procedure to first estimate model parameters and then estimate heteroscedasticity and
200 autocorrelation parameters. While this study uses the method of Schoups and Vrugt (2010),
201 exploring other methods is warranted in future studies.

202 After investigating the impacts of the data models on Bayesian inference, this study evaluates
203 the impacts of the data models on predictive performance of the three soil respiration models.
204 Using random samples generated during the Bayesian inference, a prediction ensemble is produced
205 for each soil respiration model. The ensemble is used to evaluate predictive performance of the
206 models in a stochastic sense by estimating to what extent the models can predict future events. The
207 evaluation in this study is done in a cross-validation manner ~~by~~ splitting a the dataset of CO₂
208 efflux into two parts for Bayesian inference and cross-validation, respectively. The evaluation of

209 predictive performance is important because different data models may give different parameter
210 distributions and accordingly different predictive performance. For example, the study of van Wijk
211 et al. (2008) concluded that the choice of the residual function is crucial to achieve accurate model
212 prediction and parameter estimation. Shi et al. (2014) showed that the posterior parameter
213 distributions and predictive performance given by two data models (weighted least square and
214 skew exponential power distribution after removing heteroscedasticity and autocorrelation) are
215 dramatically different, and a definitive conclusion was drawn that one data model is better than
216 the other. The evaluation of predictive analysis is conducted for the following two cases: (1) the
217 prediction ensemble is generated by random samples of the soil respiration models only (i.e.
218 credible interval), and (2) the prediction ensemble is generated by random samples of not only the
219 soil respiration models but also the data models (i.e. predictive interval). The two cases lead to
220 different conclusions about the predictive performance. It is expected that the evaluation of
221 predictive performance conducted in this study can help select the most appropriate data model to
222 achieve optimal model predictions.

223 The remainder of the paper is organized as follows. Section 2 starts with a description of the
224 evolving data models and their corresponding likelihood functions used in Bayesian inference,
225 followed by a brief summary of the three soil respiration models. The results of Bayesian inference
226 are discussed in Section 3 and Section 4, addressing the data model implications on parameter
227 estimation and predictive performance, respectively. Section 5 summarizes the key findings and
228 limitations of this study, and provides recommendations for approaching data model selection.

229 2 Methodology

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

236 ~~21 Methodology~~

237 ~~This section starts with a descriptions of the eight data models that account for the three pairs~~
238 ~~of assumptions about residuals in a stepwise manner in Section 2.1. The data models are used to~~
239 ~~build the likelihood functions used in Section 2.2 for Bayesian inference. The three soil respiration~~
240 ~~models and observations of CO₂ efflux are described in Sections 2.3 and 2.4, respectively.~~

241 **2.1 Data models**

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

$$245 \quad a_t = \frac{\varepsilon_t}{\sigma_t} \quad a_t \sim X, \quad (1)$$

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

$$250 \quad a_t = \frac{\varepsilon_t}{\sigma_0} \quad a_t \sim N(0,1), \quad (2)$$

251 where $\sigma_t = \sigma_0$ is a constant for all the data (i.e., homoscedasticity), and X is the standard normal
 252 distribution, $N(0,1)$. The unknown parameter σ_0 is estimated jointly with unknown physical
 253 model parameters. If σ_t is not a constant (i.e., heteroscedasticity), SLS becomes the
 254 weighted least squared (WLS) data model. While heteroscedasticity can be accounted for through
 255 residuals transformation (e.g. Thiemann et al., 200; Smith et al., 2010b) or other similar approaches
 256 (Gragne et al., 2015) a linear heteroscedastic model $\sigma_t = \sigma_0 + \sigma_1 Y_t$ is assumed following other
 257 studies (Thyer et al., 2009; Schoups and Vrugt, 2010; Evin et al., 2013, 2014). With the linear
 258 model, there is no need to estimate σ_t for each data. Instead, σ_t is calculated by estimating only
 259 two parameters, σ_0 and σ_1 . The WSL data model is written as

$$260 \quad a_t = \frac{\varepsilon_t}{\sigma_0 + \sigma_1 Y_t} \quad a_t \sim N(0,1). \quad (3)$$

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

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

$$270 \quad a_t = \frac{\varepsilon_t - \sum_{i=1}^p \phi_i \varepsilon_{t-i}}{\sigma_0} \quad a_t \sim N(0,1) \quad (4)$$

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

$$274 \quad a_t = \frac{\varepsilon_t - \sum_{i=1}^p \phi_i \varepsilon_{t-1}}{\sigma_0 + \sigma_1 Y_t} \quad a_t \sim N(0,1) \quad (5)$$

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

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

$$280 \quad p(a_t | \xi, \beta) = \frac{2\sigma_\xi}{\xi + \xi^{-1}} \omega_\beta \exp\left[-c_\beta |a_{\xi,t}|^{2/(1+\beta)}\right], \quad (6)$$

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

$$282 \quad a_{\xi,t} = (\mu_\xi + \sigma_\xi a_t) / \xi^{\text{sign}(\mu_\xi + \sigma_\xi a_t)}, \quad \mu_\xi = M (\xi - \xi^{-1}), \quad \omega_\beta = \frac{\Gamma^{1/2}[3(1+\beta)/2]}{(1+\beta)\Gamma^{3/2}[(1+\beta)/2]},$$

$$283 \quad \sigma_\xi = \sqrt{(1-M^2)(\xi^2 + \xi^{-2}) + 2M^2 - 1}, \quad M = \frac{\Gamma[1+\beta]}{\Gamma^{1/2}[3(1+\beta)/2]\Gamma^{1/2}[(1+\beta)/2]}, \quad \text{and}$$

$$284 \quad c_\beta = \left(\frac{\Gamma[3(1+\beta)/2]}{\Gamma[(1+\beta)/2]}\right)^{1/(1+\beta)}$$

are derived variables of β and ξ , and $\Gamma[.]$ is the gamma function. The

285 kurtosis parameter $\{\beta \in \mathbb{R} : -1 \leq \beta \leq 1\}$ determines the peakness of the pdf such that the β values
 286 of -1, 0, and 1 give uniform, Gaussian and Laplace distributions, respectively. The skewness

287 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,
 288 1, and 10 give positively skewed, symmetric, and negatively skewed distributions, respectively.
 289 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
 290 skew exponential power distribution $SEP(0,1,\xi=1,\beta=0)$ becomes the standard normal
 291 distribution,

$$292 \quad p(a_t | \xi=1, \beta=0) = \frac{1}{\sqrt{2\pi}} \exp\left[-\frac{1}{2}(a_t)^2\right]. \quad (7)$$

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

294 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
 295 SEP, WSEP, SEP-AC, and WSEP-AC as follows,

$$296 \quad a_t = \frac{\varepsilon_t}{\sigma_0} \quad a_t \sim SEP(0,1,\xi,\beta) \quad (8)$$

$$297 \quad a_t = \frac{\varepsilon_t}{\sigma_0 + \sigma_1 Y_t} \quad a_t \sim SEP(0,1,\xi,\beta). \quad (9)$$

$$298 \quad a_t = \frac{\varepsilon_t - \sum_{i=1}^p \phi_i \varepsilon_{t-i}}{\sigma_0} \quad a_t \sim SEP(0,1,\xi,\beta) \quad (10)$$

$$299 \quad a_t = \frac{\varepsilon_t - \sum_{i=1}^p \phi_i \varepsilon_{t-i}}{\sigma_0 + \sigma_1 Y_t} \quad a_t \sim SEP(0,1,\xi,\beta) \quad (11)$$

300 In comparison with the Gaussian data models, the SEP-based data models have two more
 301 parameters (ξ and β) to be estimated jointly with physical model parameters. WSEP-AC data
 302 model, which is known as the generalized likelihood function, is the most commonly used SEP-
 303 based data model (e.g. Vrugt and Ter Braak, 2011; Hublart et al., 2016; [Scholz et al., 2018](#)). [A](#)

304 [summary table of the eight data models with corresponding parameters is provided in the](#)
 305 [supplementary materials.](#)

306 2.2 Bayesian inference and likelihood functions

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

$$312 \quad p(\boldsymbol{\theta}|\mathbf{d}) = \frac{p(\mathbf{d}|\boldsymbol{\theta})p(\boldsymbol{\theta})}{\int p(\mathbf{d}|\boldsymbol{\theta})p(\boldsymbol{\theta})d\boldsymbol{\theta}} \quad (12)$$

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

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

$$322 \quad p(\mathbf{d}|\boldsymbol{\theta}) = p(\boldsymbol{\varepsilon}_t|\boldsymbol{\theta}) = \prod_{t=1}^n \sigma_t^{-1} \frac{2\sigma_\xi}{\xi + \xi^{-1}} \omega_\beta \exp\left(-c_\beta |a_{\xi,t}|^{2/(1+\beta)}\right). \quad (13)$$

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

325 $\omega_\beta = 1/\sqrt{2\pi}$, $c_\beta = 1/2$, and $a_{\xi,t} = a_t$, equation (13) becomes the likelihood function corresponding
326 to the SLS data model. Replacing $\sigma_t = \sigma_0$ by $\sigma_t = \sigma_0 + \sigma_1 E_t$, equation (13) becomes the likelihood
327 function of the WLS data model.

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

336 2.3 Soil respiration models

337 Zhang et al. (2014) studied the Birch effect (the peak soil microbial respiration pulses in
338 response to episodic rainfall pulses), and developed five models, evolving from an existing four-
339 carbon pool model to models with additional carbon pools and/or explicit representations of soil
340 moisture controls on carbon degradation and microbial uptake rates. Three of the five models are
341 used in this study, and they are denoted as 4C, 5C, and 6C. Note that model 4C is model 4C_NOSM
342 of Zhang et al. (2014), not ~~their~~ model 4C. Figure 1 is the diagram of model 6C, the most complex
343 one among the five models. The simplest one, model 4C, has four carbon pools, i.e., soil organic
344 carbon (SOC), dissolved organic carbon (DOC), microbial biomass (MIC), and enzymes (ENZ),
345 and does not consider the soil moisture control on carbon degradation and microbial uptake rates.
346 Models 5C and 6C has an explicit representation of soil moisture controls on the rates. Based on

347 the dual Arrhenius and Michaelis–Menten kinetics model, the original SOC degradation rate,
 348 V_{decom} , is (Davidson et al., 2011; Davidson and Janssens, 2006)

$$349 \quad V_{decom} = V_{\max} C_{ENZ} \frac{C_{SOC}}{K_m + C_{SOC}} \quad (14)$$

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

$$354 \quad V_{uptake} = V_{\max_up} C_{MIC} \frac{C_{DOC}}{K_{m_up} + C_{DOC}} \frac{C_{O_2}}{K_{m_upO_2} + C_{O_2}}, \quad (15)$$

355 where V_{\max_up} [s^{-1}] is the maximum DOC uptake rate when the substrates is not limiting, C_{MIC}
 356 [gCm^{-3}] is the microbial biomass pool size, C_{DOC} [gCm^{-3}] is the DOC pool size, C_{O_2} [m^3m^{-3}] is
 357 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
 358 corresponding half-saturation constants for DOC and O_2 , respectively. With the explicit
 359 representation of soil moisture control, the two rates become (Zhang et al., 2014)

$$360 \quad V_{decom} = V_{\max} C_{ENZ} \frac{C_{SOC}}{K_m + C_{SOC}} \left(\frac{\theta}{\theta_s} \right) \quad (16)$$

$$361 \quad V_{uptake} = V_{\max_up} C_{MIC} \frac{C_{DOC}}{K_{m_up} + C_{DOC}} \frac{C_{O_2}}{K_{m_upO_2} + C_{O_2}} \left(\frac{\theta}{\theta_s} \right) \quad (17)$$

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

363 In addition to using the new rate equations, models 5C and 6C have more carbon pools. In
 364 model 5C, DOC is split into two sub-pools for wet zone and dry zone of soil pores, and only the

365 wet DOC is used by MIC, as shown in Figure 1. The moisture-controlled microbial uptake rate
 366 becomes

$$367 \quad V_{uptake} = V_{max_up} C_{MIC} \frac{C_{DOC_w}}{K_{m_up} + C_{DOC_w}} \frac{C_{O2}}{K_{m_upO2} + C_{O2}} \left(\frac{\theta}{\theta_s} \right). \quad (18)$$

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

$$372 \quad V_{decom} = V_{max} C_{ENZ_W} \frac{C_{SOC}}{K_m + C_{SOC}} \left(\frac{\theta}{\theta_s} \right) \quad (19)$$

373 for the wet ENZ and

$$374 \quad V_{decom} = V_{max} C_{ENZ_D} \frac{C_{SOC}}{K_m + C_{SOC}} \left(1 - \frac{\theta}{\theta_s} \right) \varepsilon_D \quad (20)$$

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

378 Due to considering the moisture control and adding more soil pools, model 5C is expected to
 379 be significantly better than model 4C for simulating the Birch effect. Since the accumulated ENZ
 380 in dry soil is secondary, model 6C is expected to be slightly better than model 5C. In terms of
 381 model structural error, model 4C has the largest model structure error, model 5C has significantly
 382 less model structure error, and model 6C has the smallest model structural error. As shown below,
 383 the degree of model structural error is reflected in the process of Bayesian inference and verified
 384 by the cross-validation.

385 2.4 Observations and parameter estimation

386 Figure 2 plots the time series of 17,016 observations of soil moisture and CO₂ efflux used in
387 this study. The observations were obtained during the entire year of 2007, covering a long period
388 of dry season prior to monsoon and episodic rainfall events during monsoon. The first two third of
389 this dataset is used for the Bayesian inference, and the last one third is used for cross-validation.
390 The inference and cross-validation periods have both dry and wet periods, as shown in Figure 2.
391 The observation site is located within the Santa Rita Experimental Range (SRER, 31.8214°N,
392 110.8661°W, elevation 1,116 m) outside of Tucson, Arizona (Barron-Gafford et al., 2011; Scott
393 et al., 2009). This savanna site was covered by 22% of perennial grass, forbs and subshrubs and
394 35% of mesquite. The soils are uniformly Comoro loamy sand (77.6% sand, 11.0% clay, and
395 11.4% silt). The half-hourly atmospheric forcing data were collected from measurements through
396 an eddy covariance tower (Scott et al., 2009). This includes downward shortwave, longwave,
397 precipitation, wind, air temperature, humidity, and pressure. Volumetric CO₂ concentration was
398 measured at half-hourly interval through compact probes. The CO₂ efflux was estimated from the
399 gradient of CO₂ concentration measured at two depths of 2 cm and 10 cm through Fick's first law
400 of diffusion, and the estimates were validated against measurements from a portable CO₂ gas
401 analyzer.

402 The parameters estimated in this study include the parameters of the soil respiration models
403 (4C – 6C) and the parameters of the data models described in Section 2.1. The estimated
404 parameters of models 4C and 5C include the microbial carbon use efficiency (CUE) [g/g], enzyme
405 production rate, k_e [g/m³s], microbial turnover rate, τ_m [1/s], and enzyme turnover rate τ_e [1/s].
406 Uniform distributions are used as the prior in the Bayesian inference, and the ranges of the four
407 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.

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

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

417 42.15 Metrics for evaluating predictive performance

418 Three criteria are used to evaluate the predictive performance of the soil respiration models
 419 and data models, and they are central mean tendency, dispersion, and reliability. Each
 420 criteria~~criteria~~ criterion is measured by a single metric. In addition, a newly defined metric is also used
 421 for simultaneously measuring the three criteria. The central mean tendency is measured in this
 422 study using the Nash-Sutcliffe model efficiency (NSME) coefficient (Nash and Sutcliffe, 1970),

$$423 \quad \underline{NSME = 1 - \frac{\sum_{i=1}^n (d_i - \bar{Y}_i)^2}{\sum_{i=1}^n (d_i - \bar{\mathbf{d}})^2}}, \quad (21)$$

424 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
 425 \bar{Y}_i is the mean of the prediction ensemble, \bar{Y}_i , for d_i . NSME ranges from $-\infty$ to 1, with $NSME =$
 426 1 corresponding to a perfect match between data and mean prediction, i.e., the ensemble is centered
 427 on the data. $NSME = 0$ indicates that the model predictions are as only accurate as the mean of the
 428 data, while an efficiency $NSME < 1$ indicates that the mean of data is a better prediction than the
 429 mean prediction.

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

$$433 \textit{Sharpness} = 1/n \sum_{i=1}^n [\textit{Max}(\mathbf{Y}_i) - \textit{Min}(\mathbf{Y}_i)] \quad (22)$$

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

439 To account for the trade-off between the three metrics,(Elshall et al., 2018b) defined relative
 440 model score (RMS) that simultaneously measure all the three criteria. Scoring rules are commonly
 441 used in hydrology to assess predictive performance (e.g. Weijis et al., 2010; Westerberg et al.,
 442 2011). RMS is used in this study to measure the relative predictive performance of the
 443 combinations of soil respiration models and data models. For combination M_j , RMS is defined as

$$444 \textit{RMS}(M_j) = \frac{\sum_{i=1}^n p(d_i | \mathbf{Y}_{ij}, M_j)}{\sum_{j=1}^m p(d_i | \mathbf{Y}_{ij}, M_j)} \times 100 \quad (23)$$

445 where $m = 24$ is the number of combinations; and the ensemble prediction $\mathbf{X}\mathbf{Y}_{ij}$ is similar to $\mathbf{X}\mathbf{Y}_i$
 446 above where is i index over time and specific to the j -th combination. The density function,
 447 $p(d_i | \mathbf{X}\mathbf{Y}_{ij})$, can be evaluated by first obtaining the density function $p(\mathbf{X}\mathbf{Y}_{ij})$ of the ensemble
 448 prediction $\mathbf{X}\mathbf{Y}_{ij}$ (e.g., by using the kernel density function) and then evaluating $p(d_i | \mathbf{X}\mathbf{Y}_{ij})$ using
 449 interpolation methods based on the intersection of $\mathbf{X}\mathbf{Y}_{ij}$ and d_i . This evaluation is based purely on
 450 the model predictions, and does not involve any assumptions on the models, their parameters,

451 and likelihood functions. Larger RMS values indicate better overall predictive performance. A
452 figure of our workflow scheme is presented in the supplementary materials.

453

454 **3 Results of Bayesian Inverse Modeling**

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

458 **3.1 Residual characterization**

459 Figure 3 shows residual plots for model 6C based on data models SLS and WSEP-AC. SLS is
460 the simplest one with the assumptions of homoscedastic, independent, and Gaussian residuals, and
461 the WSEP-AC is the most complex one without the assumptions. Model 6C is the most complex
462 model and also the best one as ranked by Zhang et al. (2014) using Bayesian model selection. The
463 variable a_t plotted in Figures 3a-3c and Figures 3d-3f is defined in equations (2) and (11),
464 respectively. Figures 3a – 3c show that the three residual assumptions are violated when SLS is
465 used because (i) the residual variance is not constant, but increases as a function of the simulated
466 CO₂ efflux (Figure 3a); (ii) the autocorrelation function at most lags is beyond the 95% confidence
467 interval (Figure 3b); (iii) and the standard normal density function cannot adequately characterize
468 the residuals (Figure 3c). Figures 3d-f show that, after relaxing the three assumptions, the
469 processed residuals, a_t , can be well characterized by WSEP-AC. Figure 3d shows that, after
470 normalizing ε_t with the linear variance ($\sigma_t = 0.034 + 0.099 E_t$), the variation of the variance of
471 a_t becomes significantly smaller, although the variance is still not a-constant. Figure 3e shows that,
472 after removing a first-order autoregressive model from ε_t , a_t becomes less correlated, although the
473 correlation is not fully removed. The two coefficients of the autoregressive model are $\phi_1 = 0.989$

474 and $\phi_2 = 4.5 \times 10^{-6}$; the small value of ϕ_2 indicates that there is no need to attempt an autoregressive
475 model of higher order. Figure 3f shows that a_t follows the SEP distribution with the estimated
476 skewness coefficient of $\xi = 0.933$ and kurtosis coefficient of $\beta = 0.998$. As a summary, Figure
477 3 shows that it is important to examine the residuals and to determine whether a data model is
478 adequate for charactering the residuals. Although WSEP-AC still cannot perfectly characterize ε_t ,
479 it is significantly better than SLS.

480 Although the Gaussian assumption used in SLS is violated for model 4C (Figure 3c), this is
481 not generally the case for other data models and physical-soil respiration models. This is shown in
482 Figure 4, which presents the quantile-quantile (Q-Q) plot for the eight data models and the three
483 soil respiration models. For SLS, WLS, SLS-AC, and WLS-AC, the theoretical quantiles are based
484 on the standard normal distribution, $N(0,1)$; for SEP, WSEP, SEP-AC, and WSEP-AC, the
485 theoretical quantiles are based on the standard skew exponential power distribution, $SEP(0,1,1,0)$.
486 If the residuals follow the assumed standard distributions, the Q-Q plots fall on the 1:1 line, which
487 is marked as the theoretical lines in Figure 4. If the residuals are Gaussian or SEP but not standard,
488 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
489 respiration models, the Q-Q plots of SLS and SEP deviate significantly from the theoretical lines
490 and exhibit fat-tail behaviors, which is an indication of outliers (Thyer et al., 2009). The deviation
491 is reduced after accounting for autocorrelation in SLS-AC and SEP-AC, as shown in Figures 4c
492 and 4g. ~~(It is interesting to observe from the two figures that the Q-Q plots of the three models~~
493 ~~are almost visually identical).~~ The deviation is almost fully removed after accounting for
494 heteroscedasticity in WLS and WSEP in that their corresponding Q-Q plots fall on the 1:1 lines,
495 especially for models 5C and 6C, as shown in Figures 4b and 4f. However, the Q-Q plots start
496 deviating from the 1:1 lines as shown in Figures 4d and 4h, after accounting for both

497 heteroscedasticity and autocorrelation in WLS-AC and WSEP-AC. As a summary, Figure 4 shows
498 that, for the numerical example of this study, either the Gaussian or the SEP distribution is valid if
499 heteroscedasticity is accounted for in the data models. However, accounting for autocorrelation in
500 the data models does not help improve the characterization of the residual distribution.

501 **3.2 Posterior parameter distributions**

502 While Figures 3 and 4 help understand validity of the three assumptions used in the data
503 models, the impacts of the data models on estimating model parameter distributions must be
504 evaluated separately. This section discusses the impact of the data model selection on parameter
505 estimation with the objective of understanding if incorrect specification of the data model, will
506 necessarily lead to biased parameter estimates. Such assessment is not a trivial task for three main
507 reasons. First, microbial soil respiration models aggregate complex natural processes and spatial
508 details into simpler conceptual representations. As a results several model parameters are effective
509 values of several complex natural processes that cannot be actually measured in the field as
510 discussed by Vrugt et al. (2013). Second, even for model parameter that can be measured in the
511 field, since the model structure is imperfect, it can be the case that parameter values can be
512 accepted beyond their physically reasonable range as discussed by Pappenberg and Beven
513 (2006). This is often undesirable, if we seek to make the models more mechanistically descriptive.

514 We focus our discussion on carbon use efficiency (CUE) for microbial growth since CUE is a
515 fundamental parameter in microbial soil respiration models, and a reasonable physical range for
516 CUE can estimated. The concept of microbial CUE(Allison et al., 2010; Bradford et al., 2008;
517 Manzoni et al., 2012; Wieder et al., 2013) has been used to present fundamental microbial
518 processes recent microbial enzyme models(Allison et al., 2010; German et al., 2011; Schimel and
519 Weintraub, 2003; Wang et al., 2013). The microbial CUE, which is marked between MIC and CO₂

520 in Figure 1, controls microbial growth, enzyme production and microbial respiration. A reasonable
521 range of CUE can be estimated from the physical viewpoint(Tang and Riley, 2014). Sinsabaugh
522 et al. (2013) study shows that the thermodynamic calculations support a maximum CUE of 0.60
523 and that methods used to estimate CUE in terrestrial systems report a mean value of 0.55.
524 Theoretically, there no lower limit for CUE as it can approach zero, and $CUE < 0.1$ are reported
525 for terrestrial ecosystems (e.g. Fernández-Martínez et al., 2014) and used in modeling studies (Li
526 et al., 2014).

527 Figure 5 plots the CUE posterior marginal density of the three soil respiration models obtained
528 using the eight data models. The physical range between zero and 0.6 is marked in yellow. Figure
529 5 shows that the CUE posterior parameter distribution for Model 6C for all likelihood functions
530 that does not account for autocorrelation are within a reasonable physical range. For models 4C
531 and 5C, the posterior parameter samples are outside the physical range for six data models. For
532 model 4C, the posterior parameters are within the physical range only for data models SEP and
533 WSEP; for model 5C, the two data models are WLS and WSEP. It is not surprising to find the
534 posterior parameter distribution of models 4C and 5C, which have a certain degree of model
535 structure error, to be out of the plausible physical range. This can be attributed to two reasons.
536 First, the model solution can be biased toward the missing processes in the model structure such
537 as the additional carbon pool in both 4C and 5C or the explicit accounting for soil moisture in 4C.
538 Second, biased parameter estimation can compensate for model structure inadequacy and other
539 sources of discrepancy in both the physical model and the statistical model.

540 In addition, it is important to understand how accounting for autocorrelation, heteroscedasticity
541 and non-Gaussian residuals can affect the parameter estimation. First, ~~it is not unexpected to get we~~
542 obtained biased parameter estimates that ~~can be is~~ out the reasonable physical range when

543 autocorrelation is explicitly accounted for as shown in Figure 5e-h. This may suggest again that
544 accounting for heteroscedasticity is desirable but accounting for autocorrelation is not. A possible
545 reason is that filtering autocorrelation may reduce the residual space such that the transformed
546 residual space cannot correspond to the parameter space of the models. In other words, parameter
547 information may be lost due to filtering out autocorrelation. However, it is not fully understood
548 why this does not occur for the model 6C under data model SLS-AC, and more research is
549 warranted. Second, unlike accounting for auto-correlation, accounting only for heteroscedasticity
550 (i.e. WLS and WSEP) since this will only amplify or reduce the variance without affecting the
551 structure of the residual space. Figure 5c-d shows that account for heteroscedasticity (i.e. WLS
552 and WSEP) tends to improve the parameter estimation in comparison with homoscedastic data
553 models (i.e. SLS and SEP) shown in Figure 5a-b. Finally, with respect to non-Gaussian residuals,
554 Schoups and Vrugt (2010) proposes that the peaked pdf of the SEP with heavier tails compared to
555 Gaussian pdf is useful for making parameter inference robust against outliers. To a certain degree,
556 this can be substantiated by the results in Figure 5a-d, such that SEP and WSEP provide more
557 favorable parameter estimates than SLS and WLS.

558 Finally, from Figure 5 we can also notice that the posterior parameter distribution of SLS
559 (Figure 5a) is very narrow. This narrow posterior parameter distribution of SLS compared to other
560 likelihood functions can be attributed to several reasons. Since SEP can have heavier tails than
561 Gaussian distribution, this can further increase the samples acceptance ratio from tails resulting in
562 wider distribution (Figure 5b). In addition, accounting for heteroscedasticity will wider the
563 posterior parameter distribution (Figure 5c) due to accepting higher variances at peak effluxes.
564 Moreover, filtering correlation (Figure 5e-h) increases the entropy.

565 4. Results of Predictive Performance

566 Based on the last one third of the CO₂ efflux observations, a cross-validation test was
567 conducted for all the 24 models, -the combinations of three soil respiration models and eight data
568 models. Given the cross-validation ~~dataperiod~~, the predictive performance is examined using the
569 four statistical metrics ~~that are~~ defined in Section ~~24.51~~. The metrics are also calculated for the
570 calibration ~~dataperiod~~. This is not to perform Bayesian model selection given the calibration data,
571 but to better understand the impact of data models. For each calibration and each cross-validation
572 data, a prediction ensemble is generated from the two perspectives of parametric uncertainty only
573 and total uncertainty, as presented in Section ~~4.21~~ and ~~4.23~~, respectively.

574 ~~4.1 Metrics for evaluating predictive performance~~

575 ~~Three criteria are used to evaluate the predictive performance of the soil respiration models~~
576 ~~and data models, and they are central mean tendency, dispersion, and reliability. Each criteria is~~
577 ~~measured by a single metric. In addition, a newly defined metric is also used for simultaneously~~
578 ~~measuring the three criteria. The central mean tendency is measured in this study using the Nash-~~
579 ~~Sutcliffe model efficiency (NSME) coefficient (Nash and Sutcliffe, 1970),~~

$$580 \quad NSME = 1 - \frac{\sum_{i=1}^n (d_i - \bar{X}_i)^2}{\sum_{i=1}^n (d_i - \bar{\mathbf{d}})^2}, \quad (21)$$

581 ~~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~~
582 ~~\bar{X}_i is the mean of the prediction ensemble, X_i , for d_i . NSME ranges from $-\infty$ to 1, with $NSME = 1$~~
583 ~~corresponding to a perfect match between data and mean prediction, i.e., the ensemble is centered~~
584 ~~on the data. $NSME = 0$ indicates that the model predictions are as only accurate as the mean of the~~
585 ~~data, while an efficiency $NSME < 1$ indicates that the mean of data is a better prediction than the~~
586 ~~mean prediction.~~

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

$$590 \text{Sharpness} = 1/n \sum_{i=1}^n [\text{Max}(X_i) - \text{Min}(X_i)] \quad (22)$$

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

596 To account for the trade-off between the three metrics, (Elshall et al., 2018) defined relative
 597 model score (RMS) that simultaneously measure all the three criteria. Scoring rules are commonly
 598 used in hydrology to assess predictive performance (e.g. Weijis et al., 2010; Westerberg et al.,
 599 2011). RMS is used in this study to measure the relative predictive performance of the
 600 combinations of soil respiration models and data models. For combination M_j , RMS is defined as

$$601 \text{RMS}(M_j) = \frac{\sum_{i=1}^n p(d_i | X_{ij}, M_j)}{\sum_{j=1}^m p(d_i | X_{ij}, M_j)} \times 100 \quad (23)$$

602 where $m=24$ is the number of combinations, and X_{ij} is similar to X_i above and specific to the j -th
 603 combination. The density function, $p(d_i | X_{ij})$, can be evaluated by first obtaining the density function
 604 $p(X_{ij})$ of the ensemble prediction X_{ij} (e.g., by using the kernel density function) and then evaluating
 605 $p(d_i | X_{ij})$ using interpolation methods based on the intersection of X_{ij} and d_i . This evaluation is based
 606 purely on the model predictions, and does not involve any assumptions on the models, their
 607 parameters, and likelihood functions. Larger RMS values indicate better overall predictive
 608 performance.

609 **4.21 Predictive performance with parametric uncertainty of soil respiration models**

610 In this section the ensemble is generated by running the soil respiration models with the
611 posterior samples (obtained from the Bayesian inference) of the physical model parameters. In
612 other words, the ensemble addresses parametric uncertainty of the soil respiration models only.
613 Considering the relative contribution of parametric uncertainty only will provide insights for
614 modeling approaches that attempt to segregate various sources of uncertainty (e.g. Thyer et al.,
615 2009; Elshall and Tsai, 2014); (Tsai and Elshall, 2013).

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

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

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

632 6h for data models SLS and WSEP-AC. For models 5C and 6C, their NSME values vary with the
633 data models; ~~with and~~ the central mean accuracy ~~is being~~ the worst for SLS-AC ~~that~~which
634 considers only autocorrelation.

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

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

646 As a summary, while it is necessary to account for heteroscedasticity in a data model, caution
647 is needed when accounting for autocorrelation in the manner described in Section 2.1. In addition,
648 after comparing the RMS values of the residuals using the Gaussian and SEP distributions, ~~t.~~The
649 conclusion is that the SEP distribution outperforms the Gaussian distribution with respect to
650 predictive performance. Finally, uncertainty underestimation as evident by the very small
651 predictive coverage. The underestimation of uncertainty for all the physical models with all
652 likelihood functions ~~s~~ makes ~~s~~ sense because only parametric uncertainty is considered. Considering
653 the overall predictive uncertainty is the subject of the next section.

4.3.2 Predictive performance with ~~parametric uncertainty of soil respiration models and uncertainty from data model~~total uncertainty s

The simulated output $Y(\theta_p)$ will generally not be equal to the observed output \mathbf{d} and we have a residual ~~ue error~~ term \mathbf{e} due to measurement, input and model structure errors such that $\mathbf{d} = Y(\theta_p) + \mathbf{e}$. Accounting for the error term \mathbf{e} can be through separating various error terms. For example, in section 4.2-1 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 measurement, model input, and model structure uncertainty ~~which include measurement, model input, model structure and parameter estimation errors~~ (e.g. Thyer et al., 2009; Schoups and Vrugt, 2010). This lumped approach is based on sampling the ~~residual error model~~ residuals model $\mathbf{e}(\theta_e)$ with parameters θ_e . SLS has one fixed parameter that is the constant variance and other data models have two to six parameters. Thus in ~~Section 4.3~~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 4.12.

We start by 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

677 of sampling the residuals ~~s_errors~~-model. Figure 8b shows that despite the wide prediction interval
678 of model 4C, which has significant model structure error, it could not capture the birch pulse
679 around day 180. This clearly indicates that proper modeling of the residuals ~~s_error~~ will not make-
680 up for of significant model structure error.

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

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

694 With respect to the overall predictive performance, the RMS values are largely determined by
695 mean accuracy and sharpness as the predictive coverage is similar for different data models.
696 Figures 9g and 9h of RMS show that the predictive performance of the four data models that
697 account for autocorrelation is worse than that of the other four data models. This suggests again
698 that one needs to be cautious when building autocorrelation into a data model. This is consistent
699 with the finding of Evin et al. (2013, 2014) that accounting for autocorrelation before accounting

700 for heteroscedasticity or jointly accounting for autocorrelation and heteroscedasticity can result in
701 poor predictive performance. In summary, Figures 9g and 9h show for both the calibration and
702 prediction periods that accounting for heteroscedasticity (i.e. WLS and WSEP) will give the best
703 overall predictive skillperformance, and accounting for autocorrelation without heteroscedasticity
704 (i.e. SLS-AC and SEP-AC) will give the worst overall predictive skillperformance. Finally, for the
705 three soil respiration models, RMS shows that model 4C has the worst predictive performance for
706 both the calibration and cross-validation data. Generally speaking, the high fidelity model 6C
707 outperforms model 5C for both the calibration and cross-validation data, which justifies the
708 complexity of model 6C.

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

712 In Figure 10 we try to understand the predictive performance characteristics of the different
713 data models by looking at the predictive performance of model 6C. Specific predictive
714 performance patterns can be identified. Figures 10-a-d show that SLS and SEP have similar
715 predictive performance with SEP generally having better predictive skillperformance especially
716 during the validation period. Accounting for heteroscedasticity using WLS as shown in Figures
717 10e and 10h will make the predictions more sensitive to peak carbon effluxes and will generally
718 improve the predictive coverage on the expense of sharpness and the central mean tendency. WLS
719 and WSEP have similar predictive performance. However, WSEP maintains slightly better central
720 mean tendency and overall predictive performance than WLS. Accounting for autocorrelation
721 using SLS-AC and SEP-AC as shown in Figures 10i and 10l reduces the information content of
722 the residuals, and thus resulting in wider uncertainty bands and insensitivity to peak carbon

723 effluxes as compared to SLS and SEP (Figures 10a-d). This resulted in deteriorating the sharpness,
724 the central mean tendency and the capturing of peak carbon fluxes, especially during the validation
725 period. Accounting for both heteroscedasticity and autocorrelation using WLS-AC and WSEP-
726 AC will make the inference robust against peak carbon effluxes, yet due to the loss of information
727 content uncertainty bands are still wider and uncertainty becomes overestimated especially during
728 validation period as compared to WLS and WSEP. The results of Models 4C and 5C, which are
729 not shown here, also show the same prediction patterns with respect to non-Gaussian residuals,
730 heteroscedasticity and autocorrelation.

731 From figure 10 we also notice that data models that have good overall predictive performance
732 as measured by RMS during the calibration period will maintain this good predictive performance
733 during the validation period. For model 6C, RMS values for the calibration and validation periods
734 are very well correlated with a correlation coefficient of 0.92. However, we note that for models
735 4C and 5C the overall predictive performance during the calibration and validation periods are not
736 that well correlated as 6C, with correlation coefficients of 0.52 for model 4C and 0.61 for model
737 5C. This suggests that model 6C is more robust than 4C and 5C for forecasting and hindcasting.

738 5. Conclusions

739 In parameter estimation and prediction of soil carbon fluxes to the atmosphere we often
740 assume that residuals, which include observation, model input, model ~~structure and parameter~~
741 ~~estimation~~ errors, are normally distributed, homoscedastic and uncorrelated. We studied these
742 assumptions by calibrating three microbial enzyme models, which have varying degrees of model
743 structure errors. We tested eight data model~~ing~~ starting with the standard least squares (SLS) and
744 skew exponential power (SEP) data models that assume homoscedast~~ictic~~ and non-correlated
745 residuals. Given these two distributions, we evaluated six other data models that account for

746 heteroscedasticity (WLS and WSEP), autocorrelation (SLS-AC and SEP-AC) and joint inversion
747 of heteroscedasticity and autocorrelation (WLS-AC and WSEP-AC). To our knowledge this is the
748 first study that provide such detailed analysis soil respiration inverse modeling. We also used three
749 solid respiration models with different degrees of model fidelity (i.e. model realism) and model
750 complexity (i.e. number of model parameters), to understand the impact of model discrepancy on
751 the calibration results under different data models. We analyzed the calibration results with respect
752 to (i) residual characterization, (ii) parameter estimation, (iii) predictive performance and (iv)
753 impact of model discrepancy. The main findings of this study can be ~~synthesized~~ summarized as
754 follows:

755 (i) With respect to residual characterization, residual analysis results suggest that the common
756 assumption of not accounting for heteroscedasticity and autocorrelation of residuals (i.e. SLS and
757 SEP) results in poor characterization of residuals. Explicit accounting for heteroscedasticity (i.e.
758 WLS and WSEP) can result in good characterization of the residuals, and is followed by joint the
759 inversion of heteroscedasticity and autocorrelation (i.e. WSL-AC and WSEP-AC). Accounting for
760 autocorrelation only (i.e. SLS-AC and SEP-AC) may not improve much the characterization of the
761 residuals.

762 (ii) With respect to parameter estimation, we focused on carbon use efficiency (CUE), which
763 is a central parameter in soil respiration modeling. We found the SLS with relatively reasonable
764 posterior parameter distribution for CUE, yet very narrow posterior. Data models consider
765 autocorrelation (i.e. SLS-AC, SEP-AC, WLS-AC and WSEP-AC) tend to generally yield CUE
766 estimates that are physically non-reasonable. We speculate that filtering correlation can affect the
767 mapping of the model physics (as implicitly included in the residuals) into the likelihood space,
768 which might result in biased parameter estimates that are physically unreasonable.

769 (iii) With respect to predictive performance, we assessed the central mean tendency,
770 uncertainty bands and the overall predictive performance for both the calibration and the cross-
771 validation periods. Results show that accounting for autocorrelation (i.e. SLS-AC, SEP-AC, WLS-
772 AC, and WSEP-AC) deteriorate^s the predicative performance, such that the predictive
773 performance is inferior to SLS in terms of the central mean tendency and overall predictive
774 skillperformance, especially during the cross-validation period. Results also indicates that using a
775 SEP distribution can potentially improve the predictive performance. The same is true for
776 accounting for heteroscedasticity. Using SEP distribution and accounting for heteroscedasticity
777 (i.e. WSEP) can potentially improve the predictive performance.

778 (iv) With respect to the impact of model discrepancy, the high fidelity complex model (6C)
779 gives the best results with respect to parameter estimation and predictive performance. Model 6C
780 generally maintained its superior performance under different data models. This justifies the
781 complexity of model 6C relative to model 5C that has one less carbon pool. Model 4C that has a
782 low fidelity model with only four carbon pools and lacks the explicit representation of soil moisture
783 control, maintains its poor performance for different data models.

784 From the empirical findings of this research we conclude the following:

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

789 (ii) Using a non-Gaussian residual error model can improve the parameter estimates, and the
790 predictive performance with respect to central mean tendency and uncertainty estimation.

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

793 (iv) This study confirms the empirical findings and theoretical analysis ~~of~~ (Evin et al., (2013;
794 2014; [Ammann et al. 2018](#))) that separate accounting for autocorrelation or joint inversion of
795 correlation and heteroscedasticity can be problematic. ~~Relatively poor performance with respect~~
796 ~~to autocorrelation can be due to our implementation scheme.~~ By drawing on similarity from surface
797 hydrology, the study of Ammann et al. (2018) suggests that this might be attributed to non-
798 stationarity due to wet-dry periods with half-hourly data. Accounting for non-stationarity (Smith
799 et al., 2010b, Ammann et al. 2018) could address this problem. Relatively poor performance with
800 respect to autocorrelation can be also attributed to the implementation scheme. The inference
801 scheme such as joint inference as in this study, which can be improved by using the post-processing
802 inference approach for autocorrelation (Evin et al., 2013; 2014), residuals transformation approach
803 (e.g. Lu et al., 2013) or similar other strategies (Li et al., 2015, 2016a) could have an impact. Yet
804 (Ammann et al., (2018) study states that the joint inversion is still preferred, and understanding the
805 conditions where accounting for auto-correlation can be achieved remain poorly understood.
806 Further investigation of this point is warranted in a future study.

807 The above conclusions ~~above~~ are subject to several limitations. First, the conclusions are
808 specific to the soil respiration models developed and validated for semi-arid savannah.
809 Performance variations across different soil respiration models with different levels of
810 complexities is possible. Second, the conclusions are conditioned on the data that were obtained
811 at the half-hour interval over a one-year period. Different conclusions are possible if the data are
812 thinned to daily or weekly scales or data of longer observation periods are used. Third, the study
813 investigates effects of the residual assumptions of formal likelihood functions through direct

814 conditioning of the error-residuals model parameters, yet this can also be done through other
815 approaches such as residuals transformation (Thiemann et al., 2001), autogressive bias model
816 (Del Giudice et al., 2013), approximate Bayesian computation (Sadegh and Vrugt, 2013), data
817 assimilation (Spaaks and Bouten, 2013). Comparing different methods for accounting the residual
818 assumptions are beyond the scope of this work. Fourth, this study focuses on formal Bayesian
819 computation using formal likelihood functions, and comparison with other inference functions
820 such as informal likelihood functions or approximate Bayesian computation is warranted in a
821 future study.

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

830 Acronyms

831	<u>4C</u>	<u>Four carbon pool model</u>
832	<u>5C</u>	<u>Five carbon pool model</u>
833	<u>6C</u>	<u>Six carbon pool model</u>
834	<u>CUE</u>	<u>Microbial carbon use efficiency</u>
835	<u>DOC</u>	<u>Dissolved organic carbon</u>
836	<u>ENZ</u>	<u>Enzymes</u>
837	<u>MCMC</u>	<u>Markov chain Monte Carlo</u>
838	<u>MIC</u>	<u>Microbial biomass</u>
839	<u>NSME</u>	<u>Nash-Sutcliffe model efficiency</u>
840	<u>PDF</u>	<u>Probability density function</u>
841	<u>RMS</u>	<u>Relative model score</u>
842	<u>SEP</u>	<u>Skew exponential power distribution</u>

843	<u>SEP-AC</u>	<u>Skew exponential power distribution with autocorrelation</u>
844	<u>SLS</u>	<u>Standard least square</u>
845	<u>SLS-AC</u>	<u>Standard least square with autocorrelation</u>
846	<u>SOC</u>	<u>Soil organic carbon</u>
847	<u>WLS</u>	<u>Weighted least squared</u>
848	<u>WLS-AC</u>	<u>Weight least square with autocorrelation</u>
849	<u>WSEP</u>	<u>Weighted skew exponential power distribution</u>
850	<u>WSEP-AC</u>	<u>Weighted skew exponential power distribution with autocorrelation</u>

851

852 **Code and data availability**

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

857 **Author contributions**

858 ASE developed and implemented the code for the eight data models for soil respiration modeling,
859 and prepared the manuscript with contribution of all co-authors. MY developed the research idea
860 and outline, and supervised the research implementation. GN developed the soil respiration
861 models. GAB collected and processed the eddy-covariance data used for model calibration.

862 **Competing interests**

863 The authors declare that they have no conflict of interest.

864 **Acknowledgement**

865 This work was supported by the Department of Energy Early Career Award, DE-SC0008272 and
866 U.S. National Science Foundation Award# OIA-1557349.

867 **References**

868 Ahrens, B., Reichstein, M., Borken, W., Muhr, J., Trumbore, S. E. and Wutzler, T.: Bayesian
869 calibration of a soil organic carbon model using ΔC measurements of soil organic carbon
870 and heterotrophic respiration as joint constraints, *Biogeosciences*, 11(8), 2147–2168,

871 doi:10.5194/bg-11-2147-2014, 2014.

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

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

877 Bagnara, M., Sottocornola, M., Cescatti, A., Minerbi, S., Montagnani, L., Gianelle, D. and
878 Magnani, F.: Bayesian optimization of a light use efficiency model for the estimation of
879 daily gross primary productivity in a range of Italian forest ecosystems, *Ecol. Modell.*, 306,
880 57–66, doi:10.1016/j.ecolmodel.2014.09.021, 2015.

881 Bagnara, M., Oijen, M. Van, Cameron, D., Gianelle, D., Magnani, F. and Sottocornola, M.:
882 Bayesian calibration of simple forest models with multiplicative mathematical structure : A
883 case study with two Light Use Efficiency models in an alpine forest, *Ecol. Modell.*,
884 371(January), 90–100, doi:10.1016/j.ecolmodel.2018.01.014, 2018.

885 Barr, J. G., Engel, V., Fuentes, J. D., Fuller, D. O. and Kwon, H.: Modeling light use efficiency in
886 a subtropical mangrove forest equipped with CO₂ eddy covariance, *Biogeosciences*, 10(3),
887 2145–2158, doi:10.5194/bg-10-2145-2013, 2013.

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

891 Barron-Gafford, G. A., Scott, R. L., Jenerette, G. D. and Huxman, T. E.: The relative controls of
892 temperature, soil moisture, and plant functional group on soil CO₂ efflux at diel, seasonal,
893 and annual scales, *J. Geophys. Res. Biogeosciences*, 116(1), 1–16,

894 doi:10.1029/2010JG001442, 2011.

895 Berryman, E. M., Frank, J. M., Massman, W. J. and Ryan, M. G.: Agricultural and Forest
896 Meteorology Using a Bayesian framework to account for advection in seven years of
897 snowpack CO₂ fluxes in a mortality-impacted subalpine forest, *Agric. For. Meteorol.*,
898 249(April 2017), 420–433, doi:10.1016/j.agrformet.2017.11.004, 2018.

899 Box, E. P. and Tiao, G. C. (1992), *Bayesian Inference in Statistical Analysis*, 588 pp., Wiley, New
900 York.

901 Braakhekke, M. C., Beer, C., Schrumpf, M., Ekici, A., Ahrens, B., Hoosbeek, M. R., Kruijt, B.,
902 Kabat, P. and Reichstein, M.: The use of radiocarbon to constrain current and future soil
903 organic matter turnover and transport in a temperate forest, *J. Geophys. Res. Biogeosciences*,
904 372–391, doi:10.1002/2013JG002420.Received, 2014.

905 Bradford, M. A., Davies, C. A., Frey, S. D., Maddox, T. R., Melillo, J. M., Mohan, J. E., Reynolds,
906 J. F., Treseder, K. K. and Wallenstein, M. D.: Thermal adaptation of soil microbial
907 respiration to elevated temperature, *Ecol. Lett.*, 11(12), 1316–1327, doi:10.1111/j.1461-
908 0248.2008.01251.x, 2008.

909 Braswell, B. H., Sacks, W. J., Linder, E. and Schimel, D. S.: Estimating diurnal to annual
910 ecosystem parameters by synthesis of a carbon flux model with eddy covariance net
911 ecosystem exchange observations, *Glob. Chang. Biol.*, 335–355, doi:10.1111/j.1365-
912 2486.2005.00897.x, 2015.

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

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

917 Ewers, B. E., Pendall, E., Welker, J. M., Charlet, T. N., Cleary, M., Griffith, A., Nowak, R.
918 S., Rogers, M., Steltzer, H., Sullivan, P. F. and Gestel, N. C. Van: The temperature responses
919 of soil respiration in deserts: a seven desert synthesis, *Biogeochemistry*, 71–90,
920 doi:10.1007/s10533-010-9448-z, 2011.

921 Chevallier, F. and O’Dell, C. W.: Error statistics of Bayesian CO₂ flux inversion schemes as seen
922 from GOSAT, *Geophys. Res. Lett.*, 40(6), 1252–1256, doi:10.1002/grl.50228, 2013.

923 Correia, A. C., Minunno, F., Caldeira, M. C., Banza, J., Mateus, J., Carneiro, M., Wingate, L.,
924 Shvaleva, A., Ramos, A., Jongen, M., Bugalho, M. N., Nogueira, C., Lecomte, X. and
925 Pereira, J. S.: Agriculture, Ecosystems and Environment Soil water availability strongly
926 modulates soil CO₂ efflux in different Mediterranean ecosystems: Model calibration using
927 the Bayesian approach, *Agric. Ecosyst. Environ.*, 161, 88–100,
928 doi:10.1016/j.agee.2012.07.025, 2012.

929 Davidson, E. A. and Janssens, I. A.: Temperature sensitivity of soil carbon decomposition and
930 feedbacks to climate change, *Nature*, 440, 165 [online] Available from:
931 <http://dx.doi.org/10.1038/nature04514>, 2006.

932 Davidson, E. A., Samanta, S., Caramori, S. S. and Savage, K.: The Dual Arrhenius and Michaelis–
933 Menten kinetics model for decomposition of soil organic matter at hourly to seasonal time
934 scales, *Glob. Chang. Biol.*, 18(1), 371–384, doi:10.1111/j.1365-2486.2011.02546.x, 2011.

935 Du, Z., Nie, Y., He, Y., Yu, G. and Wang, H.: Tellus B: Chemical and Physical Meteorology
936 Complementarity of flux- and biometric-based data to constrain parameters in a terrestrial
937 carbon model Complementarity of flux- and biometric-based data to constrain parameters in
938 a terrestrial carbon model, *Tellus B Chem. Phys. Meteorol.*, 0889,
939 doi:10.3402/tellusb.v67.24102, 2015.

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

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

946 Elshall, A. S., Ye, M., Pei, Y., Zhang, F., Niu, G.-Y. and Barron-Gafford, G. A.: Relative model
947 score: a scoring rule for evaluating ensemble simulations with application to microbial soil
948 respiration modeling, *Stoch. Environ. Res. Risk Assess.*, doi:10.1007/s00477-018-1592-3,
949 2018a.

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

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

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

961 Fernández-Martínez, M., Vicca, S., Janssens, I. A., Sardans, J., Luysaert, S., Campioli, M.,
962 Chapin III, F. S., Ciais, P., Malhi, Y., Obersteiner, M., Papale, D., Piao, S. L., Reichstein,

963 M., Rodà, F. and Peñuelas, J.: Nutrient availability as the key regulator of global forest
964 carbon balance, *Nat. Clim. Chang.*, 4, 471 [online] Available from:
965 <http://dx.doi.org/10.1038/nclimate2177>, 2014.

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

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

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

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

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

980 Hashimoto, S., Morishita, T., Sakata, T., Ishizuka, S., Kaneko, S. and Takahashi, M.: Simple
981 models for soil CO₂, CH₄, and N₂O fluxes calibrated using a Bayesian approach and multi-
982 site data, *Ecol. Modell.*, 222(7), 1283–1292, doi:10.1016/j.ecolmodel.2011.01.013, 2011.

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

986 Hilton, T. W., Davis, K. J. and Keller, K.: Evaluating terrestrial CO₂ flux diagnoses and
987 uncertainties from a simple land surface model and its residuals, *Biogeosciences*, 11(2), 217–
988 235, doi:10.5194/bg-11-217-2014, 2014.

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

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

994 Hublart, P., Ruelland, D., De Cortázar-Atauri, I. G., Gascoin, S., Lhermitte, S. and Ibacache, A.:
995 Reliability of lumped hydrological modeling in a semi-arid mountainous catchment facing
996 water-use changes, *Hydrol. Earth Syst. Sci.*, 20(9), 3691–3717, doi:10.5194/hess-20-3691-
997 2016, 2016.

998 Ishikura, K., Yamada, H., Toma, Y., Takakai, F., Darung, U., Limin, A. and Limin, S. H.: Soil
999 Science and Plant Nutrition Effect of groundwater level fluctuation on soil respiration rate
1000 of tropical peatland in Central Kalimantan , Indonesia, *Soil Sci. Plant Nutr.*, 63(1), 1–13,
1001 doi:10.1080/00380768.2016.1244652, 2017.

1002 Janssens, I. A., Freibauer, A., Ciais, P., Smith, P., Nabuurs, G.-J., Folberth, G., Schlamadinger, B.,
1003 Hutjes, R. W. A., Ceulemans, R., Schulze, E.-D., Valentini, R. and Dolman, A. J.: Europe’s
1004 terrestrial biosphere absorbs 7 to 12% of European anthropogenic CO₂ emissions., *Science*,
1005 300(5625), 1538–42, doi:10.1126/science.1083592, 2003.

1006 Katz, R. W., Craigmile, P. F., Guttorp, P., Haran, M., Sansó, B. and Stein, M. L.: Uncertainty
1007 analysis in climate change assessments, *Nat. Clim. Chang.*, 3, 769 [online] Available from:
1008 <http://dx.doi.org/10.1038/nclimate1980>, 2013.

1009 Kavetski, D., Franks, S. W. and Kuczera, G.: Confronting Input Uncertainty in Environmental
1010 Modelling, Calibration Watershed Model., doi:doi:10.1029/WS006p0049, 2013.

1011 Keenan, T. F., Davidson, E., Moffat, A. M., Munger, W. and Richardson, A. D.: Using model-data
1012 fusion to interpret past trends, and quantify uncertainties in future projections, of terrestrial
1013 ecosystem carbon cycling, *Glob. Chang. Biol.*, 18(8), 2555–2569, doi:10.1111/j.1365-
1014 2486.2012.02684.x, 2012.

1015 Kim, Y., Nishina, K., Chae, N., Park, S. J., Yoon, Y. J. and Lee, B. Y.: Constraint of soil moisture
1016 on CO₂ efflux from tundra lichen, moss, and tussock in Council, Alaska, using a
1017 hierarchical Bayesian model, *Biogeosciences*, 5567–5579, doi:10.5194/bg-11-5567-2014,
1018 2014.

1019 Klemedtsson, L., Jansson, P. E., Gustafsson, D., Karlberg, L., Weslien, P., Von Arnold, K.,
1020 Ernfors, M., Langvall, O. and Lindroth, A.: Bayesian calibration method used to elucidate
1021 carbon turnover in forest on drained organic soil, *Biogeochemistry*, 89(1), 61–79,
1022 doi:10.1007/s10533-007-9169-0, 2008.

1023 Laloy, E. and Vrugt, J. A.: High-dimensional posterior exploration of hydrologic models using
1024 multiple-try DREAM(ZS) and high-performance computing, *Water Resour. Res.*, 48(1),
1025 doi:10.1029/2011WR010608, 2012.

1026 Li, J., Wang, G., Allison, S. D., Mayes, M. A. and Luo, Y.: Soil carbon sensitivity to temperature
1027 and carbon use efficiency compared across microbial-ecosystem models of varying
1028 complexity, *Biogeochemistry*, 119, 67–84 [online] Available from:
1029 <http://www.jstor.org/stable/24716883>, 2014.

1030 Li, M., Wang, Q. J., Bennett, J. C. and Robertson, D. E.: A strategy to overcome adverse effects
1031 of autoregressive updating of streamflow forecasts, *Hydrol. Earth Syst. Sci.*, 19(1), 1–15,

1032 doi:10.5194/hess-19-1-2015, 2015.

1033 Li, M., Wang, Q. J., Bennett, J. C. and Robertson, D. E.: Error reduction and representation in
1034 stages (ERRIS) in hydrological modelling for ensemble streamflow forecasting, *Hydrol.*
1035 *Earth Syst. Sci.*, 20(9), 3561–3579, doi:10.5194/hess-20-3561-2016, 2016a.

1036 Li, Q., Xia, J., Shi, Z., Huang, K., Du, Z. and Lin, G.: Variation of parameters in a Flux-Based
1037 Ecosystem Model across 12 sites of terrestrial ecosystems in the conterminous USA, *Ecol.*
1038 *Modell.*, 336, 57–69, doi:10.1016/j.ecolmodel.2016.05.016, 2016b.

1039 Lu, D., Ye, M., Meyer, P. D., Curtis, G. P., Shi, X., Niu, X.-F. and Yabusaki, S. B.: Effects of error
1040 covariance structure on estimation of model averaging weights and predictive performance,
1041 *Water Resour. Res.*, 49(9), 6029–6047, doi:10.1002/wrcr.20441, 2013.

1042 Luo, Y., Ogle, K., Tucker, C., Fei, S., Gao, C., LaDeau, S., Clark, J. S. and Schimel, D. S.:
1043 Ecological forecasting and data assimilation in a data-rich era, *Ecol. Appl.*, 21(5), 1429–
1044 1442, doi:10.1890/09-1275.1, 2011.

1045 Luo, Y., Keenan, T. F. and Smith, M.: Predictability of the terrestrial carbon cycle, *Glob. Chang.*
1046 *Biol.*, 21(5), 1737–1751, doi:10.1111/gcb.12766, 2014.

1047 Manzoni, S., Taylor, P., Richter, A., Porporato, A. and Ågren, G. I.: Environmental and
1048 stoichiometric controls on microbial carbon-use efficiency in soils, *New Phytol.*, 196(1), 79–
1049 91, doi:10.1111/j.1469-8137.2012.04225.x, 2012.

1050 McInerney, D., Thyer, M., Kavetski, D., Lerat, J. and Kuczera, G.: Improving probabilistic
1051 prediction of daily streamflow by identifying Pareto optimal approaches for modeling
1052 heteroscedastic residual errors, *Water Resour. Res.*, 53, 2199–2239,
1053 doi:10.1002/2016WR019168.Received, 2017.

1054 Menichetti, L., Kätterer, T. and Leifeld, J.: Parametrization consequences of constraining soil

1055 organic matter models by total carbon and radiocarbon using long-term field data,
1056 Biogeosciences, 3003–3019, doi:10.5194/bg-13-3003-2016, 2016.

1057 Nash, J. E. and Sutcliffe, J. V.: River flow forecasting through conceptual models part I — A
1058 discussion of principles, *J. Hydrol.*, 10(3), 282–290, doi:https://doi.org/10.1016/0022-
1059 1694(70)90255-6, 1970.

1060 Ogle, K., Ryan, E., Dijkstra, F. A. and Pendall, E.: *Journal of Geophysical Research:*
1061 *Biogeosciences*, *J. Geophys. Res. Biogeosciences*, 1–14, doi:10.1002/2016JG003385, 2016.

1062 Pappenberger, F. and Beven, K. J.: Ignorance is bliss: Or seven reasons not to use uncertainty
1063 analysis, *Water Resour. Res.*, 42(5), doi:10.1029/2005WR004820, 2006.

1064 Peters, W., Jacobson, A. R., Sweeney, C., Andrews, A. E., Conway, T. J., Masarie, K., Miller, J.
1065 B., Bruhwiler, L. M. P., Pétron, G., Hirsch, A. I., Worthy, D. E. J., van der Werf, G. R.,
1066 Randerson, J. T., Wennberg, P. O., Krol, M. C. and Tans, P. P.: An atmospheric perspective
1067 on North American carbon dioxide exchange: CarbonTracker., *Proc. Natl. Acad. Sci. U. S.*
1068 *A.*, 104(48), 18925–30, doi:10.1073/pnas.0708986104, 2007.

1069 Le Quéré, C., Peters, G. P., Andres, R. J., Andrew, R. M., Boden, T. A., Ciais, P., Friedlingstein,
1070 P., Houghton, R. A., Marland, G., Moriarty, R., Sitch, S., Tans, P., Arneeth, A., Arvanitis, A.,
1071 Bakker, D. C. E., Bopp, L., Canadell, J. G., Chini, L. P., Doney, S. C., Harper, A., Harris, I.,
1072 House, J. I., Jain, A. K., Jones, S. D., Kato, E., Keeling, R. F., Klein Goldewijk, K.,
1073 Körtzinger, A., Koven, C., Lefèvre, N., Maignan, F., Omar, A., Ono, T., Park, G.-H., Pfeil,
1074 B., Poulter, B., Raupach, M. R., Regnier, P., Rödenbeck, C., Saito, S., Schwinger, J.,
1075 Segschneider, J., Stocker, B. D., Takahashi, T., Tilbrook, B., van Heuven, S., Viovy, N.,
1076 Wanninkhof, R., Wiltshire, A. and Zaehle, S.: Global carbon budget 2013, *Earth Syst. Sci.*
1077 *Data*, 6(1), 235–263, doi:10.5194/essd-6-235-2014, 2014.

1078 Raich, J. W. J. W., Potter, C. S. C. and Bhagawati, D.: Interannual variability in global soil
1079 respiration, 1980-94, *Glob. Chang. Biol.*, 8, 800–812, doi:10.1046/j.1365-
1080 2486.2002.00511.x, 2002.

1081 Ren, X., He, H., Moore, D. J. P., Zhang, L., Liu, M., Li, F., Yu, G. and Wang, H.: Uncertainty
1082 analysis of modeled carbon and water fluxes in a subtropical coniferous plantation, *J.*
1083 *Geophys. Res. Biogeosciences*, 118(4), 1674–1688, doi:10.1002/2013JG002402, 2013.

1084 Ricciuto, D. M., King, A. W., Dragoni, D. and Post, W. M.: Parameter and prediction uncertainty
1085 in an optimized terrestrial carbon cycle model: Effects of constraining variables and data
1086 record length, *J. Geophys. Res. Biogeosciences*, 116(1), 1–17, doi:10.1029/2010JG001400,
1087 2011.

1088 Richardson, A. D. and Hollinger, D. Y.: Statistical modeling of ecosystem respiration using eddy
1089 covariance data: Maximum likelihood parameter estimation, and Monte Carlo simulation of
1090 model and parameter uncertainty, applied to three simple models, *Agric. For. Meteorol.*,
1091 131(3–4), 191–208, doi:10.1016/j.agrformet.2005.05.008, 2005.

1092 Sadegh, M. and Vrugt, J. A.: Bridging the gap between GLUE and formal statistical approaches:
1093 Approximate Bayesian computation, *Hydrol. Earth Syst. Sci.*, 17(12), 4831–4850,
1094 doi:10.5194/hess-17-4831-2013, 2013.

1095 Scharnagl, B., Vrugt, J. A., Vereecken, H. and Herbst, M.: Inverse modelling of in situ soil water
1096 dynamics: Investigating the effect of different prior distributions of the soil hydraulic
1097 parameters, *Hydrol. Earth Syst. Sci.*, 15(10), 3043–3059, doi:10.5194/hess-15-3043-2011,
1098 2011.

1099 Schimel, J. P. and Weintraub, M. N.: The implications of exoenzyme activity on microbial carbon
1100 and nitrogen limitation in soil: a theoretical model, *Soil Biol. Biochem.*, 35(4), 549–563,

1101 doi:10.1016/S0038-0717(03)00015-4, 2003.

1102 Schmidt, M. W. I., Torn, M. S., Abiven, S., Dittmar, T., Guggenberger, G., Janssens, I. A., Kleber,
1103 M., Kögel-Knabner, I., Lehmann, J., Manning, D. A. C., Nannipieri, P., Rasse, D. P., Weiner,
1104 S. and Trumbore, S. E.: Persistence of soil organic matter as an ecosystem property, *Nature*,
1105 478(7367), 49–56, doi:10.1038/nature10386, 2011.

1106 Scholz, K., Hammerle, A., Hiltbrunner, E. and Wohlfahrt, G.: Analyzing the Effects of Growing
1107 Season Length on the Net Ecosystem Production of an Alpine Grassland Using Model – Data
1108 Fusion, *Ecosystems*, 21(5), 982–999, doi:10.1007/s10021-017-0201-5, 2018.

1109 Schoups, G. and Vrugt, J. A.: A formal likelihood function for parameter and predictive inference
1110 of hydrologic models with correlated, heteroscedastic, and non-Gaussian errors, *Water*
1111 *Resour. Res.*, 46(10), 1–17, doi:10.1029/2009WR008933, 2010.

1112 Scott, R. L., Jenerette, G. D., Potts, D. L. and Huxman, T. E.: Effects of seasonal drought on net
1113 carbon dioxide exchange from a woody-plant-encroached semiarid grassland, *J. Geophys.*
1114 *Res. Biogeosciences*, 114(4), doi:10.1029/2008JG000900, 2009.

1115 Shi, X., Ye, M., Curtis, G. P., Miller, G. L., Meyer, P. D., Kohler, M., Yabusaki, S. and Wu, J.:
1116 Assessment of parametric uncertainty for groundwater reactive transport modeling, *Water*
1117 *Resour. Res.*, 50(5), 4416–4439, doi:10.1002/2013WR013755, 2014.

1118 Sinsabaugh, R. L., Manzoni, S., Moorhead, D. L. and Richter, A.: Carbon use efficiency of
1119 microbial communities: stoichiometry, methodology and modelling, *Ecol. Lett.*, 16(7), 930–
1120 939, doi:10.1111/ele.12113, 2013.

1121 Smith, M. W., Bracken, L. J. and Cox, N. J.: Toward a dynamic representation of hydrological
1122 connectivity at the hillslope scale in semiarid areas, *Water Resour. Res.*, 46(12),
1123 doi:10.1029/2009WR008496, 2010a.

1124 Smith, T., Sharma, A., Marshall, L., Mehrotra, R. and Sisson, S.: Development of a formal
1125 likelihood function for improved Bayesian inference of ephemeral catchments, *Water*
1126 *Resour. Res.*, 46(12), 1–11, doi:10.1029/2010WR009514, 2010b.

1127 Smith, T., Marshall, L. and Sharma, A.: Modeling residual hydrologic errors with Bayesian
1128 inference, *J. Hydrol.*, 528, 29–37, doi:10.1016/j.jhydrol.2015.05.051, 2015.

1129 Spaaks, J. H. and Bouten, W.: Resolving structural errors in a spatially distributed hydrologic
1130 model using ensemble Kalman filter state updates, *Hydrol. Earth Syst. Sci.*, 17(9), 3455–
1131 3472, doi:10.5194/hess-17-3455-2013, 2013.

1132 Steinacher, M. and Joos, F.: Transient Earth system responses to cumulative carbon dioxide
1133 emissions: Linearities, uncertainties, and probabilities in an observation-constrained model
1134 ensemble, *Biogeosciences*, 13(4), 1071–1103, doi:10.5194/bg-13-1071-2016, 2016.

1135 Tang, J. and Riley, W. J.: Weaker soil carbon–climate feedbacks resulting from microbial and
1136 abiotic interactions, *Nat. Clim. Chang.*, 5, 56 [online] Available from:
1137 <http://dx.doi.org/10.1038/nclimate2438>, 2014.

1138 Tang, J. and Zhuang, Q.: A global sensitivity analysis and Bayesian inference framework for
1139 improving the parameter estimation and prediction of a process-based Terrestrial Ecosystem
1140 Model, *J. Geophys. Res. Atmos.*, 114(D15), doi:10.1029/2009JD011724, 2009.

1141 Thiemann, M., Trosset, M., Gupta, H. and Sorooshian, S.: Bayesian recursive parameter estimation
1142 for hydrologic models, *Water Resour. Res.*, 37(10), 2521–2535,
1143 doi:10.1029/2000WR900405, 2001.

1144 Thyer, M., Renard, B., Kavetski, D., Kuczera, G., Franks, S. W. and Srikanthan, S.: Critical
1145 evaluation of parameter consistency and predictive uncertainty in hydrological modeling: A
1146 case study using Bayesian total error analysis, *Water Resour. Res.*, 45(12), 1–22,

1147 doi:10.1029/2008WR006825, 2009.

1148 Tiedeman, C. R. and Green, C. T.: Effect of correlated observation error on parameters,
1149 predictions, and uncertainty, *Water Resour. Res.*, 49(10), 6339–6355,
1150 doi:10.1002/wrcr.20499, 2013.

1151 Tsai, F. T.-C. and Elshall, A. S.: Hierarchical Bayesian model averaging for hydrostratigraphic
1152 modeling: Uncertainty segregation and comparative evaluation, *Water Resour. Res.*, 49(9),
1153 doi:10.1002/wrcr.20428, 2013.

1154 Tucker, C. L., Bell, J., Pendall, E. and Ogle, K.: Does declining carbon-use efficiency explain
1155 thermal acclimation of soil respiration with warming?, *Glob. Chang. Biol.*, 252–263,
1156 doi:10.1111/gcb.12036, 2013.

1157 Tucker, C. L., Young, J. M., Williams, D. G. and Ogle, K.: Process-based isotope partitioning of
1158 winter soil respiration in a subalpine ecosystem reveals importance of rhizospheric
1159 respiration, *Biogeochemistry*, 121, 389–408 [online] Available from:
1160 <http://www.jstor.org/stable/24717586>, 2014.

1161 Tuomi, M., Vanhala, P., Karhu, K., Fritze, H. and Liski, J.: Heterotrophic soil respiration-
1162 Comparison of different models describing its temperature dependence, *Ecol. Modell.*,
1163 211(1–2), 182–190, doi:10.1016/j.ecolmodel.2007.09.003, 2008.

1164 Vargas, R., Carbone, M. S., Reichstein, M. and Baldocchi, D. D.: Frontiers and challenges in soil
1165 respiration research: from measurements to model-data integration, *Biogeochemistry*,
1166 102(1), 1–13, doi:10.1007/s10533-010-9462-1, 2011.

1167 Vrugt, J. A. and Ter Braak, C. J. F.: DREAM(D): An adaptive Markov Chain Monte Carlo
1168 simulation algorithm to solve discrete, noncontinuous, and combinatorial posterior parameter
1169 estimation problems, *Hydrol. Earth Syst. Sci.*, 15(12), 3701–3713, doi:10.5194/hess-15-

1170 3701-2011, 2011.

1171 Vrugt, J. A., ter Braak, C. J. F., Diks, C. G. H. and Schoups, G.: Hydrologic data assimilation using
1172 particle Markov chain Monte Carlo simulation: Theory, concepts and applications, *Adv.*
1173 *Water Resour.*, 51, 457–478, doi:10.1016/j.advwatres.2012.04.002, 2013.

1174 Wang, G., Post, W. M. and Mayes, M. A.: Development of microbial-enzyme-mediated
1175 decomposition model parameters through steady-state and dynamic analyses, *Ecol. Appl.*,
1176 23(1), 255–272, doi:10.1890/12-0681.1, 2013.

1177 Weijs, S. V., Schoups, G. and Van De Giesen, N.: Why hydrological predictions should be
1178 evaluated using information theory, *Hydrol. Earth Syst. Sci.*, 14(12), 2545–2558,
1179 doi:10.5194/hess-14-2545-2010, 2010.

1180 Westerberg, I. K., Guerrero, J. L., Younger, P. M., Beven, K. J., Seibert, J., Halldin, S., Freer, J.
1181 E. and Xu, C. Y.: Calibration of hydrological models using flow-duration curves, *Hydrol.*
1182 *Earth Syst. Sci.*, 15(7), 2205–2227, doi:10.5194/hess-15-2205-2011, 2011.

1183 Wieder, W. R., Bonan, G. B. and Allison, S. D.: Global soil carbon projections are improved by
1184 modelling microbial processes, *Nat. Clim. Chang.*, 3, 909 [online] Available from:
1185 <http://dx.doi.org/10.1038/nclimate1951>, 2013.

1186 Wieder, W. R., Allison, S. D., Davidson, E. A., Georgiou, K., Hararuk, O., He, Y., Hopkins, F.,
1187 Luo, Y., Smith, M. J., Sulman, B., Todd-Brown, K., Wang, Y.-P., Xia, J. and Xu, X.:
1188 Explicitly representing soil microbial processes in Earth system models, *Global*
1189 *Biogeochem. Cycles*, 29(10), 1782–1800, doi:10.1002/2015GB005188, 2015.

1190 Van Wijk, M. T., Van Putten, B., Hollinger, D. Y. and Richardson, A. D.: Comparison of different
1191 objective functions for parameterization of simple respiration models, *J. Geophys. Res.*
1192 *Biogeosciences*, 113(3), 1–11, doi:10.1029/2007JG000643, 2008.

1193 Xu, T., White, L., Hui, D. and Luo, Y.: Probabilistic inversion of a terrestrial ecosystem model:
1194 Analysis of uncertainty in parameter estimation and model prediction, *Global Biogeochem.*
1195 *Cycles*, 20(2), 1–15, doi:10.1029/2005GB002468, 2006.

1196 Xu, X., Schimel, J. P., Thornton, P. E., Song, X., Yuan, F. and Goswami, S.: Substrate and
1197 environmental controls on microbial assimilation of soil organic carbon: a framework for
1198 Earth system models, *Ecol. Lett.*, 17(5), 547–555, doi:10.1111/ele.12254, 2014.

1199 Yeluripati, J. B., van Oijen, M., Wattenbach, M., Neftel, A., Ammann, A., Parton, W. J. and Smith,
1200 P.: Bayesian calibration as a tool for initialising the carbon pools of dynamic soil models,
1201 *Soil Biol. Biochem.*, 41(12), 2579–2583, doi:10.1016/j.soilbio.2009.08.021, 2009.

1202 Yuan, W., Liang, S., Liu, S., Weng, E., Luo, Y. and Hollinger, D.: Improving model parameter
1203 estimation using coupling relationships between vegetation production and ecosystem
1204 respiration, *Ecol. Modell.*, 240, 29–40, doi:10.1016/j.ecolmodel.2012.04.027, 2012.

1205 Yuan, W., Xu, W., Ma, M., Chen, S. and Liu, W.: Agricultural and Forest Meteorology Improved
1206 snow cover model in terrestrial ecosystem models over the Qinghai – Tibetan Plateau, *Agric.*
1207 *For. Meteorol.*, 218–219, 161–170, doi:10.1016/j.agrformet.2015.12.004, 2016.

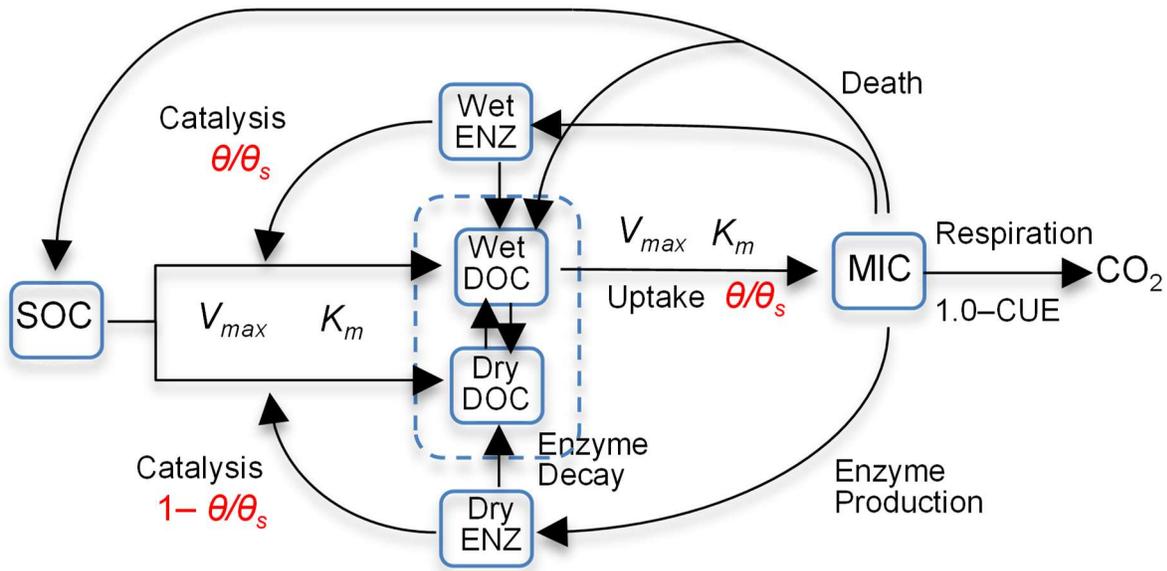
1208 Zhang, X., Niu, G.-Y., Elshall, A. S., Ye, M., Barron-Gafford, G. A. and Pavao-Zuckerman, M.:
1209 Assessing five evolving microbial enzyme models against field measurements from a
1210 semiarid savannah - What are the mechanisms of soil respiration pulses?, *Geophys. Res.*
1211 *Lett.*, 41(18), doi:10.1002/2014GL061399, 2014.

1212 Zhou, X., Luo, Y., Gao, C., Verburg, P. S. J., Arnone, J. A., Darrouzet-Nardi, A. and Schimel, D.
1213 S.: Concurrent and lagged impacts of an anomalously warm year on autotrophic and
1214 heterotrophic components of soil respiration: A deconvolution analysis, *New Phytol.*, 187(1),
1215 184–198, doi:10.1111/j.1469-8137.2010.03256.x, 2010.

1216

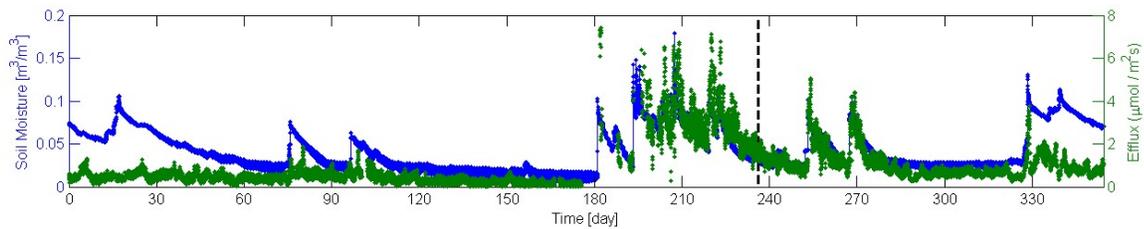
1217 Figure 1. Diagram of model 6C representing the processes of (1) degradation of soil organic carbon
 1218 (SOC) to dissolved organic carbon (DOC) through catalysis of enzymes (ENZ) produced by
 1219 microbes (MIC), (2) MIC uptake of DOC, and (3) microbial (MIC) respiration to produce CO₂
 1220 (CUE is the carbon use efficiency). SOC degradation and microbial uptake rates are controlled by
 1221 water saturation (θ / θ_s). The DOC and ENZ pools are split into two subpools, one for the wet
 1222 zone and the other for the dry zone of the soil pore space. Microbial uptake of DOC occurs only
 1223 in the wet zone, and the uptake rate is linearly related to θ / θ_s . Catalysis through ENZ in the wet
 1224 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
 1225 maximum rate, and K_m is the half-saturation concentration.

1226



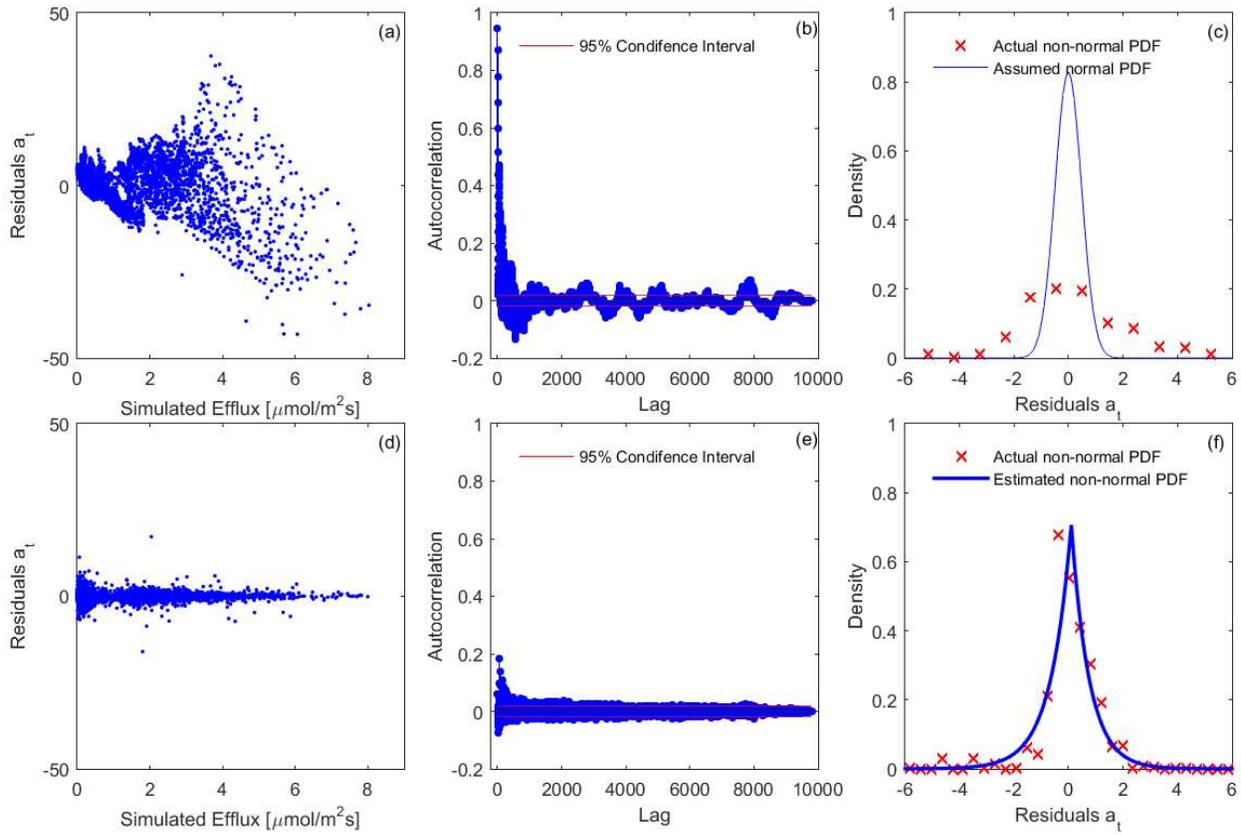
1227
1228

1229 Figure 2. Time series of soil moisture and efflux observations. The dashed line marks the divide
1230 of the dataset into calibration and validation periods.
1231



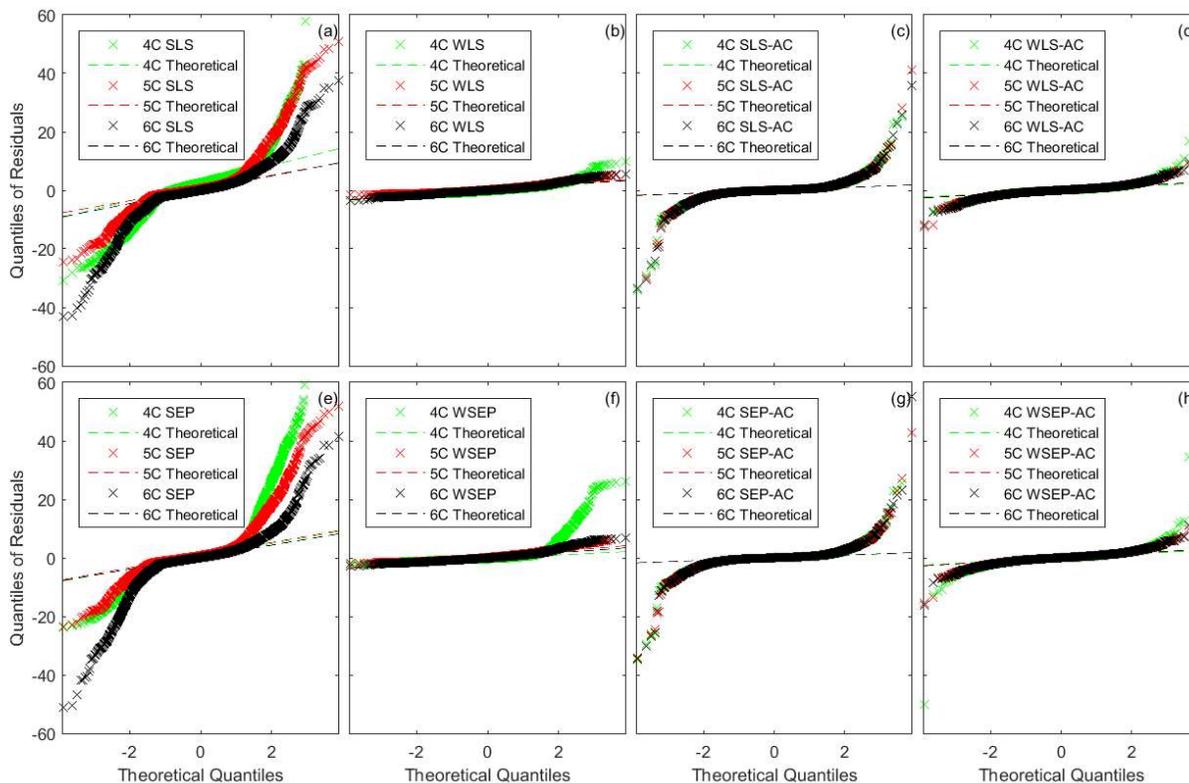
1232

1233 Figure 3. Residual analysis of the best realization (among multiple MCMC realizations) for model
1234 6C using data models (a-c) SLS and (d-f) WSEP-AC.
1235



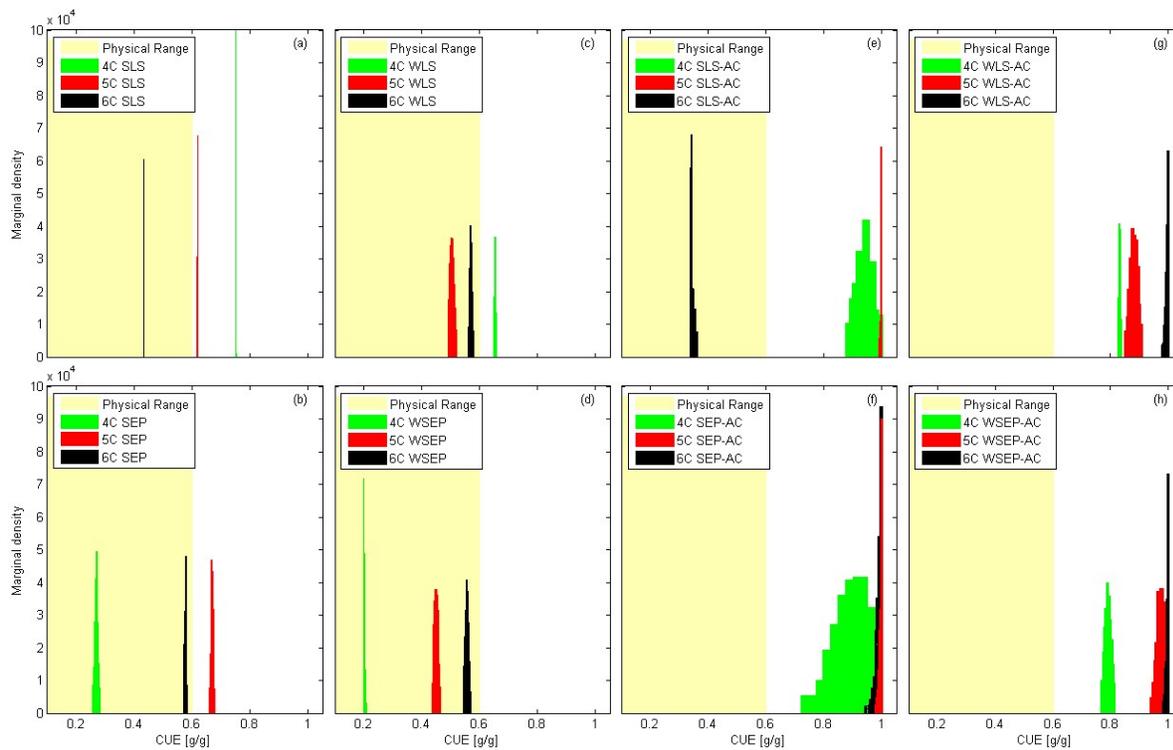
1236

1237 Figure 4. Residual quantile-quantile (Q-Q) plots of the best realization (among multiple MCMC
1238 realizations) for the three soil respiration models and eight data models.
1239



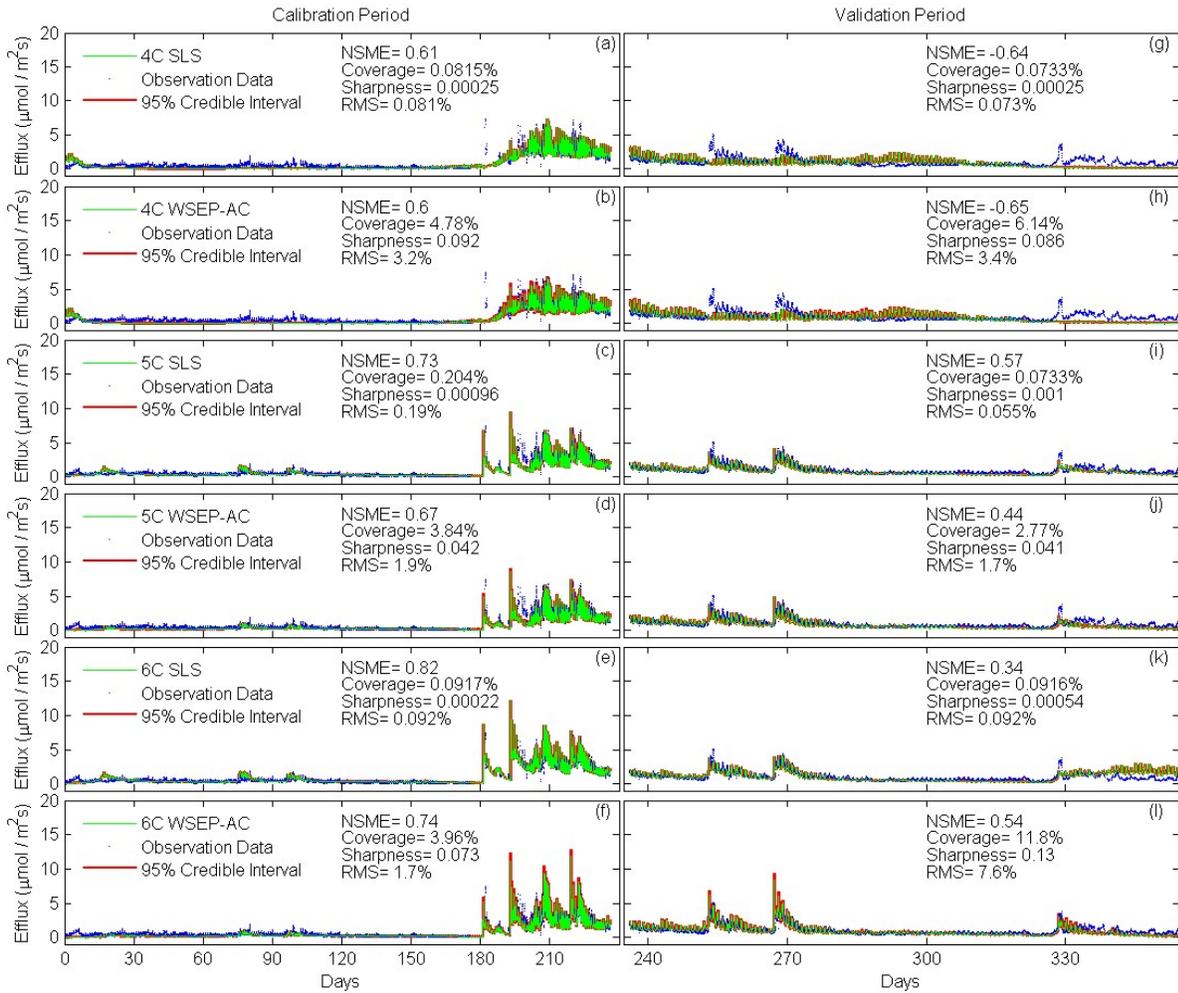
1240

1241 Figure 5. Marginal posterior parameter density of carbon use efficiency (CUE) for the three soil
 1242 respiration models and eight data models.
 1243



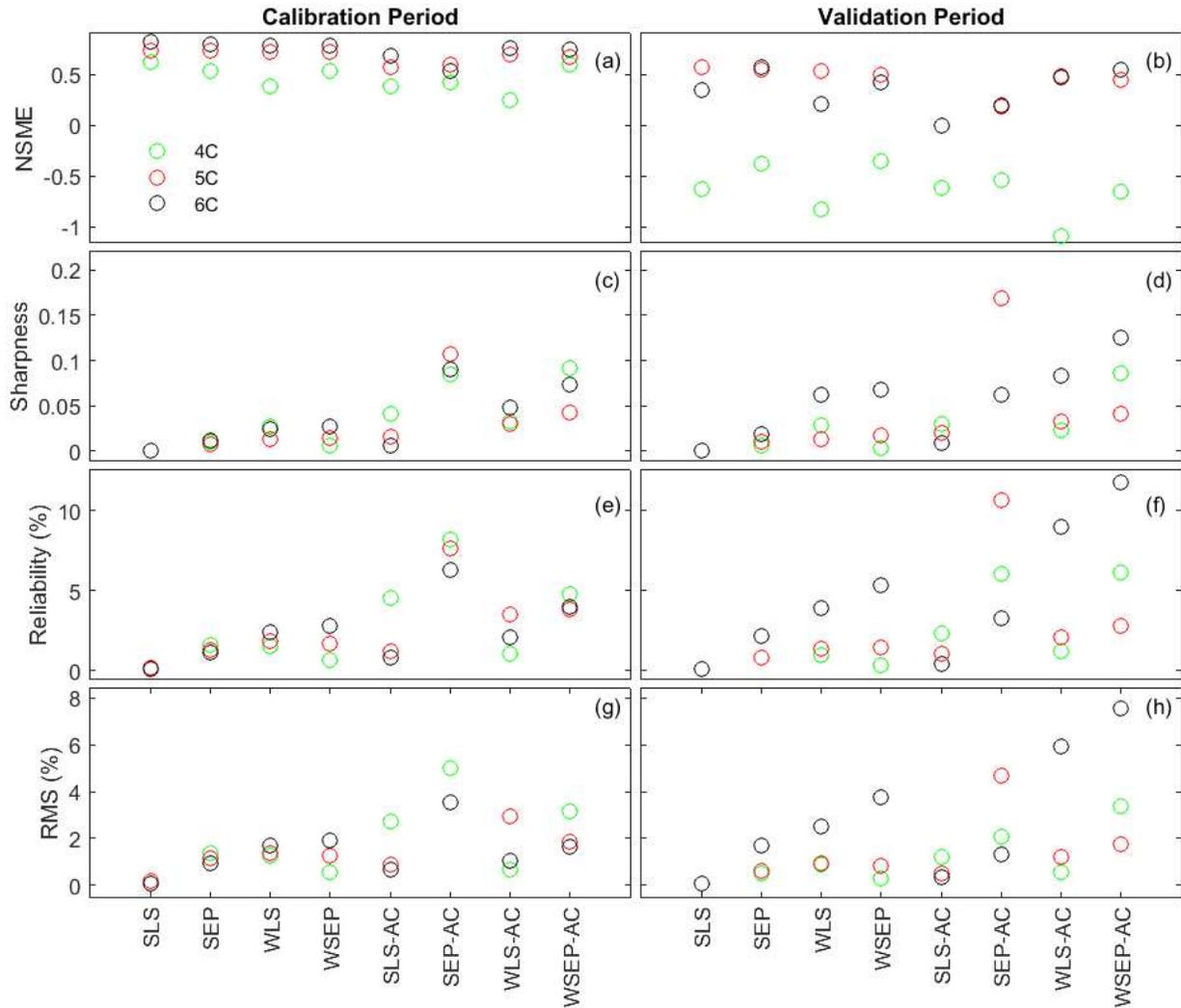
1244

1245 Figure 6. Observation data (blue dots) and mean prediction (green line) and 95% credible intervals
 1246 (red line) of prediction ensembles for (a)-(f) the calibration period and (g)-(l) the validation period.
 1247 The plots are for the three soil respiration models using data models SLS and WSEP-AC. *The*
 1248 *prediction ensembles are generated to consider parametric uncertainty of the soil respiration*
 1249 *models only.*
 1250



1251

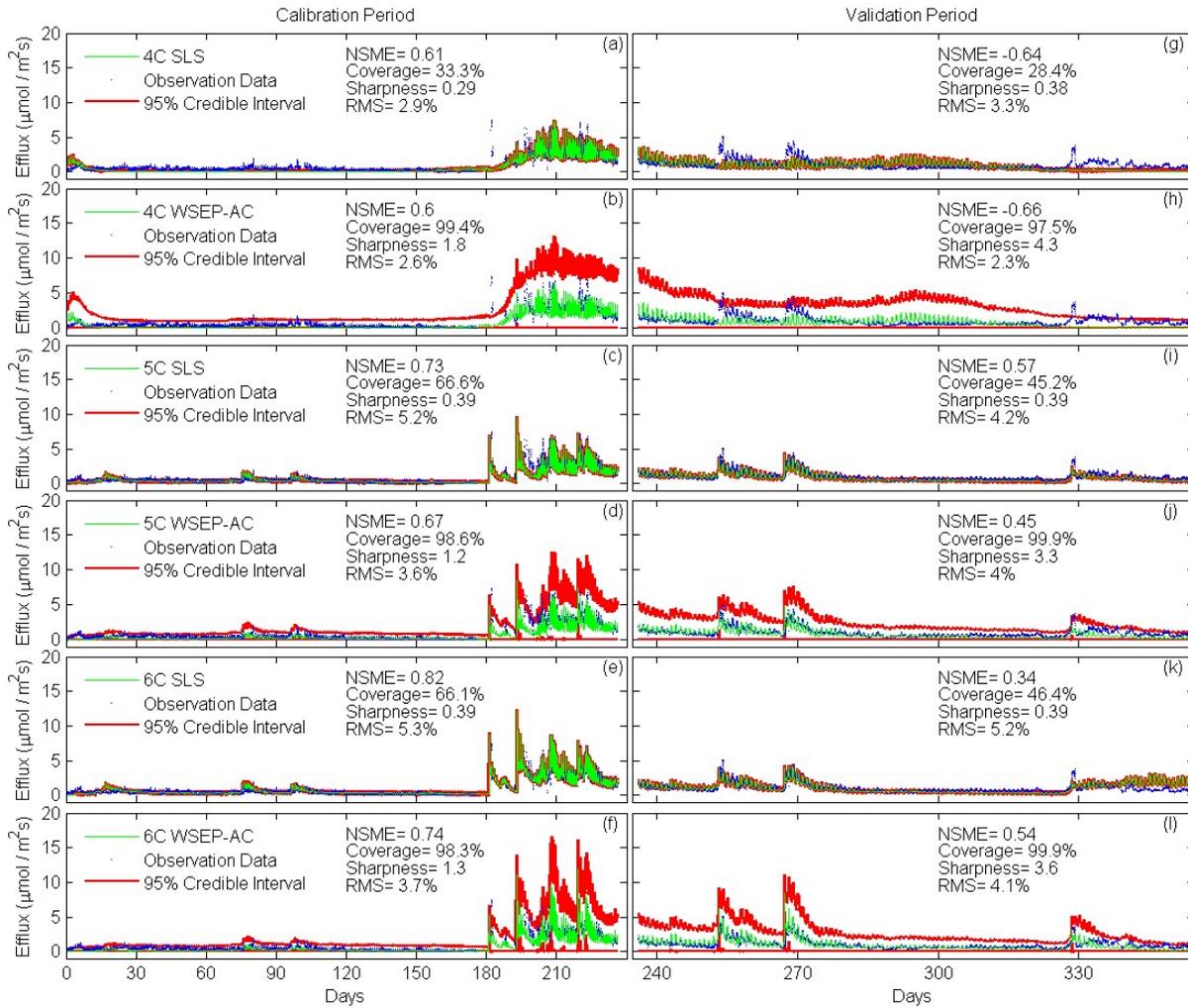
1252 Figure 7. (a-b) Nash-Sutcliffe model efficiency (NSME), (c)-(d) sharpness, (e)-(f) predictive
 1253 coverage, and (g)-(h) relative model score for measuring predictive performance of the three soil
 1254 respiration models and the eight data models during the calibration and cross-validation periods.
 1255 *The statistics are evaluated from the prediction ensembles generated to consider parametric*
 1256 *uncertainty of the soil respiration models only.*
 1257



1258

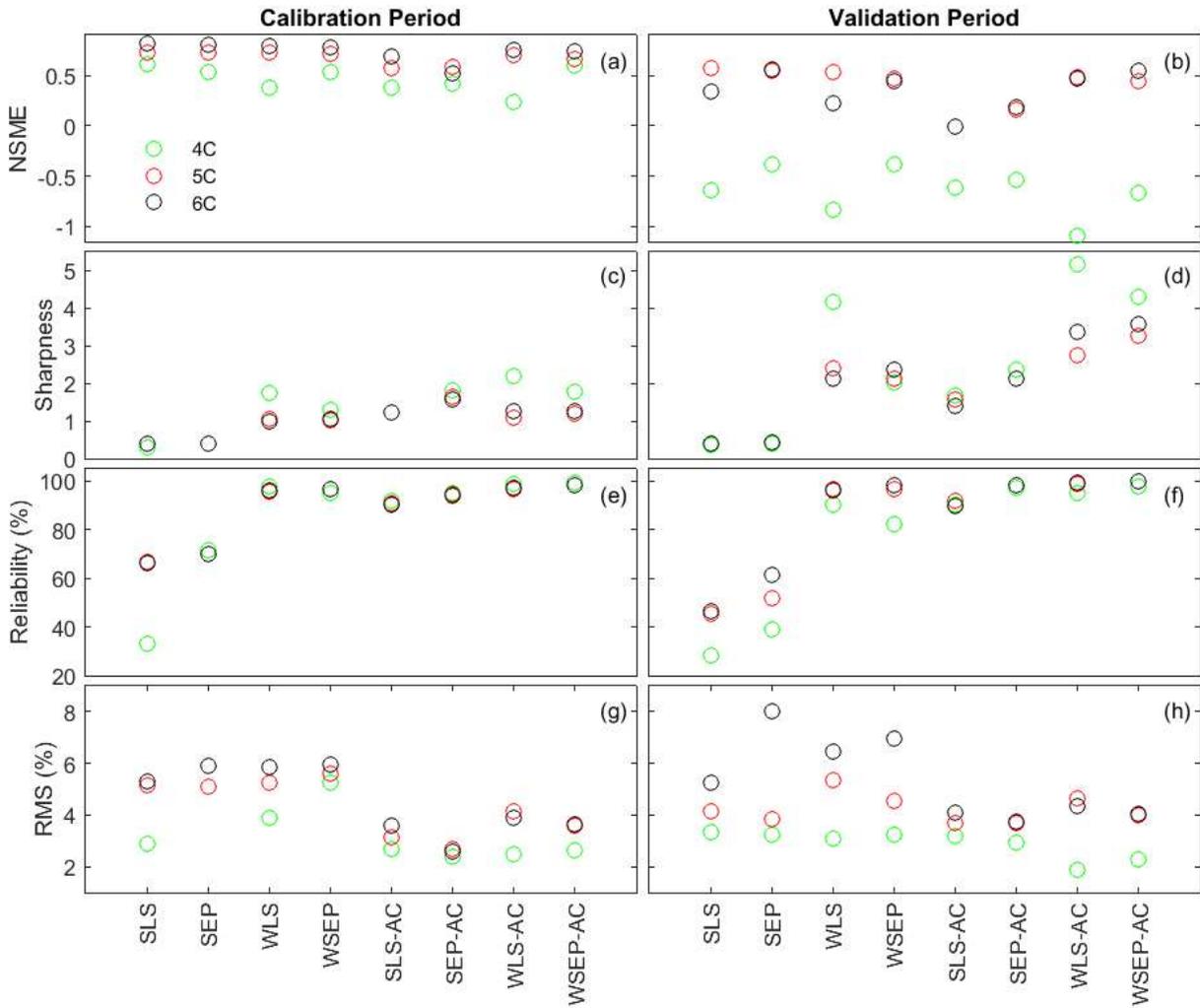
1259

1260 Figure 8. Observation data (blue dots) and mean prediction (green line) and 95% credible intervals
 1261 (red line) of prediction ensembles for (a)-(f) the calibration period and (g)-(l) the validation period.
 1262 The plots are for the three soil respiration models using data models SLS and WSEP-AC. *The*
 1263 *prediction ensembles are generated to consider parametric uncertainty of not only the soil*
 1264 *respiration models but also the data models.*
 1265



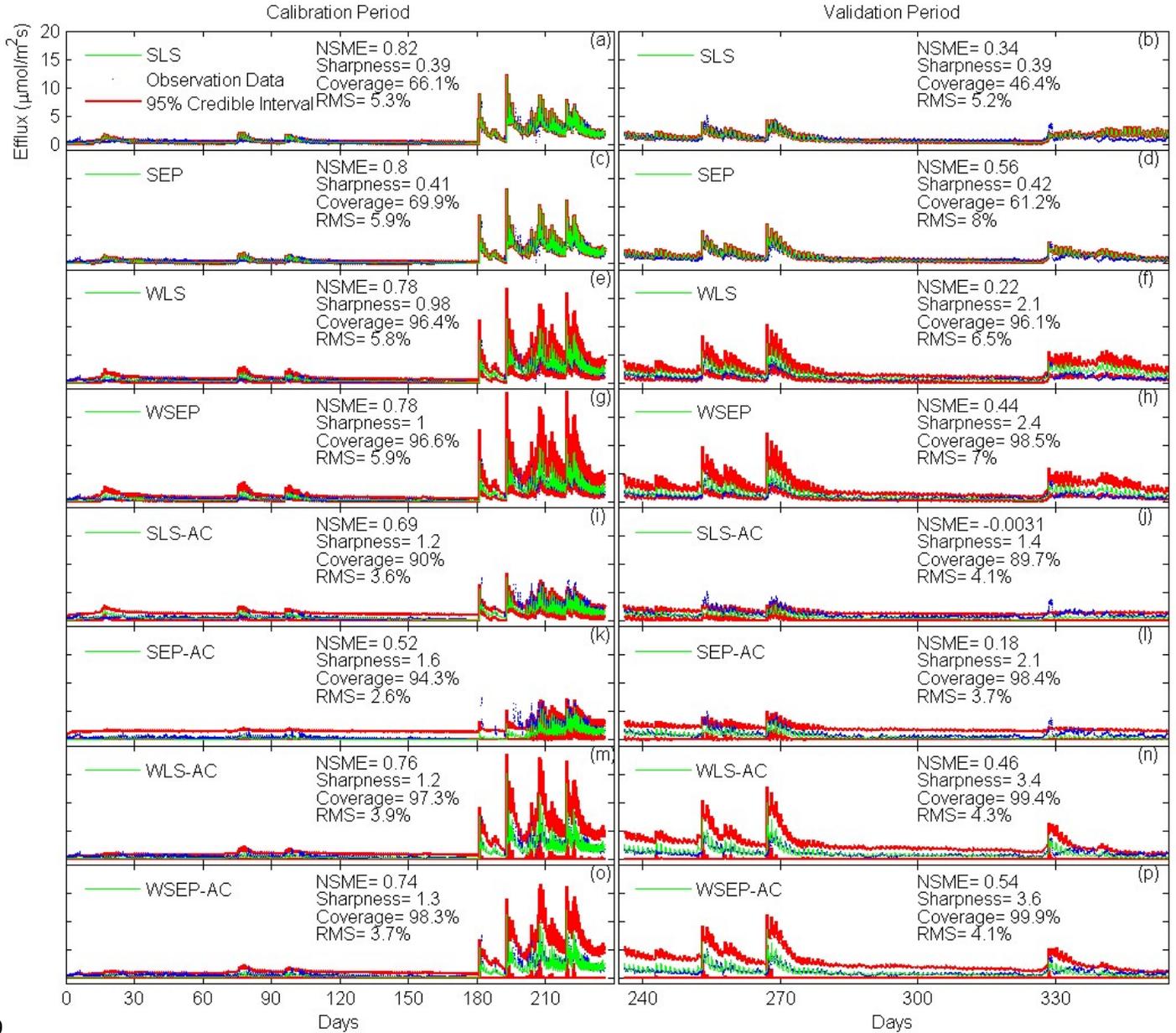
1266

1267 Figure 9. (a-b) Nash-Sutcliffe model efficiency (NSME), (c)-(d) sharpness, (e)-(f) predictive
 1268 coverage, and (g)-(h) relative model score for measuring predictive performance of the three soil
 1269 respiration models and the eight data models during the calibration and cross-validation periods.
 1270 *The statistics are evaluated from the prediction ensembles generated to consider parametric*
 1271 *uncertainty of not only the soil respiration models but also the data models.*
 1272



1273

1274 Figure 10. Observation data (blue dots) and mean prediction (green line) and 95% credible
 1275 intervals (red line) for 6C for the eight likelihood functions during the calibration period (a)-(h)
 1276 and the validation period (i)-(p). The prediction ensembles are generated to consider parametric
 1277 uncertainty of not only the soil respiration models but also the data models.
 1278



1279