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 Geosciences, University of Hawaii Manoa, Honolulu, Hawaii, USA
7  [2] Water Resources Research Center, University of Hawaii Manoa, Honolulu, Hawaii, USA
8  [3] Department of Earth, Ocean, and Atmospheric Science, Florida State University, Tallahassee,
9  Florida
10 [4] Biosphere 2, University of Arizona, Tucson, Arizona
11 [5] Department of Hydrology and Water Resources, University of Arizona, Tucson, Arizona
12 [6] School of Geography and Development, University of Arizona, Tucson, Arizona
13
14
15 *Corresponding Author: Ming Ye, Telephone: (850) 644-4587, Email: mye@fsu.edu
16
17
18
19
20
21
22
23
24
25
26
27
28
29
30
31
32
33
34
35
36
37
38
39
40
41
42
43
44
45
46

Geoscientific
Model Development
Discussions

**Key Points**

47

48

49  (1)  Bayesian inference and prediction are useful to evaluate multiple soil respiration models

50       with different levels of 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

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 $CO_2$ 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 models.

## 1   Introduction

94

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 $CO_2$ 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

109   characterizing residuals, i.e., the difference between data and corresponding model simulations.

110   While a large number of data models have been used, to our knowledge, comprehensive and

111   systematic evaluation of data models for soil respiration models has not been reported in literature.

112   The goal of this study is to evaluate the impacts of data models on Bayesian inference and

113   predictive performance of three mechanistic soil respiration models, and use these findings to

114   make broader recommendations. The three models were developed by Zhang et al. (2014) to

115   simulate the Birch effect (the peak soil microbial respiration pulses in response to episodic rainfall

116   pulses) at a site scale and a short temporal scale, which are important for gaining mechanistic

Geoscientific
Model Development
Discussions

117   understanding of $CO_2$ efflux production (Högberg and Read, 2006; Vargas et al., 2011). Zhang et

118   al. (2014) developed a total five models, including an existing four-carbon pool model and four

119   new models with additional carbon pools and/or explicit representations of soil moisture controls

120   on carbon degradation and microbial uptake rates. The models Zhang et al. (2014) were calibrated,

121   and Bayesian model selection was used to select and the best model. However, this effort was

122   based on a single data model. It is unknown whether the best model still remains the best (in terms

123   of reproducing the both calibration data and the cross-validation data) if a different data model is

124   used. In addition, since predictive performance of the models was not evaluated in Zhang et al.

125   (2014), it is unknown whether the best model will give the best predictions. These two questions

126   are addressed in this study by considering eight data models and by evaluating predictive

127   performance in a manner of cross-validation. The top two models (also the two most high fidelity

128   models) ranked by Zhang et al. (2014) are considered in this study, and the worst model (also the

129   low fidelity model) is also considered in this study for comparison. Model fidelity refers to the

130   degree of realism of representing our scientific knowledge with respect to the real world system.

131   That is model with less discrepancy. Conducting Bayesian inference and evaluating predictive

132   performance for the three models with different degrees of fidelity provides more insights than for

133   a single model.

134       Bayesian inference in general uses the Bayes' theorem to update the distributions of model

135   parameters to posterior parameter distributions given a likelihood function. The mathematical

136   formulation of the (formal and informal) likelihood function requires a probabilistic data model

137   that however is intrinsically unknown due to unknown errors in all model components such as

138   observation data, model structures, parameters, and driving forces. Bayesian inference of soil

139   respiration models often adopts the assumption of independent, normally distributed and

140  homoscedastic residuals (e.g. Ahrens et al., 2014; Barr et al., 2013; Hararuk et al., 2014;

141  Hashimoto et al., 2011; Klemedtsson et al., 2008; Raich et al., 2002; Ren et al., 2013; Ricciuto et

142  al., 2011; Richardson and Hollinger, 2005; Steinacher and Joos, 2016; Tucker et al., 2014; Tuomi

143  et al., 2008; Xu et al., 2006; Yeluripati et al., 2009; Zhang et al., 2014; Zhou et al., 2010). These

144  assumptions are conveniently adopted since the requirement of using an unknown probability

145  model in Bayesian statistics is called "a basic dilemma" by Box and Tiao (1992). Postulating the

146  data models is always based on assumptions about residual statistics, and the most widely used

147  assumptions are paired as follows: (i) independent vs. correlated residuals, (ii) homoscedastic vs.

148  heteroscedastic residuals, and (iii) Gaussian vs. non-Gaussian residuals.  There are many

149  diagnostics available to assess these choices (a number of them is used in this paper). However,

150  few studies have focused on investigating appropriateness of the assumptions for soil respiration

151  modeling by relaxing the independent residuals assumption (Chevallier and O'Dell, 2013; Ricciuto

152  et al., 2011) and the Gaussian residuals assumption ( Ricciuto et al., 2011; Van Wijk et al., 2008).

153      This study evaluates the above assumptions by considering eight data models which relaxes

154  these three assumptions stepwise as shown in Section 2. For example, combining the assumptions

155  of independent, homoscedastic, and Gaussian residuals leads to the standard least squares data

156  model. This model is the simplest one among the eight data models, since it requires only one

157  parameter, i.e., the constant variance of the Gaussian distribution. Note that there is a difference

158  between the physical model parameters and data model parameters. They technically can be

159  estimated together, but one arises from assumptions about process, and the other assumptions

160  about the data models. Relaxing the homoscedastic assumption to heteroscedastic gives the

161  weighted least squares data model. It is more complex, because it requires multiple variances for

162  multiple data. Whenever one or combinations of the three assumptions (independence,

163    homoscedasticity, and normality) are relaxed, the resulting data models become more complex and

164    require more parameters. This systematic way of formulating data models is similar to that of

165    Smith et al. (2010b, 2015), and it is necessary to evaluate appropriateness of the three basic

166    assumptions and their impacts on Bayesian inference.

167         The assumptions of heteroscedastic, correlated, and non-Gaussian residuals are accounted for

168    using the method of Schoups and Vrugt (2010) in the following procedure: (i) the correlation is

169    removed from the residuals by using an autoregressive model; (ii) the resulting residuals are

170    normalized by a linear model of variance; and (iii) the normalized residuals are characterized by

171    using the skew exponential power distribution. The data model parameters (i.e., coefficients of the

172    autoregressive model, the linear variance model, and the skew exponential power distribution) are

173    not specified by users, but estimated together with physical model parameters during the Bayesian

174    inference. The skew exponential power distribution is general in that by adjusting the values of its

175    kurtosis and skewness parameters the distribution can produce other distributions such as the

176    Laplace distribution used by (Van Wijk et al., 2008) and (Ricciuto et al., 2011), and other

177    distributions given by using different kurtosis parameters of an exponential model (Tang and

178    Zhuang, 2009). It is worth pointing out that there exist other methods to account for the three

179    assumptions. Evin et al. (2013) suggested accounting for residual heteroscedasticity before

180    accounting for residual autocorrelation. Lu et al. (2013) developed an iterative two-stage procedure

181    to separately estimate physical model parameters and data model parameters. Evin et al. (2014)

182    developed a similar procedure to first estimate model parameters and then estimate

183    heteroscedasticity and autocorrelation parameters. While this study uses the method of Schoups

184    and Vrugt (2010), exploring other methods is warranted in future studies.

185    After investigating the impacts of the data models on Bayesian inference, this study evaluates

186    the impacts of the data models on predictive performance of the three soil respiration models.

187    Using random samples generated during the Bayesian inference, a prediction ensemble is produced

188    for each soil respiration model. The ensemble is used to evaluate predictive performance of the

189    models in a stochastic sense by estimating to what extent the models can predict future events. The

190    evaluation in this study is done in a cross-validation manner to split a dataset of $CO_2$ efflux into

191    two parts for Bayesian inference and cross-validation, respectively. The evaluation of predictive

192    performance is important because different data models may give different parameter distributions

193    and accordingly different predictive performance. For example, the study of van Wijk et al. (2008)

194    concluded that the choice of the residual function is crucial to achieve accurate model prediction

195    and parameter estimation. Shi et al. (2014) showed that the posterior parameter distributions and

196    predictive performance given by two data models (weighted least square and skew exponential

197    power distribution after removing heteroscedasticity and autocorrelation) are dramatically

198    different, and a definitive conclusion was drawn that one data model is better than the other. The

199    evaluation of predictive analysis is conducted for the following two cases: (1) the prediction

200    ensemble is generated by random samples of the soil respiration models only (i.e. credible

201    interval), and (2) the prediction ensemble is generated by random samples of not only the soil

202    respiration models but also the data models (i.e. predictive interval). The two cases lead to different

203    conclusions about the predictive performance. It is expected that the evaluation of predictive

204    performance conducted in this study can help select the most appropriate data model to achieve

205    optimal model predictions.

206    The remainder of the paper is organized as follows. Section 2 starts with a description of the

207    evolving data models and their corresponding likelihood functions used in Bayesian inference,

208 followed by a brief summary of the three soil respiration models. The results of Bayesian inference

209 are discussed in Section 3 and Section 4, addressing the data model implications on parameter

210 estimation and predictive performance, respectively. Section 5 summarizes the key findings and

211 limitations of this study, and provides recommendations for approaching data model selection.

212 **2    Methodology**

213   This section starts with a descriptions of the eight data models that account for the three pairs

214 of assumptions about residuals in a stepwise manner in Section 2.1. The data models are used to

215 build the likelihood functions used in Section 2.2 for Bayesian inference. The three soil respiration

216 models and observations of $CO_2$ efflux are described in Sections 2.3 and 2.4, respectively.

217 **2.1    Data models**

218   This study considers eight evolving data models starting from a data model that assumes

219 independent, homoscedastic, and Gaussian residuals to a data model that relaxes all the three

220 assumptions. The eight data models are based on the generic normalized residual,

221 $$a_t = \frac{\varepsilon_t}{\sigma_t} \qquad a_t \sim X \,, \tag{1}$$

222 where $\varepsilon_t = d_t - E_t$ is the residual (the difference between data $d_t$ and its corresponding model

223 simulation $E_t$) at time or location $t$, $\sigma_t$ is the standard deviation of the residual, and $X$ is the

224 probability density function (PDF) of $a_t$. The eight data models are formulated with different forms

225 of $\varepsilon_t$, $\sigma_t$, and $X$. The standard least square (SLS) data model is

226 $$a_t = \frac{\varepsilon_t}{\sigma_0} \qquad a_t \sim N(0,1) \,, \tag{2}$$

227 where $\sigma_t = \sigma_0$ is a constant for all the data (i.e., homoscedasty), and $X$ is the standard normal

228 distribution, $N(0,1)$. The unknown parameter $\sigma_0$ is estimated jointly with unknown physical

229  model parameters. If $\sigma_t$ is not a constant (i.e., heteroscedasty), SLS becomes the weighted least

230  squared (WLS) data model. While heteroscedasty can be accounted for through residuals

231  transformation (e.g. Thiemann et al., 200; Smith et al., 2010b) or other similar approaches (Gragne

232  et al., 2015) a linear heteroscedastic model $\sigma_t = \sigma_0 + \sigma_1 E_t$ is assumed following other studies

233  (Thyer et al., 2009; Schoups and Vrugt, 2010; Evin et al., 2013, 2014). With the linear model,

234  there is no need to estimate $\sigma_t$ for each data. Instead, $\sigma_t$ is calculated by estimating only two

235  parameters, $\sigma_0$ and $\sigma_1$. The WSL data model is written as

236
$$a_t = \frac{\varepsilon_t}{\sigma_0 + \sigma_1 E_t} \qquad a_t \sim N(0,1) \,. \tag{3}$$

237  The two unknown parameters $\sigma_0$ and $\sigma_1$ are estimated jointly with unknown physical model

238  parameters. The linear model assigns smaller weight to the data with larger simulation, $E_t$. If the

239  simulation is small and $\sigma_0 \gg \sigma_1 E_t$, the weight becomes constant for all data. Both SLS and WLS

240  assume that $a_t$ is independently and identically distributed.

241      It is not uncommon that residuals are correlated in space and time, due to propagation of

242  measurement errors (Tiedeman and Green, 2013) and model structure errors (Evin et al., 2014;

243  Kavetski et al., 2013; Lu et al., 2013). The temporal correlation that occurs in the numerical

244  example of this study can be accounted for using a *p*-order autoregressive model. This leads to the

245  data model of standard least square with autocorrelation (SLS-AC),

246
$$a_t = \frac{\varepsilon_t - \sum_{i=1}^{p} \phi_i \varepsilon_{t-i}}{\sigma_0} \qquad a_t \sim N(0,1) \tag{4}$$

247    where $p$ is the order of autocorrelation, and $\phi_i$ is an autocorrelation coefficient. The unknown $\phi_i$

248    and $\sigma_0$ are estimated together with unknown model parameters. By extending the concept of

249    correlated residuals to WLS leads to the weight least square with autocorrelation (WLS-AC),

250    $$a_t = \frac{\varepsilon_t - \sum_{i=1}^{p} \phi_i \varepsilon_{t-1}}{\sigma_0 + \sigma_1 E_t} \qquad a_t \sim N(0,1) \tag{5}$$

251    The unknown parameters of $\sigma_0$, $\sigma_1$, and $\phi_i$ are estimated jointly with physical model

252    parameters. Equations (2) – (5) assume that the residuals are Gaussian.

253    The next four data models are similar to the previous four models except that the standard

254    normal distribution of $a_t$ is replaced by the skew exponential power distribution, $SEP(0,1,\xi,\beta)$,

255    (Schoups and Vrugt, 2010)

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

257    where zero is mean, one is standard deviation, $\xi$ is skewness, $\beta$ is kurtosis,

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

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

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

261    kurtosis parameter $\{\beta \in \mathbb{R} : -1 \leq \beta \leq 1\}$ determines the peakness of the pdf such that the $\beta$ values

262    of -1, 0, and 1 give uniform, Gaussian and Laplace distributions, respectively. The skewness

263    parameter $\{\xi \in \mathbb{R} : 0.1 \leq \xi \leq 10\}$ determines the skewness of the pdf such that the $\xi$ values of 0.1,

264    1, and 10 give positively skewed, symmetric, and negatively skewed distributions, respectively.

Geoscientific
Model Development
Discussions

265    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

266    skew exponential power distribution $SEP(0,1,\xi=1,\beta=0)$ becomes the standard normal

267    distribution,

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

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

270        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

271    SEP, WSEP, SEP-AC, and WSEP-AC as follows,

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

273    $$a_t = \frac{\varepsilon_t}{\sigma_0 + \sigma_1 E_t} \qquad a_t \sim SEP(0,1,\xi,\beta).$$    (9)

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

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

276    In comparison with the Gaussian data models, the SEP-based data models have two more

277    parameters ($\xi$ and $\beta$) to be estimated jointly with physical model parameters. WSEP-AC data

278    model, which is known as the generalized likelihood function, is the most commonly used SEP-

279    based data model (e.g. Vrugt and Ter Braak, 2011; Hublart et al., 2016).

280    **2.2    Bayesian inference and likelihood functions**

281        Consider a Bayesian inference problem for a nonlinear model, $f$, used to simulate state

282    variables (e.g., $CO_2$ efflux), $\boldsymbol{d} = f(\boldsymbol{\theta}) + \boldsymbol{\varepsilon}$, where $\boldsymbol{d}$ is a vector of data, $\boldsymbol{\theta}$ is a vector of model

283    parameters, and $\boldsymbol{\varepsilon}$ is a vector of residuals that may include errors in data, model parameters, and

284    model structures. The goal of Bayesian inference is to estimate the posterior distributions, $p(\boldsymbol{\theta}|\boldsymbol{d})$,

285    of model parameters, $\boldsymbol{\theta}$, given data, $\boldsymbol{d}$, using Bayes' theorem (Box and Tiao, 1992)

286 
$$p(\boldsymbol{\theta}\,|\,\boldsymbol{d}) = \frac{p(\boldsymbol{d}\,|\,\boldsymbol{\theta})\,p(\boldsymbol{\theta})}{\int p(\boldsymbol{d}\,|\,\boldsymbol{\theta})\,p(\boldsymbol{\theta})\,d\boldsymbol{\theta}} \tag{12}$$

287    where $p(\boldsymbol{\theta})$ is the prior distribution, and $p(\boldsymbol{d}|\boldsymbol{\theta})$ is the likelihood function to measure goodness-of-

288    fit between model simulations, $f(\boldsymbol{\theta})$, and data, $\boldsymbol{d}$. The prior distribution can be obtained from data

289    of previous studies or expert judgment. When prior information is lacking, a common practice is

290    to assume uniform distributions with relatively large parameter ranges so that the prior

291    distributions do not affect the estimation of posterior distributions.

292      The data models above can be used to construct the likelihood functions. For the Gaussian data

293    models given in equations (2) – (5), the corresponding Gaussian likelihood functions are

294    straightforward, and an example is equation (7). For the SEP data models, the corresponding

295    likelihood that is called generalized likelihood function is (Schoups and Vrugt, 2010)

296 
$$p(\boldsymbol{d}\,|\,\boldsymbol{\theta}) = p(\boldsymbol{\varepsilon}_t\,|\,\boldsymbol{\theta}) = \prod_{t=1}^{n} \sigma_t^{-1} \frac{2\sigma_\xi}{\xi + \xi^{-1}} \omega_\beta \exp\left(-c_\beta \left|a_{\xi,t}\right|^{2/(1+\beta)}\right). \tag{13}$$

297    where $n$ is the dimension of $\boldsymbol{d}$. The Gaussian likelihood functions are special case of the generalized

298    likelihood functions. For example, by setting $\beta = 0$, $\xi = 1$, $\phi_i = 0$, $\sigma_t = \sigma_0$, $\sigma_\xi = 1$, $\mu_\xi = 0$,

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

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

301    function of the WLS data model.

302      In this study, the distributions of the data model parameters are obtained jointly with the

303    physical model parameters using the MT-DREAM$_{(ZS)}$ code (Laloy and Vrugt, 2012), which

Geoscientific
Model Development
Discussions

304    implements a Markov chain Monte Carlo (MCMC) algorithm by running multiple Markov chains

305    in parallel with discrete proposal distribution, multiple-try sampling,  and sampling from an

306    archive of past states. These state-of-the-art features assist in overcoming common challenges in

307    the sampling landscape such as multimodality, ill-conditioning, and high dimensionality, and thus

308    allow for accurate exploration of the targeted distributions.

309    **2.3    Soil respiration models**

310       Zhang et al. (2014) studied the Birch effect (the peak soil microbial respiration pulses in

311    response to episodic rainfall pulses), and developed five models, evolving from an existing four-

312    carbon pool model to models with additional carbon pools and/or explicit representations of soil

313    moisture controls on carbon degradation and microbial uptake rates. Three of the five models are

314    used in this study, and they are dented as 4C, 5C, and 6C. Note that model 4C is model 4C_NOSM

315    of Zhang et al. (2014), not their model 4C. Figure 1 is the diagram of model 6C, the most complex

316    one among the five models. The simplest one, model 4C, has four carbon pools, i.e., soil organic

317    carbon (SOC), dissolved organic carbon (DOC), microbial biomass (MIC), and enzymes (ENZ),

318    and does not consider the soil moisture control on carbon degradation and microbial uptake rates.

319    Models 5C and 6C has an explicit representation of soil moisture controls on the rates. Based on

320    the dual Arrhenius and Michaelis–Menten kinetics model, the original SOC degradation rate,

321    $V_{decom}$, is (Davidson et al., 2011; Davidson and Janssens, 2006)

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

323    where $V_{\max}$ [s$^{-1}$] is the maximum SOC degradation rate per unit enzyme when the substrates is not

324    limiting, $C_{ENZ}$ [gCm$^{-3}$] is enzyme pool size, $C_{SOC}$ [gCm$^{-3}$] is SOC pool size,  and $K_m$ is the half-

325    saturation for SOC. The original microbial uptake rate, $V_{uptake}$, is (Davidson et al., 2011; Davidson

326    and Janssens, 2006)

327    $$V_{uptake} = V_{max\_up} C_{MIC} \frac{C_{DOC}}{K_{m\_up} + C_{DOC}} \frac{C_{O2}}{K_{m\_upO2} + C_{O2}},$$    (15)

328    where $V_{max\_up}$ [s$^{-1}$] is the maximum DOC uptake rate when the substrates is not limiting, $C_{MIC}$

329    [gCm$^{-3}$] is the microbial biomass pool size, $C_{DOC}$ [gCm$^{-3}$] is the DOC pool size, $C_{O2}$ [m$^3$m$^{-3}$] is

330    the gas concentration of O$_2$ in the soil pore, and $K_{m\_up}$ [gCm$^{-3}$] and $K_{m\_upO2}$ [m$^3$m$^{-3}$] are the

331    corresponding half-saturation constants for DOC and O$_2$, respectively. With the explicit

332    representation of soil moisture control, the two rates become (Zhang et al., 2014)

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

334    $$V_{uptake} = V_{max\_up} C_{MIC} \frac{C_{DOC}}{K_{m\_up} + C_{DOC}} \frac{C_{O2}}{K_{m\_upO2} + C_{O2}} \left( \frac{\theta}{\theta_s} \right)$$    (17)

335    where $\theta$ [-] is the volumetric soil moisture, and $\theta_s$ [-] is the porosity.

336      In addition to using the new rate equations, models 5C and 6C have more carbon pools. In

337    model 5C, DOC is split into two sub-pools for wet zone and dry zone of soil pores, and only the

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

339    becomes

340    $$V_{uptake} = V_{max\_up} C_{MIC} \frac{C_{DOC\_w}}{K_{m\_up} + C_{DOC\_w}} \frac{C_{O2}}{K_{m\_upO2} + C_{O2}} \left( \frac{\theta}{\theta_s} \right).$$    (18)

341    where $C_{DOC\_w}$ [gCm$^{-3}$] is the DOC pool size in the wet soil pores. Model 6C is more complex in

342    that ENZ is further split into two sub-pools for wet and dry pores, and both the wet and dry ENZ

343    are subject to degradation, as shown in Figure 1. The moisture-controlled SOC degradation rate

344    becomes

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

346    for the wet ENZ and

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

348    for the dry ENZ, where $C_{ENZ\_W}$ [gCm$^{-3}$] is the wet soil pores enzyme pool size, $C_{ENZ\_D}$ [gCm$^{-3}$] is

349    the enzyme pool size in the dry soil pores, and $\varepsilon_D$ is the catalysis efficiency of the dry zone enzyme.

350        Due to considering the moisture control and adding more soil pools, model 5C is expected to

351    be significantly better than model 4C for simulating the Birch effect. Since the accumulated ENZ

352    in dry soil is secondary, model 6C is expected to be slightly better than model 5C. In terms of

353    model structural error, model 4C has the largest model structure error, model 5C has significantly

354    less model structure error, and model 6C has the smallest model structural error. As shown below,

355    the degree of model structural error is reflected in the process of Bayesian inference and verified

356    by the cross-validation.

357    **2.4      Observations and parameter estimation**

358        Figure 2 plots the time series of 17,016 observations of soil moister and $CO_2$ efflux used in

359    this study. The observations were obtained during the entire year of 2007, covering a long period

360    of dry season prior to monsoon and episodic rainfall events during monsoon. The first two third of

361    this dataset is used for the Bayesian inference, and the last one third is used for cross-validation.

362    The inference and cross-validation periods have both dry and wet periods, as shown in Figure 2.

363    The observation site is located within the Santa Rita Experimental Range (SRER, 31.8214°N,

Geoscientific
Model Development
Discussions

364   110.8661°W, elevation 1,116 m) outside of Tucson, Arizona (Barron-Gafford et al., 2011; Scott

365   et al., 2009). This savanna site was covered by 22% of perennial grass, forbs and subshrubs and

366   35% of mesquite. The soils are uniformly Comoro loamy sand (77.6% sand, 11.0% clay, and

367   11.4% silt). The half-hourly atmospheric forcing data were collected from measurements through

368   an eddy covariance tower (Scott et al., 2009). This includes downward shortwave, longwave,

369   precipitation, wind, air temperature, humidity, and pressure. Volumetric $CO_2$ concentration was

370   measured at half-hourly interval through compact probes. The $CO_2$ efflux was estimated from the

371   gradient of $CO_2$ concentration measured at two depths of 2 cm and 10 cm through Fick's first law

372   of diffusion, and the estimates were validated against measurements from a portable $CO_2$ gas

373   analyzer.

374       The parameters estimated in this study include the parameters of the soil respiration models

375   (4C – 6C) and the parameters of the data models described in Section 2.1. The estimated

376   parameters of models 4C and 5C include the microbial carbon use efficiency (CUE) [g/g], enzyme

377   production rate, $k_e$ [g/m$^3$s], microbial turnover rate, $\tau_m$ [1/s], and enzyme turnover rate $\tau_e$ [1/s].

378   Uniform distributions are used as the prior in the Bayesian inference, and the ranges of the four

379   parameters are 0.2 – 1.00, $1\times10^{-12}$ – $1\times10^{-7}$, $1\times10^{-12}$ – $1\times10^{-5}$ and $1\times10^{-11}$ – $1\times10^{-6}$, respectively.

380   The values of other parameters are fixed at the values used in Allison et al. (2010). Model 6C has

381   two more parameters, and they are the catalysis efficiency $\varepsilon_D$ [-] and the turnover rate of the dry-

382   zone enzymes $\tau_{en}$ [1/s]. The prior of the two parameters are uniform distributions with the ranges

383   of 0.2 – 0.8 and $1\times10^{-12}$ – $1\times10^{-8}$, respectively.

384       The DREAM-based MCMC simulation is conducted for a total of 24 cases, the combinations

385   of eight data models and three physical models. For each case, the parameter distributions are

386   obtained after drawing a total of $5\times10^5$ samples using five Markov chains. The Gelman and Rubin

387  (1992) R-statistic is used for convergence diagnostic, and it approaches one in less than $4 \times 10^4$

388  samples. The initial 50% of the samples are discarded during the burn-in period.

389  **3    Results of Bayesian Inverse Modeling**

390     This section analyzes the residuals of the best realization (with the highest likelihood value) of

391  the MCMC simulation to understand whether the assumptions of the eight data models hold. The

392  impacts of the data models on the posterior parameter distributions are also analyzed.

393  **3.1    Residual characterization**

394     Figure 3 shows residual plots for model 6C based on data models SLS and WSEP-AC. SLS is

395  the simplest one with the assumptions of homoscedastic, independent, and Gaussian residuals, and

396  the WSEP-AC is the most complex one without the assumptions. Model 6C is the most complex

397  model and also the best one as ranked by Zhang et al. (2014) using Bayesian model selection. The

398  variable $a_t$ plotted in Figures 3a-3c and Figures 3d-3f is defined in equations (2) and (11),

399  respectively. Figures 3a – 3c show that the three residual assumptions are violated when SLS is

400  used because (i) the residual variance is not constant, but increases as a function of the simulated

401  $CO_2$ efflux (Figure 3a); (ii) the autocorrelation function at most lags is beyond the 95% confidence

402  interval (Figure 3b); (iii) and the standard normal density function cannot adequately characterize

403  the residuals (Figure 3c). Figures 3d-f show that, after relaxing the three assumptions, the

404  processed residuals, $a_t$, can be well characterized by WSEP-AC. Figure 3d shows that, after

405  normalizing $\varepsilon_t$ with the linear variance ($\sigma_t = 0.034 + 0.099 E_t$), the variation of the variance of

406  $a_t$ becomes significantly smaller, although the variance is still not a constant. Figure 3e shows that,

407  after removing a first-order autoregressive model from $\varepsilon_t$, $a_t$ becomes less correlated, although the

408  correlation is not fully removed. The two coefficients of the autoregressive model are $\phi_1 = 0.989$

409  and $\phi_2 = 4.5 \times 10^{-6}$; the small value of $\phi_2$ indicates that there is no need to attempt an autoregressive

18

410   model of higher order. Figure 3f shows that $a_t$ follows the SEP distribution with the estimated

411   skewness coefficient of $\xi = 0.933$ and kurtosis coefficient of $\beta = 0.998$. As a summary, Figure

412   3 shows that it is important to examine the residuals and to determine whether a data model is

413   adequate for charactering the residuals. Although WSEP-AC still cannot perfectly characterize $\varepsilon_t$,

414   it is significantly better than SLS.

415        Although the Gaussian assumption used in SLS is violated for model 4C (Figure 3c), this is

416   not generally the case for other data models and physical models. This is shown in Figure 4, which

417   presents the quantile-quantile (Q-Q) plot for the eight data models and the three soil respiration

418   models. For SLS, WLS, SLS-AC, and WLS-AC, the theoretical quantiles are based on the standard

419   normal distribution, $N(0,1)$; for SEP, WSEP, SEP-AC, and WSEP-AC, the theoretical quantiles

420   are based on the standard skew exponential power distribution, $SEP(0,1,1,0)$. If the residuals

421   follow the assumed standard distributions, the Q-Q plots fall on the 1:1 line, which is marked as

422   the theoretical lines in Figure 4. If the residuals are Gaussian or SEP but not standard, the Q-Q

423   plots fall on a straight line but not the 1:1 line. Figures 4a and 4e show that, for all the soil

424   respiration models, the Q-Q plots of SLS and SEP deviate significantly from the theoretical lines

425   and exhibit fat-tail behaviors, which is an indication of outliers (Thyer et al., 2009). The deviation

426   is reduced after accounting for autocorrelation in SLS-AC and SEP-AC, as shown in Figures 4c

427   and 4g (it is interesting to observe from the two figures that the Q-Q plots of the three models are

428   almost visually identical). The deviation is almost fully removed after accounting for

429   heteroscedasticity in WLS and WSEP in that their corresponding Q-Q plots fall on the 1:1 lines,

430   especially for models 5C and 6C, as shown in Figures 4b and 4f. However, the Q-Q plots start

431   deviating from the 1:1 lines as shown in Figures 4d and 4h, after accounting for both

432   heteroscedasticity and autocorrelation in WLS-AC and WSEP-AC. As a summary, Figure 4 shows

433  that, for the numerical example of this study, either the Gaussian or the SEP distribution is valid if

434  heteroscedasticity is accounted for in the data models. However, accounting for autocorrelation in

435  the data models does not help improve the characterization of the residual distribution.

436  **3.2     Posterior parameter distributions**

437  While Figures 3 and 4 help understand validity of the three assumptions used in the data

438  models, the impacts of the data models on estimating model parameter distributions must be

439  evaluated separately. This section discusses the impact of the data model selection on parameter

440  estimation with the objective of understanding if incorrect specification of the data model, will

441  necessarily lead to biased parameter estimates. Such assessment is not a trivial task for three main

442  reasons. First, microbial soil respiration models aggregate complex natural processes and spatial

443  details into simpler conceptual representations. As a results several model parameters are effective

444  values of several complex natural processes that cannot be actually measured in the field as

445  discussed by Vrugt et al. (2013). Second, even for model parameter that can be measured in the

446  field, since the model structure is imperfect, it can be the case that parameter values can be

447  accepted beyond their physically reasonable range as discussed by Pappenberger and Beven

448  (2006). This is often undesirable, if we seek to make the models more mechanistically descriptive.

449  We focus our discussion on carbon use efficiency (CUE) for microbial growth since CUE is a

450  fundamental parameter in microbial soil respiration models, and a reasonable physical range for

451  CUE can estimated.  The concept of microbial CUE(Allison et al., 2010; Bradford et al., 2008;

452  Manzoni et al., 2012; Wieder et al., 2013) has been used to present fundamental microbial

453  processes recent microbial enzyme models(Allison et al., 2010; German et al., 2011; Schimel and

454  Weintraub, 2003; Wang et al., 2013). The microbial CUE, which is marked between MIC and CO2

455  in Figure 1, controls microbial growth, enzyme production and microbial respiration. A reasonable

456 range of CUE can be estimated from the physical viewpoint(Tang and Riley, 2014). Sinsabaugh

457 et al. (2013) study shows that the thermodynamic calculations support a maximum CUE of 0.60

458 and that methods used to estimate CUE in terrestrial systems report a mean value of 0.55.

459 Theoretically, there no lower limit for CUE as it can approach zero, and CUE< 0.1 are reported

460 for terrestrial ecosystems (e.g. Fernández-Martínez et al., 2014) and used in modeling studies (Li

461 et al., 2014).

462    Figure 5 plots the CUE posterior marginal density of the three soil respiration models obtained

463 using the eight data models. The physical range between zero and 0.6 is marked in yellow. Figure

464 5 shows that the CUE posterior parameter distribution for Model 6C for all likelihood functions

465 that does not account for autocorrelation are within a reasonable physical range. For models 4C

466 and 5C, the posterior parameter samples are outside the physical range for six data models. For

467 model 4C, the posterior parameters are within the physical range only for data models SEP and

468 WSEP; for model 5C, the two data models are WLS and WSEP. It is not surprising to find the

469 posterior parameter distribution of models 4C and 5C, which have a certain degree of model

470 structure error, to be out of the plausible physical range. This can be attributed to two reasons.

471 First, the model solution can be biased toward the missing processes in the model structure such

472 as the additional carbon pool in both 4C and 5C or the explicit accounting for soil moister in 4C.

473 Second, biased parameter estimation can compensate for model structure inadequacy and other

474 sources of discrepancy in both the physical model and the statistical model.

475    In addition, it is important to understand how accounting for autocorrelation, heteroscedasty

476 and non-Gaussian residuals can affect the parameter estimation.  First, it is not unexpected to get

477 biased parameter estimates that can be out the reasonable physical range when autocorrelation is

478 explicitly accounted for as shown in Figure 5e-h. This may suggest again that accounting for

479  heteroscedasticity is desirable but accounting for autocorrelation is not. A possible reason is that

480  filtering autocorrelation may reduce the residual space such that the transformed residual space

481  cannot correspond to the parameter space of the models. In other words, parameter information

482  may be lost due to filtering out autocorrelation. However, it is not fully understood why this does

483  not occur for the model 6C under data model SLS-AC, and more research is warranted.  Second,

484  unlike accounting for auto-correlation, accounting only for heteroscedasty (i.e. WLS and WSEP)

485  since this will only amplify or reduce the variance without affecting the structure of the residual

486  space. Figure 5c-d shows that account for heteroscedasty (i.e. WLS and WSEP) tends to improve

487  the parameter estimation in comparison with homoscedastic data models (i.e. SLS and SEP) shown

488  in Figure 5a-b. Finally, with respect to non-Gaussian residuals, Schoups and Vrugt (2010)

489  proposes that the peaked pdf of the SEP with heavier tails compared to Gaussian pdf is useful for

490  making parameter inference robust against outliers. To a certain degree, this can be substantiated

491  by the results in Figure 5a-d, such that SEP and WSEP provide more favorable parameter estimates

492  than SLS and WLS.

493      Finally, from Figure 5 we can also notice that the posterior parameter distribution of SLS

494  (Figure 5a) is very narrow. This narrow posterior parameter distribution of SLS compared to other

495  likelihood functions can be attributed to several reasons. Since SEP can have heavier tails than

496  Gaussian distribution, this can further increase the samples acceptance ratio from tails resulting in

497  wider distribution (Figure 5b). In addition, accounting for heteroscedasticity will wider the

498  posterior parameter distribution (Figure 5c) due to accepting higher variances at peak effluxes.

499  Moreover, filtering correlation (Figure 5e-h) increases the entropy.

Geoscientific
Model Development
Discussions

## 4.    Results of Predictive Performance

Based on the last one third of the $CO_2$ efflux observations, a cross-validation test was conducted

for all the 24 models,  the combinations of three soil respiration models and eight data models.

Given the cross-validation data, the predictive performance is examined using the four statistical

metrics defined in Section 4.1. The metrics are also calculated for the calibration data. This is not

to perform Bayesian model selection given the calibration data, but to better understand the impact

of data models. For each calibration and each cross-validation data, a prediction ensemble is

generated from the two perspectives of parametric uncertainty only and total uncertainty, as

presented in Section 4.2 and 4.3, respectively.

### 4.1    Metrics for evaluating predictive performance

Three criteria are used to evaluate the predictive performance of the soil respiration models

and data models, and they are central mean tendency, dispersion, and reliability. Each criteria is

measured by a single metric. In addition, a newly defined metric is also used for simultaneously

measuring the three criteria. The central mean tendency is measured in this study using the Nash-

Sutcliffe model efficiency (NSME) coefficient (Nash and Sutcliffe, 1970),

$$NSME = 1 - \sum_{i=1}^{n}(d_i - \overline{X_i})^2 \bigg/ \sum_{i=1}^{n}(d_i - \overline{\mathbf{d}})^2 \,, \tag{21}$$

where $n$ is the number of cross-validation data, $d_i$ is the $i$-th data, $\overline{\mathbf{d}}$ is the mean of the data, and

$\overline{X_i}$ is the mean of the prediction ensemble, $X_i$, for $d_i$. NSME ranges from $-\infty$ to 1, with $NSME = 1$

corresponding to a perfect match between data and mean prediction, i.e., the ensemble is centered

on the data. $NSME = 0$ indicates that the model predictions are as only accurate as the mean of the

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

mean prediction.

522    In addition to the central mean tendency, it is also desirable that the ensemble is precise with

523    small dispersion and reliable to cover all the data. This study uses a nonparametric metric for

524    dispersion, and it is the sharpness of a prediction interval (e.g. Smith et al., 2010a)

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

526    where $X_i$ is the prediction ensemble within the 95% prediction interval (the Bayesian credible

527    interval, not the confidence interval used in nonlinear regression (Lu et al., 2013). Smaller values

528    of sharpness indicate better prediction precision. Reliability is measured using predictive coverage.

529    (e.g. Hoeting et al., 1999), which is the percentages of data contained in the prediction interval.

530    Larger predictive coverage values are preferred.

531    To account for the trade-off between the three metrics,(Elshall et al., 2018) defined relative

532    model score (RMS) that simultaneously measure all the three criteria. Scoring rules are commonly

533    used in hydrology to assess predictive performance (e.g. Weijs et al., 2010; Westerberg et al.,

534    2011). RMS is used in this study to measure the relative predictive performance of the

535    combinations of soil respiration models and data models. For combination $M_j$, RMS is defined as

536    $RMS(M_j) = \sum_{i=1}^{n} \dfrac{p(d_i \mid X_{ij}, M_j)}{\sum_{j=1}^{m} p(d_i \mid X_{ij}, M_j)} \times 100$                (23)

537    where $m = 24$ is the number of combinations, and $X_{ij}$ is similar to $X_i$ above and specific to the $j$-th

538    combination. The density function, $p(d_i|X_{ij})$, can be evaluated by first obtaining the density function

539    $p(X_{ij})$ of the ensemble prediction $X_{ij}$ (e.g., by using the kernel density function) and then evaluating

540    $p(d_i|X_{ij})$ using interpolation methods based on the intersection of $X_{ij}$ and $d_i$. This evaluation is based

541    purely on the model predictions, and does not involve any assumptions on the models, their

542    parameters, and likelihood functions. Larger RMS values indicate better overall predictive

543    performance.

Geoscientific
Model Development
Discussions

## 4.2 Predictive performance with parametric uncertainty of soil respiration models

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

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

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

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

567    6h for data models SLS and WSEP-AC. For models 5C and 6C, their NSME values vary with the

568    data models with the central mean accuracy being the worst for SLS-AC which considers only

569    autocorrelation.

570    With respect to parametric uncertainty estimation, Figures 7c and 7d show sharpness generally

571    increases when the three assumptions in the data models are gradually relaxed from SLS to WSEP-

572    AC. This is even more obvious during the validation period. Given that the prediction ensemble

573    does not center on the data, the increasing sharpness is desirable as it improves reliability. This is

574    confirmed by the reliability plots in Figures 7e and 7f. The exceptions are again SLS-AC and SEP-

575    AC that generally have the lowest coverage.

576    With respect to the overall predictive performance, the same variation pattern and exception

577    are also observed in the RMS plots in Figures 7g and 7h. This is not surprising because RMS is

578    the metric that can be used to measure all the three criteria (central mean tendency, sharpness, and

579    reliability). Since the prediction ensemble is not centered on the data, the sharpness and reliability

580    are the decisive factors for evaluating the predictive performance.

581    As a summary, while it is necessary to account for heteroscedasticity in a data model, caution

582    is needed when accounting for autocorrelation in the manner described in Section 2.1. In addition,

583    after comparing the RMS values of the residuals using the Gaussian and SEP distributions. The

584    conclusion is that the SEP distribution outperforms the Gaussian distribution with respect to

585    predictive performance. Finally, uncertainty underestimation as evident by the very small

586    predictive coverage. The underestimation of uncertainty for all the physical models with all

587    likelihood function make sense because only parametric uncertainty is considered.  Considering

588    the overall predictive uncertainty is the subject of the next section.

589     **4.3     Predictive performance with parametric uncertainty of soil respiration models and**

590     **data models**

591        The simulated output $\mathbf{Y}(\theta_p)$ will generally not be equally to the observed output $\mathbf{D}$ and we

592     have a residue error term $\mathbf{e}$ due to measurement, input and model structure errors such that

593     $\mathbf{D} = \mathbf{Y}(\theta_p) + \mathbf{e}$. Accounting for error term $\mathbf{e}$ can be through separating various error terms. For

594     example, in section 4.2 we obtained uncertainty due to the physical model parameters. Accounting

595     for other sources of uncertainty can be done using a single model approach (e.g. Thyer et al., 2009)

596     or a multimodel approach (e.g. Tsai and Elshall, 2013). Alternatively, we can quantify the

597     uncertainty based on total residuals, which include measurement, model input, model structure and

598     parameter estimation errors (e.g. Thyer et al., 2009; Schoups and Vrugt, 2010). This lumped

599     approach is based on sampling the residual error model $\mathbf{e}(\theta_e)$ with parameters $\theta_e$. SLS has one

600     fixed parameter that is the constant variance and other data models have two to six parameters.

601     Thus in Section 4.3, the prediction ensemble addresses parametric uncertainty of not only the soil

602     respiration models but also the data models. When generating the prediction ensemble in the

603     procedure described by Schoups and Vrugt (2010), an ensemble of residuals is first generated by

604     running the data models with posterior samples of the data model parameters for the positive

605     carbon efflux domain; the residual ensemble is then added to the prediction ensemble generated in

606     Section 4.2.

607        We start by the visual assessment of the predictive performance. Figure 8 is similar to Figure

608     6 with the exception that Figure 8 considers the overall all predictive uncertainty (i.e. parametric

609     and output uncertainty), while Figure 6 considers the parametric uncertainty only. Figure 8 reveals

610     a practical observation about accounting for the overall uncertainty through the lumped approach

611     of sampling the residual errors model. Figure 8b shows that desp8ite the wide prediction interval

Geoscientific
Model Development
Discussions

612   of model 4C, which has significant model structure error, it could not capture the birch pulse

613   around day 180. This clearly indicates that proper modeling of the residual error will not make-up

614   for of significant model structure error.

615   Figure 9 plots the four statistics (NSME, sharpness, predictive coverage, and RMS) of the three

616   models under the eight data models to assess the predictive performance. First with respect to

617   central mean tendency, The NSME values in Figures 9a-9b are visually the same as those in

618   Figures 7a-7b, indicating that the central mean accuracy under parametric uncertainty is the same

619   as that under predictive uncertainty.

620   With respect to uncertainty, the values of sharpness and predictive coverage increase

621   substantially (Figures 9c – 9f). In particular, Figures 9e and 9f show that, except for SLS and SEP,

622   the predictive coverage of the rest six data models are close to 100% for all the three models,

623   indicating that the prediction intervals cover almost all the data. This is demonstrated in Figures 6

624   for WSEP-AC. Similar to Figures 7c and 7d, Figures 9c and 9d also show a general pattern that

625   the sharpness increases when the three assumptions in the data models are gradually relaxed from

626   SLS to WSEP-AC. The data models that account for autocorrelation are still the exceptions.

627   With respect to the overall predictive performance, the RMS values are largely determined by

628   mean accuracy and sharpness as the predictive coverage is similar for different data models.

629   Figures 9g and 9h of RMS show that the predictive performance of the four data models that

630   account for autocorrelation is worse than that of the other four data models. This suggests again

631   that one needs to be cautious when building autocorrelation into a data model. This is consistent

632   with the finding of Evin et al. (2013, 2014) that accounting for autocorrelation before accounting

633   for heteroscedasticity or jointly accounting for autocorrelation and heteroscedasticity can result in

634   poor predictive performance. In summary, Figures 9g and 9h show for both the calibration and

635    prediction periods that accounting for heteroscedasticity (i.e. WLS and WSEP) will give the best

636    overall predictive skill, and accounting for autocorrelation without heteroscedasticity (i.e. SLS-

637    AC and SEP-AC) will give the worst overall predictive skill. Finally, for the three soil respiration

638    models, RMS shows that model 4C has the worst predictive performance for both the calibration

639    and cross-validation data. Generally speaking, the high fidelity model 6C outperforms model 5C

640    for both the calibration and cross-validation data, which justifies the complexity of model 6C.

641        To demonstrate the impacts of the data models on predictive performance of the soil respiration

642    models, Figure 10 plots the model simulations and predictions given by model 6C during the

643    calibration and cross-validation periods using all the eight data models.

644        In Figure 10 we try to understand the predictive performance characteristics of the different

645    data models by looking at the predictive performance of model 6C. Specific predictive

646    performance patterns can be identified. Figures 10-a-d show that SLS and SEP have similar

647    predictive performance with SEP generally having better predictive skill especially during the

648    validation period. Accounting for heteroscedasticity using WLS as shown in Figures 10e and 10h

649    will make the predictions more sensitive to peck carbon effluxes and will generally improve the

650    predictive coverage on the expense of sharpness and the central mean tendency. WLS and WSEP

651    have similar predictive performance. However, WSEP maintains slightly better central mean

652    tendency and overall predictive performance than WLS. Accounting for autocorrelation using

653    SLS-AC and SEP-AC as shown in Figures 10i and 10l reduces the information content of the

654    residuals thus resulting in wider uncertainty bands and insensitivity to peak carbon effluxes as

655    compared to SLS and SEP (Figures 10a-d). This resulted in deteriorating the sharpness, the central

656    mean tendency and the capturing of peak carbon fluxes, especially during the validation period.

657    Accounting for both heteroscedasticity and autocorrelation using WLS-AC and WSEP-AC will

658    make the inference robust against peck carbon effluxes, yet due to the loss of information content

659    uncertainty bands are still wider and uncertainty becomes overestimated especially during

660    validation period as compared to WLS and WSEP. The results of Models 4C and 5C, which are

661    not shown here, also show the same prediction patterns with respect to non-Gaussian residuals,

662    heteroscedasticity and autocorrelation.

663        From figure 10 we also notice that data models that have good overall predictive performance

664    as measured by RMS during the calibration period will maintain this good predictive performance

665    during the validation period. For model 6C, RMS values for the calibration and validation periods

666    are very well correlated with a correlation coefficient of 0.92. However, we note that for models

667    4C and 5C the overall predictive performance during the calibration and validation periods are not

668    that well correlated as 6C, with correlation coefficients of 0.52 for model 4C and 0.61 for model

669    5C. This suggests that model 6C is more robust than 4C and 5C for forecasting and hindcasting.

670    **5.    Conclusions**

671        In parameter estimation and prediction of soil carbon fluxes to the atmosphere we often

672    assume that residuals, which include observation, model input, model structure and parameter

673    estimation errors, are normally distributed, homoscedastic and uncorrelated. We studied these

674    assumptions by calibrating three microbial enzyme models, which have varying degrees of model

675    structure errors. We tested eight data modeling starting with the standard least squares (SLS) and

676    skew exponential power (SEP) data models that assume homoscedasictic and non-correlated

677    residuals. Given these two distributions, we evaluated six other data models that account for

678    heteroscedasicty (WLS and WSEP), autocorrelation (SLS-AC and SEP-AC) and joint inversion of

679    heteroscedasicty and autocorrelation (WLS-AC and WSEP-AC). To our knowledge this is the first

680    study that provide such detailed analysis soil reparation inverse modeling. We also used three solid

681   respiration models with different degrees of model realism and model complexity (i.e. number of

682   model parameters), to understand the impact of model discrepancy on the calibration results under

683   different data models. We analyzed the calibration results with respect to (i) residual

684   characterization, (ii) parameter estimation, (iii) predictive performance and (iv) impact of model

685   discrepancy. The main findings of this study can be summarized as follows:

686       (i) With respect to residual characterization, residual analysis results suggest that the common

687   assumption of not accounting for heteroscedasicty and autocorrelation of residuals (i.e. SLS and

688   SEP) results in poor characterization of residuals. Explicit accounting for heteroscedasicty (i.e.

689   WLS and WSEP) can result in good characterization of the residuals, and is followed by joint the

690   inversion of heteroscedasicty and autocorrelation (i.e. WSL-AC and WSEP-AC). Accounting for

691   autocorrelation only (i.e. SLS-AC and SEP-AC) may not improve much the characterization of the

692   residuals.

693       (ii) With respect to parameter estimation, we focused on carbon use efficiency (CUE), which

694   is a central parameter in soil respiration modeling. We found the SLS with relatively reasonable

695   posterior parameter distribution for CUE, yet very narrow posterior. Data models consider

696   autocorrelation (i.e. SLS-AC, SEP-AC, WLS-AC and WSEP-AC) tend to generally yield CUE

697   estimates that are physically non-reasonable. We speculate that filtering correlation can affect the

698   mapping of the model physics (as implicitly included in the residuals) into the likelihood space,

699   which might result in biased parameter estimates that are physically unreasonable.

700       (iii) With respect to predictive performance, we assessed the central mean tendency,

701   uncertainty bands and the overall predictive performance for both the calibration and the cross-

702   validation periods. Results show that accounting for autocorrelation (i.e. SLS-AC, SEP-AC, WLS-

703   AC, and WSEP-AC) deteriorate the predicative performance, such that the predictive performance

704 is inferior to SLS in terms of the central mean tendency and overall predictive skill, especially

705 during the cross-validation period. Results also indicates that using a SEP distribution can

706 potentially improve the predictive performance. The same is true for accounting for

707 heteroscedasticity. Using SEP distribution and accounting for heteroscedasticity (i.e. WSEP) can

708 potentially improve the predictive performance.

709 (iv) With respect to the impact of model discrepancy, the high fidelity complex model (6C)

710 gives the best results with respect to parameter estimation and predictive performance. Model 6C

711 generally maintained its superior performance under different data models. This justifies the

712 complexity of model 6C relative to model 5C that has one less carbon pool. Model 4C that has a

713 low fidelity model with only four carbon pools and lacks the explicit representation of soil moisture

714 control, maintains its poor performance for different data models.

715 From the empirical findings of this research we conclude the following: (i) Not accounting for

716 heteroscedasticity and autocorrelation using a Gaussian or non-Gaussian data model might not

717 necessarily result in biased parameter estimates or biased predictions with respect to central mean

718 tendency, but will definitely underestimate uncertainty resulting in lower overall predictive

719 performance. (ii) Using a non-Gaussian residual error model can improve the parameter estimates,

720 and the predictive performance with respect to central mean tendency and uncertainty estimation.

721 (iii) Accounting for heteroscedasticity will definitely improve the uncertainty estimation with

722 respect to reliability at the cost of having a wider predictive interval. (iv) This study confirms the

723 empirical findings and theoretical analysis of Evin et al. (2013; 2014) that separate accounting for

724 autocorrelation or joint inversion of correlation and heteroscedasticity can be problematic.

725 Relatively poor performance with respect to autocorrelation can be due to our implementation

726 scheme, which can be improved by using the post-processing inference approach for

727    autocorrelation (Evin et al., 2013; 2014) or similar strategies (Li et al., 2015, 2016). Further

728    investigation of this point is warranted in a future study.

729        The conclusions above are subject to several limitations. First, the conclusions are specific to

730    the soil respiration models developed and validated for semi-arid savannah. Performance

731    variations across different soil respiration models with different levels of complexities is possible.

732    Second, the conclusions are conditioned on the data that were obtained at the half-hour interval

733    over a one-year period. Different conclusions are possible if the data are thinned to daily or weekly

734    scales or data of longer observation periods are used. Third, the study investigates effects of the

735    residual assumptions of formal likelihood functions through direct conditioning of the error model

736    parameters, yet this can also be done through other approaches such as residuals transformation

737    (Thiemann et al., 2001), autorgressive bias model (Del Giudice et al., 2013), approximate Bayesian

738    computation (Sadegh and Vrugt, 2013), data assimliation (Spaaks and Bouten, 2013). Comparing

739    different methods for accounting the residual assumptions are beyond the scope of this work.

740    Fourth, this study focuses on formal Bayesian computation using formal likelihood functions, and

741    comparison with other inference functions such as informal likelihood functions or approximate

742    Bayesian computation is warranted in a future study.

743        Based on the aforesaid conclusions and limitations, we recommend to start calibrating soil

744    respiration models with simple SLS or SEP likelihood function. If the residuals characterization is

745    adequate (e.g. Scharnagl et al., 2011), then the underlying assumptions are met. Otherwise,

746    increase complexity of the data model until satisfactory results are obtained in terms of residuals

747    characterization, posterior parameter estimation and predictive performance. Although the

748    empirical findings of this study provide general guidelines for data model selection of microbial

Geoscientific
Model Development
Discussions

749    soil respiration models, more comparative studies are needed to validate and refute the findings of

750    this study.

751    **Code and data availability**

752    The data and codes and models used to produce this paper are available on contact of the

753    corresponding author at mye@fsu.edu. We cannot publicly share the workflow because MT-

754    DREAM$_{(ZS)}$ code (Laloy and Vrugt, 2012) , which is a main component in the workflow, is in the

755    process of becoming a commercial code.

756    **Author contributions**

757    ASE developed and implemented the code for the eight data models for soil respiration modeling,

758    and prepared the manuscript with contribution of all co-authors. MY developed the research idea

759    and outline, and supervised the research implementation. GN developed the soil respiration

760    models. GAB collected and processed the eddy-covariance data used for model calibration.

761    **Competing interests**

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

763    **Acknowledgement**

766    **References**

767    Ahrens, B., Reichstein, M., Borken, W., Muhr, J., Trumbore, S. E. and Wutzler, T.: Bayesian

768        calibration of a soil organic carbon model using ΔC measurements of soil organic carbon

769        and heterotrophic respiration as joint constraints, Biogeosciences, 11(8), 2147–2168,

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

771    Allison, S. D., Wallenstein, M. D. and Bradford, M. A.: Soil-carbon response to warming

772    dependent on microbial physiology, Nat. Geosci., 3, 336 [online] Available from:

773        http://dx.doi.org/10.1038/ngeo846, 2010.

774    Barr, J. G., Engel, V., Fuentes, J. D., Fuller, D. O. and Kwon, H.: Modeling light use efficiency in

775        a subtropical mangrove forest equipped with $CO_2$ eddy covariance, Biogeosciences, 10(3),

776        2145–2158, doi:10.5194/bg-10-2145-2013, 2013.

777    Barron-Gafford, G. A., Scott, R. L., Jenerette, G. D. and Huxman, T. E.: The relative controls of

778        temperature, soil moisture, and plant functional group on soil $CO_2$ efflux at diel,

779        seasonal, and annual scales, J. Geophys. Res. Biogeosciences, 116(1), 1–16,

780        doi:10.1029/2010JG001442, 2011.

781    Bradford, M. A., Davies, C. A., Frey, S. D., Maddox, T. R., Melillo, J. M., Mohan, J. E., Reynolds,

782        J. F., Treseder, K. K. and Wallenstein, M. D.: Thermal adaptation of soil microbial

783        respiration to elevated temperature, Ecol. Lett., 11(12), 1316–1327, doi:10.1111/j.1461-

784        0248.2008.01251.x, 2008.

785    Chevallier, F. and O'Dell, C. W.: Error statistics of Bayesian $CO_2$ flux inversion schemes as seen

786        from GOSAT, Geophys. Res. Lett., 40(6), 1252–1256, doi:10.1002/grl.50228, 2013.

787    Davidson, E. A. and Janssens, I. A.: Temperature sensitivity of soil carbon decomposition and

788        feedbacks to climate change, Nature, 440, 165 [online] Available from:

789        http://dx.doi.org/10.1038/nature04514, 2006.

790    Davidson, E. A., Samanta, S., Caramori, S. S. and Savage, K.: The Dual Arrhenius and Michaelis–

791        Menten kinetics model for decomposition of soil organic matter at hourly to seasonal time

792        scales, Glob. Chang. Biol., 18(1), 371–384, doi:10.1111/j.1365-2486.2011.02546.x, 2011.

793    Elshall, A. S. and Tsai, F. T.-C.: Constructive epistemic modeling of groundwater flow with

794        geological structure and boundary condition uncertainty under the Bayesian paradigm, J.

795     Hydrol., 517, doi:10.1016/j.jhydrol.2014.05.027, 2014.

796     Elshall, A. S., Ye, M., Pei, Y., Zhang, F., Niu, G.-Y. and Barron-Gafford, G. A.: Relative model

797         score: a scoring rule for evaluating ensemble simulations with application to microbial soil

798         respiration modeling, Stoch. Environ. Res. Risk Assess., 32(10), 2809–2819,

799         doi:10.1007/s00477-018-1592-3, 2018.

800     Evin, G., Kavetski, D., Thyer, M. and Kuczera, G.: Pitfalls and improvements in the joint inference

801         of heteroscedasticity and autocorrelation in hydrological model calibration, Water Resour.

802         Res., 49(7), 4518–4524, doi:10.1002/wrcr.20284, 2013.

803     Evin, G., Thyer, M., Kavetski, D., McInerney, D. and Kuczera, G.: Comparison of joint versus

804         postprocessor approaches for hydrological uncertainty estimation accounting for error

805         autocorrelation and heteroscedasticity, Water Resour. Res., 50(3), 2350–2375,

806         doi:10.1002/2013WR014185, 2014.

807     Fernández-Martínez, M., Vicca, S., Janssens, I. A., Sardans, J., Luyssaert, S., Campioli, M.,

808         Chapin III, F. S., Ciais, P., Malhi, Y., Obersteiner, M., Papale, D., Piao, S. L., Reichstein,

809         M., Rodà, F. and Peñuelas, J.: Nutrient availability as the key regulator of global forest

810         carbon balance, Nat. Clim. Chang., 4, 471 [online] Available from:

811         http://dx.doi.org/10.1038/nclimate2177, 2014.

812     Gelman, A. and Rubin, D. B.: Inference from Iterative Simulation Using Multiple Sequences, Stat.

813         Sci., 7(4), 457–472, doi:10.1214/ss/1177011136, 1992.

814     German, D. P., Marcelo, K. R. B., Stone, M. M. and Allison, S. D.: The Michaelis–Menten kinetics

815         of soil extracellular enzymes in response to temperature: a cross-latitudinal study, Glob.

816         Chang. Biol., 18(4), 1468–1479, doi:10.1111/j.1365-2486.2011.02615.x, 2011.

817     Del Giudice, D., Honti, M., Scheidegger, A., Albert, C., Reichert, P. and Rieckermann, J.:

818    Improving uncertainty estimation in urban hydrological modeling by statistically describing

819    bias, Hydrol. Earth Syst. Sci., 17(10), 4209–4225, doi:10.5194/hess-17-4209-2013, 2013.

820  Gragne, A. S., Sharma, A., Mehrotra, R. and Alfredsen, K.: Improving real-time inflow forecasting

821    into hydropower reservoirs through a complementary modelling framework, Hydrol. Earth

822    Syst. Sci., 19(8), 3695–3714, doi:10.5194/hess-19-3695-2015, 2015.

823  Hararuk, O., Xia, J. and Luo, Y.: Evaluation and improvement of a global land model against soil

824    carbon data using a Bayesian Markov chain Monte Carlo method, J. Geophys. Res.

825    Biogeosciences, 119(3), 403–417, doi:10.1002/2013JG002535, 2014.

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

827    models for soil $CO_2$, $CH_4$, and $N_2O$ fluxes calibrated using a Bayesian approach and multi-

828    site data, Ecol. Modell., 222(7), 1283–1292, doi:10.1016/j.ecolmodel.2011.01.013, 2011.

829  Hilton, T. W., Davis, K. J. and Keller, K.: Evaluating terrestrial $CO_2$flux diagnoses and

830    uncertainties from a simple land surface model and its residuals, Biogeosciences, 11(2), 217–

831    235, doi:10.5194/bg-11-217-2014, 2014.

832  Hoeting, J. A., Madigan, D., Raftery, A. E. and Volinsky, C. T.: Bayesian model averaging: a

833    tutorial (with comments by M. Clyde, David Draper and E. I. George, and a rejoinder by the

834    authors, Stat. Sci., 14(4), 382–417, doi:10.1214/ss/1009212519, 1999.

835  Högberg, P. and Read, D. J.: Towards a more plant physiological perspective on soil ecology,

836    Trends Ecol. Evol., 21(10), 548–554, doi:10.1016/j.tree.2006.06.004, 2006.

837  Hublart, P., Ruelland, D., De Cortázar-Atauri, I. G., Gascoin, S., Lhermitte, S. and Ibacache, A.:

838    Reliability of lumped hydrological modeling in a semi-arid mountainous catchment facing

839    water-use changes, Hydrol. Earth Syst. Sci., 20(9), 3691–3717, doi:10.5194/hess-20-3691-

840    2016, 2016.

Geoscientific
Model Development
Discussions

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

842         Hutjes, R. W. A., Ceulemans, R., Schulze, E.-D., Valentini, R. and Dolman, A. J.: Europe's

843         terrestrial biosphere absorbs 7 to 12% of European anthropogenic CO2 emissions., Science,

844         300(5625), 1538–42, doi:10.1126/science.1083592, 2003.

845    Katz, R. W., Craigmile, P. F., Guttorp, P., Haran, M., Sansó, B. and Stein, M. L.: Uncertainty

846         analysis in climate change assessments, Nat. Clim. Chang., 3, 769 [online] Available from:

847         http://dx.doi.org/10.1038/nclimate1980, 2013.

848    Kavetski, D., Franks, S. W. and Kuczera, G.: Confronting Input Uncertainty in Environmental

849         Modelling, Calibration Watershed Model., doi:doi:10.1029/WS006p0049, 2013.

850    Keenan, T. F., Davidson, E., Moffat, A. M., Munger, W. and Richardson, A. D.: Using model-data

851         fusion to interpret past trends, and quantify uncertainties in future projections, of terrestrial

852         ecosystem carbon cycling, Glob. Chang. Biol., 18(8), 2555–2569, doi:10.1111/j.1365-

853         2486.2012.02684.x, 2012.

854    Klemedtsson, L., Jansson, P. E., Gustafsson, D., Karlberg, L., Weslien, P., Von Arnold, K.,

855         Ernfors, M., Langvall, O. and Lindroth, A.: Bayesian calibration method used to elucidate

856         carbon turnover in forest on drained organic soil, Biogeochemistry, 89(1), 61–79,

857         doi:10.1007/s10533-007-9169-0, 2008.

858    Laloy, E. and Vrugt, J. A.: High-dimensional posterior exploration of hydrologic models using

859         multiple-try DREAM(ZS) and high-performance computing, Water Resour. Res., 48(1),

860         doi:10.1029/2011WR010608, 2012.

861    Li, J., Wang, G., Allison, S. D., Mayes, M. A. and Luo, Y.: Soil carbon sensitivity to temperature

862         and carbon use efficiency compared across microbial-ecosystem models of varying

863         complexity,    Biogeochemistry,    119,    67–84    [online]    Available    from:

864      http://www.jstor.org/stable/24716883, 2014.

865   Li, M., Wang, Q. J., Bennett, J. C. and Robertson, D. E.: A strategy to overcome adverse effects

866      of autoregressive updating of streamflow forecasts, Hydrol. Earth Syst. Sci., 19(1), 1–15,

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

868   Li, M., Wang, Q. J., Bennett, J. C. and Robertson, D. E.: Error reduction and representation in

869      stages (ERRIS) in hydrological modelling for ensemble streamflow forecasting, Hydrol.

870      Earth Syst. Sci., 20(9), 3561–3579, doi:10.5194/hess-20-3561-2016, 2016.

871   Lu, D., Ye, M., Meyer, P. D., Curtis, G. P., Shi, X., Niu, X.-F. and Yabusaki, S. B.: Effects of error

872      covariance structure on estimation of model averaging weights and predictive performance,

873      Water Resour. Res., 49(9), 6029–6047, doi:10.1002/wrcr.20441, 2013.

874   Luo, Y., Ogle, K., Tucker, C., Fei, S., Gao, C., LaDeau, S., Clark, J. S. and Schimel, D. S.:

875      Ecological forecasting and data assimilation in a data-rich era, Ecol. Appl., 21(5), 1429–

876      1442, doi:10.1890/09-1275.1, 2011.

877   Luo, Y., Keenan, T. F. and Smith, M.: Predictability of the terrestrial carbon cycle, Glob. Chang.

878      Biol., 21(5), 1737–1751, doi:10.1111/gcb.12766, 2014.

879   Manzoni, S., Taylor, P., Richter, A., Porporato, A. and Ågren, G. I.: Environmental and

880      stoichiometric controls on microbial carbon-use efficiency in soils, New Phytol., 196(1), 79–

881      91, doi:10.1111/j.1469-8137.2012.04225.x, 2012.

882   Nash, J. E. and Sutcliffe, J. V: River flow forecasting through conceptual models part I — A

883      discussion of principles, J. Hydrol., 10(3), 282–290, doi:https://doi.org/10.1016/0022-

884      1694(70)90255-6, 1970.

885   Pappenberger, F. and Beven, K. J.: Ignorance is bliss: Or seven reasons not to use uncertainty

886      analysis, Water Resour. Res., 42(5), doi:10.1029/2005WR004820, 2006.

887   Peters, W., Jacobson, A. R., Sweeney, C., Andrews, A. E., Conway, T. J., Masarie, K., Miller, J.

888        B., Bruhwiler, L. M. P., Pétron, G., Hirsch, A. I., Worthy, D. E. J., van der Werf, G. R.,

889        Randerson, J. T., Wennberg, P. O., Krol, M. C. and Tans, P. P.: An atmospheric perspective

890        on North American carbon dioxide exchange: CarbonTracker., Proc. Natl. Acad. Sci. U. S.

891        A., 104(48), 18925–30, doi:10.1073/pnas.0708986104, 2007.

892   Le Quéré, C., Peters, G. P., Andres, R. J., Andrew, R. M., Boden, T. A., Ciais, P., Friedlingstein,

893        P., Houghton, R. A., Marland, G., Moriarty, R., Sitch, S., Tans, P., Arneth, A., Arvanitis, A.,

894        Bakker, D. C. E., Bopp, L., Canadell, J. G., Chini, L. P., Doney, S. C., Harper, A., Harris, I.,

895        House, J. I., Jain, A. K., Jones, S. D., Kato, E., Keeling, R. F., Klein Goldewijk, K.,

896        Körtzinger, A., Koven, C., Lefèvre, N., Maignan, F., Omar, A., Ono, T., Park, G.-H., Pfeil,

897        B., Poulter, B., Raupach, M. R., Regnier, P., Rödenbeck, C., Saito, S., Schwinger, J.,

898        Segschneider, J., Stocker, B. D., Takahashi, T., Tilbrook, B., van Heuven, S., Viovy, N.,

899        Wanninkhof, R., Wiltshire, A. and Zaehle, S.: Global carbon budget 2013, Earth Syst. Sci.

900        Data, 6(1), 235–263, doi:10.5194/essd-6-235-2014, 2014.

901   Raich, J. W. J. J. W., Potter, C. S. C. and Bhagawati, D.: Interannual variability in global soil

902        respiration, 1980-94, Glob. Chang. Biol., 8, 800–812, doi:10.1046/j.1365-

903        2486.2002.00511.x, 2002.

904   Ren, X., He, H., Moore, D. J. P., Zhang, L., Liu, M., Li, F., Yu, G. and Wang, H.: Uncertainty

905        analysis of modeled carbon and water fluxes in a subtropical coniferous plantation, J.

906        Geophys. Res. Biogeosciences, 118(4), 1674–1688, doi:10.1002/2013JG002402, 2013.

907   Ricciuto, D. M., King, A. W., Dragoni, D. and Post, W. M.: Parameter and prediction uncertainty

908        in an optimized terrestrial carbon cycle model: Effects of constraining variables and data

909        record length, J. Geophys. Res. Biogeosciences, 116(1), 1–17, doi:10.1029/2010JG001400,

910     2011.

911     Richardson, A. D. and Hollinger, D. Y.: Statistical modeling of ecosystem respiration using eddy

912         covariance data: Maximum likelihood parameter estimation, and Monte Carlo simulation of

913         model and parameter uncertainty, applied to three simple models, Agric. For. Meteorol.,

914         131(3–4), 191–208, doi:10.1016/j.agrformet.2005.05.008, 2005.

915     Sadegh, M. and Vrugt, J. A.: Bridging the gap between GLUE and formal statistical approaches:

916         Approximate Bayesian computation, Hydrol. Earth Syst. Sci., 17(12), 4831–4850,

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

918     Scharnagl, B., Vrugt, J. A., Vereecken, H. and Herbst, M.: Inverse modelling of in situ soil water

919         dynamics: Investigating the effect of different prior distributions of the soil hydraulic

920         parameters, Hydrol. Earth Syst. Sci., 15(10), 3043–3059, doi:10.5194/hess-15-3043-2011,

921         2011.

922     Schimel, J. P. and Weintraub, M. N.: The implications of exoenzyme activity on microbial carbon

923         and nitrogen limitation in soil: a theoretical model, Soil Biol. Biochem., 35(4), 549–563,

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

925     Schmidt, M. W. I., Torn, M. S., Abiven, S., Dittmar, T., Guggenberger, G., Janssens, I. A., Kleber,

926         M., Kögel-Knabner, I., Lehmann, J., Manning, D. A. C., Nannipieri, P., Rasse, D. P., Weiner,

927         S. and Trumbore, S. E.: Persistence of soil organic matter as an ecosystem property, Nature,

928         478(7367), 49–56, doi:10.1038/nature10386, 2011.

929     Schoups, G. and Vrugt, J. A.: A formal likelihood function for parameter and predictive inference

930         of hydrologic models with correlated, heteroscedastic, and non-Gaussian errors, Water

931         Resour. Res., 46(10), 1–17, doi:10.1029/2009WR008933, 2010.

932     Scott, R. L., Jenerette, G. D., Potts, D. L. and Huxman, T. E.: Effects of seasonal drought on net

933    carbon dioxide exchange from a woody-plant-encroached semiarid grassland, J. Geophys.

934    Res. Biogeosciences, 114(4), doi:10.1029/2008JG000900, 2009.

935    Shi, X., Ye, M., Curtis, G. P., Miller, G. L., Meyer, P. D., Kohler, M., Yabusaki, S. and Wu, J.:

936    Assessment of parametric uncertainty for groundwater reactive transport modeling, Water

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

938    Sinsabaugh, R. L., Manzoni, S., Moorhead, D. L. and Richter, A.: Carbon use efficiency of

939    microbial communities: stoichiometry, methodology and modelling, Ecol. Lett., 16(7), 930–

940    939, doi:10.1111/ele.12113, 2013.

941    Smith, M. W., Bracken, L. J. and Cox, N. J.: Toward a dynamic representation of hydrological

942    connectivity at the hillslope scale in semiarid areas, Water Resour. Res., 46(12),

943    doi:10.1029/2009WR008496, 2010a.

944    Smith, T., Sharma, A., Marshall, L., Mehrotra, R. and Sisson, S.: Development of a formal

945    likelihood function for improved Bayesian inference of ephemeral catchments, Water

946    Resour. Res., 46(12), 1–11, doi:10.1029/2010WR009514, 2010b.

947    Smith, T., Marshall, L. and Sharma, A.: Modeling residual hydrologic errors with Bayesian

948    inference, J. Hydrol., 528, 29–37, doi:10.1016/j.jhydrol.2015.05.051, 2015.

949    Spaaks, J. H. and Bouten, W.: Resolving structural errors in a spatially distributed hydrologic

950    model using ensemble Kalman filter state updates, Hydrol. Earth Syst. Sci., 17(9), 3455–

951    3472, doi:10.5194/hess-17-3455-2013, 2013.

952    Steinacher, M. and Joos, F.: Transient Earth system responses to cumulative carbon dioxide

953    emissions: Linearities, uncertainties, and probabilities in an observation-constrained model

954    ensemble, Biogeosciences, 13(4), 1071–1103, doi:10.5194/bg-13-1071-2016, 2016.

955    Tang, J. and Riley, W. J.: Weaker soil carbon–climate feedbacks resulting from microbial and

956 abiotic interactions, Nat. Clim. Chang., 5, 56 [online] Available from:

957 http://dx.doi.org/10.1038/nclimate2438, 2014.

958 Tang, J. and Zhuang, Q.: A global sensitivity analysis and Bayesian inference framework for

959 improving the parameter estimation and prediction of a process-based Terrestrial Ecosystem

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

961 Thiemann, M., Trosset, M., Gupta, H. and Sorooshian, S.: Bayesian recursive parameter estimation

962 for hydrologic models, Water Resour. Res., 37(10), 2521–2535,

963 doi:10.1029/2000WR900405, 2001.

964 Thyer, M., Renard, B., Kavetski, D., Kuczera, G., Franks, S. W. and Srikanthan, S.: Critical

965 evaluation of parameter consistency and predictive uncertainty in hydrological modeling: A

966 case study using Bayesian total error analysis, Water Resour. Res., 45(12), 1–22,

967 doi:10.1029/2008WR006825, 2009.

968 Tiedeman, C. R. and Green, C. T.: Effect of correlated observation error on parameters,

969 predictions, and uncertainty, Water Resour. Res., 49(10), 6339–6355,

970 doi:10.1002/wrcr.20499, 2013.

971 Tsai, F. T.-C. and Elshall, A. S.: Hierarchical Bayesian model averaging for hydrostratigraphic

972 modeling: Uncertainty segregation and comparative evaluation, Water Resour. Res., 49(9),

973 doi:10.1002/wrcr.20428, 2013.

974 Tucker, C. L., Young, J. M., Williams, D. G. and Ogle, K.: Process-based isotope partitioning of

975 winter soil respiration in a subalpine ecosystem reveals importance of rhizospheric

976 respiration, Biogeochemistry, 121, 389–408 [online] Available from:

977 http://www.jstor.org/stable/24717586, 2014.

978 Tuomi, M., Vanhala, P., Karhu, K., Fritze, H. and Liski, J.: Heterotrophic soil respiration-

Geoscientific
Model Development
Discussions

979        Comparison of different models describing its temperature dependence, Ecol. Modell.,

980        211(1–2), 182–190, doi:10.1016/j.ecolmodel.2007.09.003, 2008.

981    Vargas, R., Carbone, M. S., Reichstein, M. and Baldocchi, D. D.: Frontiers and challenges in soil

982        respiration research: from measurements to model-data integration, Biogeochemistry,

983        102(1), 1–13, doi:10.1007/s10533-010-9462-1, 2011.

984    Vrugt, J. A. and Ter Braak, C. J. F.: DREAM(D): An adaptive Markov Chain Monte Carlo

985        simulation algorithm to solve discrete, noncontinuous, and combinatorial posterior parameter

986        estimation problems, Hydrol. Earth Syst. Sci., 15(12), 3701–3713, doi:10.5194/hess-15-

987        3701-2011, 2011.

988    Vrugt, J. A., ter Braak, C. J. F., Diks, C. G. H. and Schoups, G.: Hydrologic data assimilation using

989        particle Markov chain Monte Carlo simulation: Theory, concepts and applications, Adv.

990        Water Resour., 51, 457–478, doi:10.1016/j.advwatres.2012.04.002, 2013.

991    Wang, G., Post, W. M. and Mayes, M. A.: Development of microbial-enzyme-mediated

992        decomposition model parameters through steady-state and dynamic analyses, Ecol. Appl.,

993        23(1), 255–272, doi:10.1890/12-0681.1, 2013.

994    Weijs, S. V., Schoups, G. and Van De Giesen, N.: Why hydrological predictions should be

995        evaluated using information theory, Hydrol. Earth Syst. Sci., 14(12), 2545–2558,

996        doi:10.5194/hess-14-2545-2010, 2010.

997    Westerberg, I. K., Guerrero, J. L., Younger, P. M., Beven, K. J., Seibert, J., Halldin, S., Freer, J.

998        E. and Xu, C. Y.: Calibration of hydrological models using flow-duration curves, Hydrol.

999        Earth Syst. Sci., 15(7), 2205–2227, doi:10.5194/hess-15-2205-2011, 2011.

1000   Wieder, W. R., Bonan, G. B. and Allison, S. D.: Global soil carbon projections are improved by

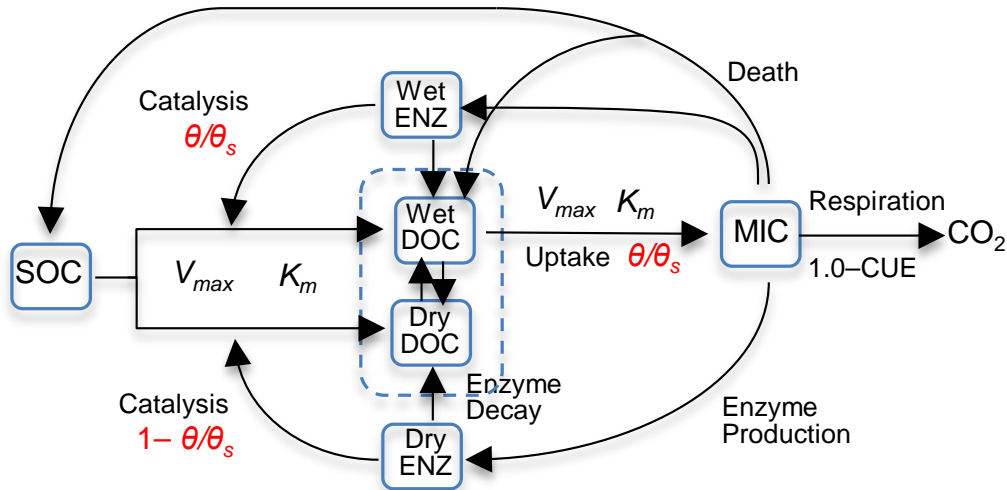1001       modelling microbial processes, Nat. Clim. Chang., 3, 909 [online] Available from:

1002    http://dx.doi.org/10.1038/nclimate1951, 2013.

1003   Wieder, W. R., Allison, S. D., Davidson, E. A., Georgiou, K., Hararuk, O., He, Y., Hopkins, F.,

1004        Luo, Y., Smith, M. J., Sulman, B., Todd-Brown, K., Wang, Y.-P., Xia, J. and Xu, X.:

1005        Explicitly representing soil microbial processes in Earth system models, Global

1006        Biogeochem. Cycles, 29(10), 1782–1800, doi:10.1002/2015GB005188, 2015.

1007   Van Wijk, M. T., Van Putten, B., Hollinger, D. Y. and Richardson, A. D.: Comparison of different

1008        objective functions for parameterization of simple respiration models, J. Geophys. Res.

1009        Biogeosciences, 113(3), 1–11, doi:10.1029/2007JG000643, 2008.

1010   Xu, T., White, L., Hui, D. and Luo, Y.: Probabilistic inversion of a terrestrial ecosystem model:

1011        Analysis of uncertainty in parameter estimation and model prediction, Global Biogeochem.

1012        Cycles, 20(2), 1–15, doi:10.1029/2005GB002468, 2006.

1013   Xu, X., Schimel, J. P., Thornton, P. E., Song, X., Yuan, F. and Goswami, S.: Substrate and

1014        environmental controls on microbial assimilation of soil organic carbon: a framework for

1015        Earth system models, Ecol. Lett., 17(5), 547–555, doi:10.1111/ele.12254, 2014.

1016   Yeluripati, J. B., van Oijen, M., Wattenbach, M., Neftel, A., Ammann, A., Parton, W. J. and Smith,

1017        P.: Bayesian calibration as a tool for initialising the carbon pools of dynamic soil models,

1018        Soil Biol. Biochem., 41(12), 2579–2583, doi:10.1016/j.soilbio.2009.08.021, 2009.

1019   Zhang, X., Niu, G.-Y., Elshall, A. S., Ye, M., Barron-Gafford, G. A. and Pavao-Zuckerman, M.:

1020        Assessing five evolving microbial enzyme models against field measurements from a

1021        semiarid savannah - What are the mechanisms of soil respiration pulses?, Geophys. Res.

1022        Lett., 41(18), doi:10.1002/2014GL061399, 2014.

1023   Zhou, X., Luo, Y., Gao, C., Verburg, P. S. J., Arnone, J. A., Darrouzet-Nardi, A. and Schimel, D.

1024        S.: Concurrent and lagged impacts of an anomalously warm year on autotrophic and

1025        heterotrophic components of soil respiration: A deconvolution analysis, New Phytol., 187(1),

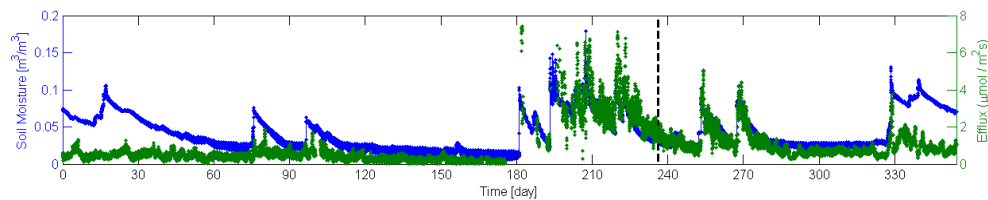1026        184–198, doi:10.1111/j.1469-8137.2010.03256.x, 2010.

1027

Figure 1. Diagram of model 6C representing the processes of (1) degradation of soil organic carbon (SOC) to dissolved organic carbon (DOC) through catalysis of enzymes (ENZ) produced by microbes (MIC), (2) MIC uptake of DOC, and (3) microbial (MIC) respiration to produce $CO_2$ (CUE is the carbon use efficiency). SOC degradation and microbial uptake rates are controlled by water saturation $(\theta / \theta_s)$. The DOC and ENZ pools are split into two subpools, one for the wet zone and the other for the dry zone of the soil pore space. Microbial uptake of DOC occurs only in the wet zone, and the uptake rate is linearly related to $\theta/\theta_s$. Catalysis through ENZ in the wet zone is proportional to $\theta/\theta_s$, while that in the dry zone is proportional to $1 - \theta/\theta_s$. $V_{max}$ (s⁻¹) is the maximum rate, and $K_m$ is the half-saturation concentration.

1040    Figure 2. Time series of soil moisture and efflux observations. The dashed line marks the divide
1041    of the dataset into calibration and validation periods.

1042



1043

Geoscientific
Model Development
Discussions

1044    Figure 3. Residual analysis of the best realization (among multiple MCMC realizations) for model
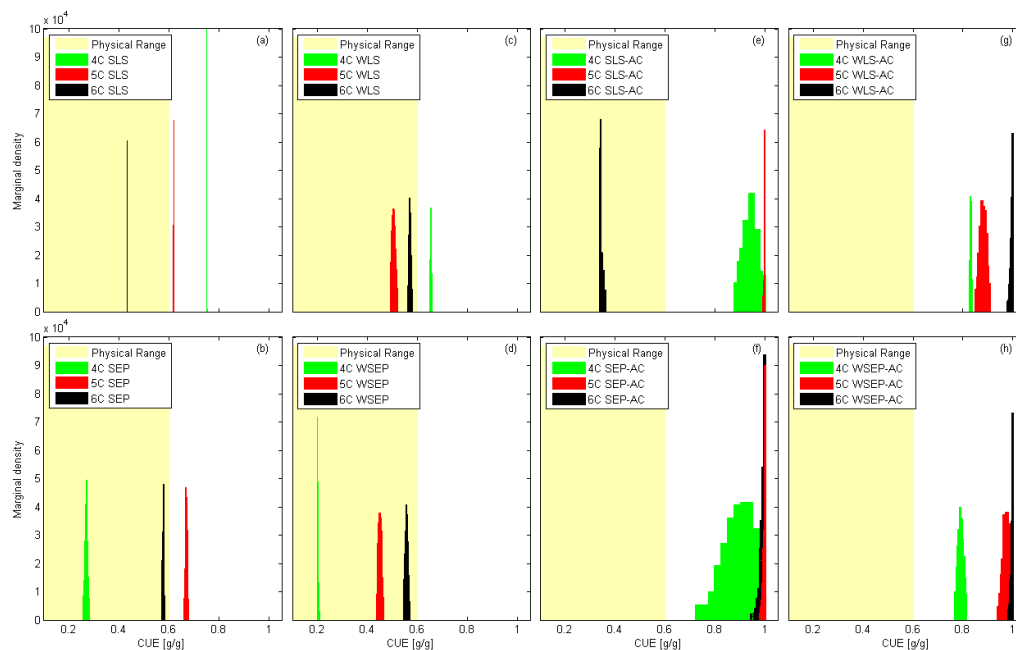1045    6C using data models (a-c) SLS and (d-f) WSEP-AC.

1046



1047

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

1052    Figure 5. Marginal posterior parameter density of carbon use efficiency (CUE) for the three soil
1053    respiration models and eight data models.

1054



1055

1056    Figure 6. Observation data (blue dots) and mean prediction (green line) and 95% credible intervals
1057    (red line) of prediction ensembles for (a)-(f) the calibration period and (g)-(l) the validation period.
1058    The plots are for the three soil respiration models using data models SLS and WSEP-AC. *The*
1059    *prediction ensembles are generated to consider parametric uncertainty of the soil respiration*
1060    *models only.*
1061



1062

Figure 7. (a-b) Nash-Sutcliffe model efficiency (NSME), (c)-(d) sharpness, (e)-(f) predictive coverage, and (g)-(h) relative model score for measuring predictive performance of the three soil respiration models and the eight data models during the calibration and cross-validation periods. *The statistics are evaluated from the prediction ensembles generated to consider parametric uncertainty of the soil respiration models only.*
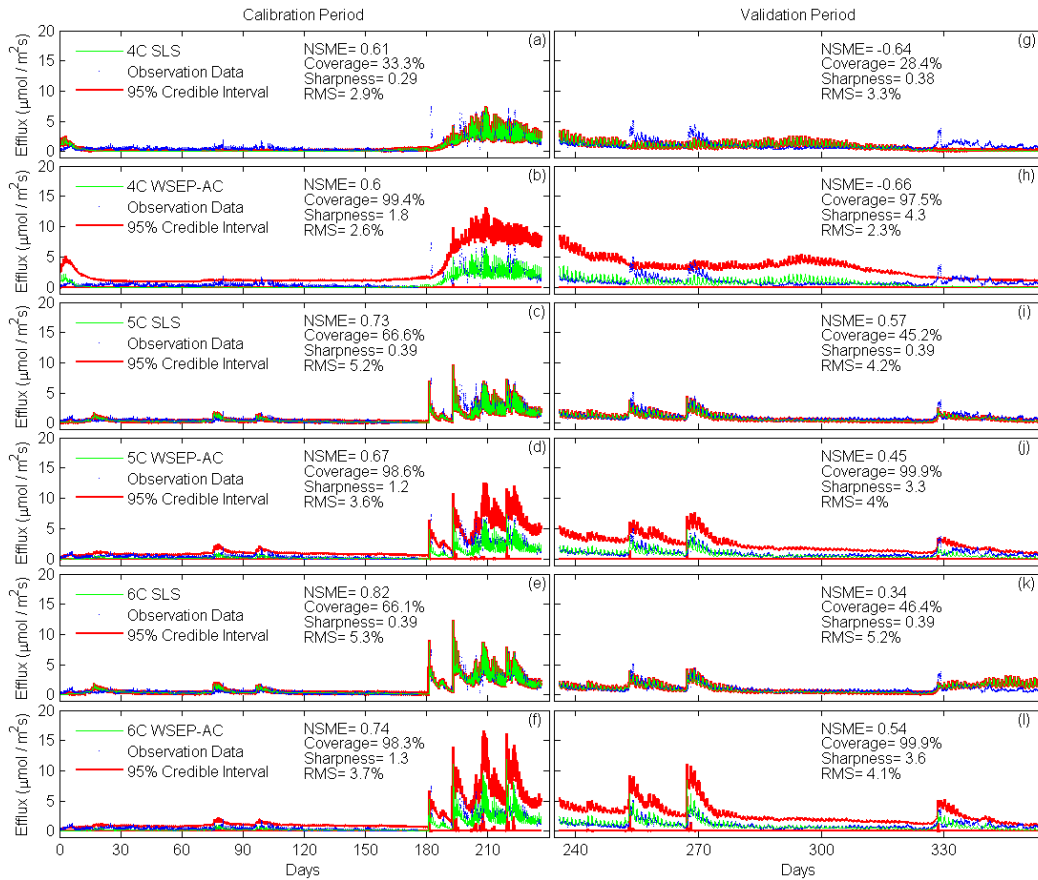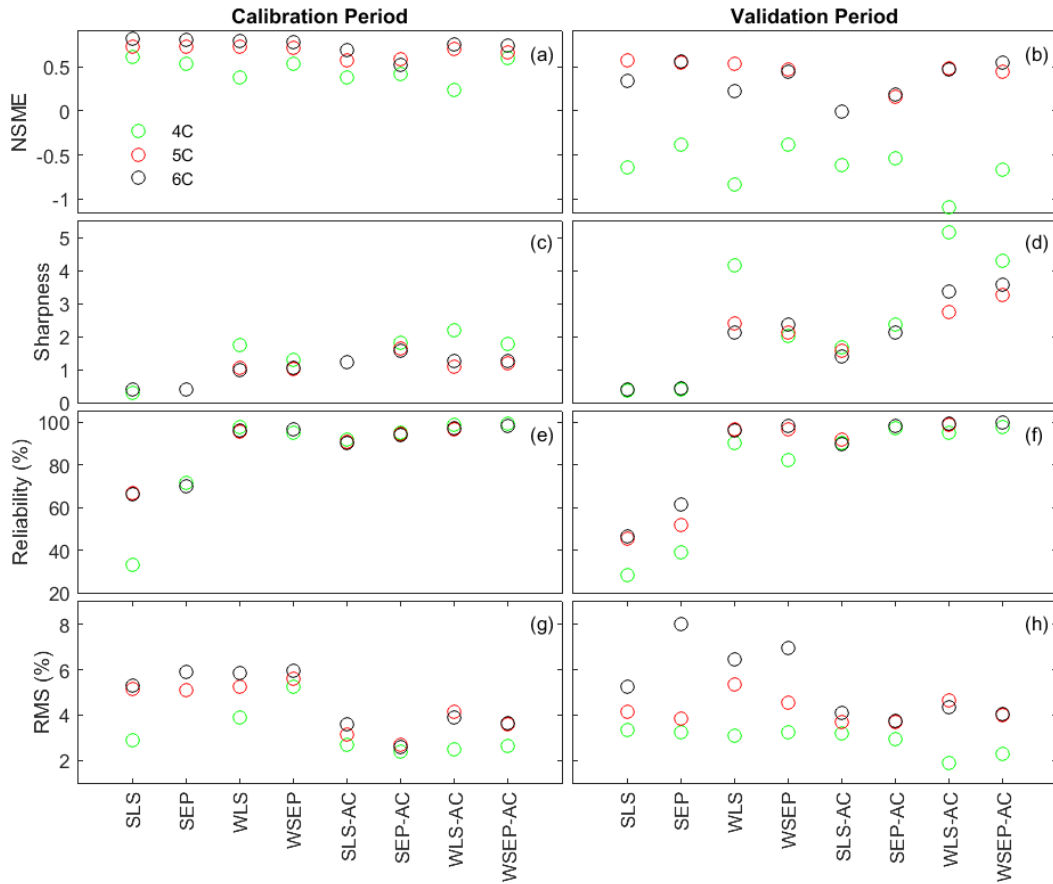
1071　　Figure 8. Observation data (blue dots) and mean prediction (green line) and 95% credible intervals
1072　　(red line) of prediction ensembles for (a)-(f) the calibration period and (g)-(l) the validation period.
1073　　The plots are for the three soil respiration models using data models SLS and WSEP-AC. *The*
1074　　*prediction ensembles are generated to consider parametric uncertainty of not only the soil*
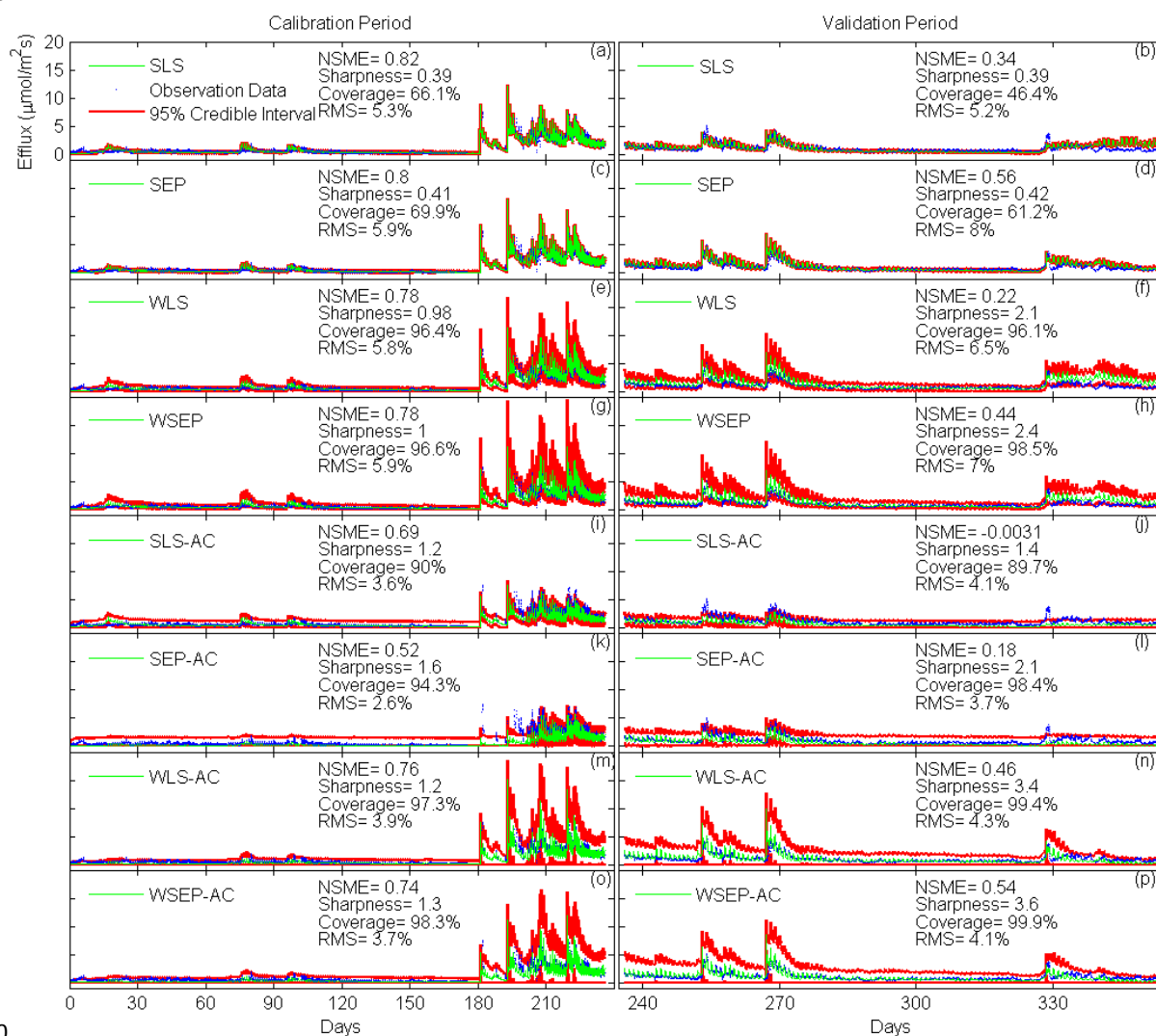1075　　*respiration models but also the data models.*
1076



1077

54

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

1085    Figure 10. Observation data (blue dots) and mean prediction (green line) and 95% credible
1086    intervals (red line) for 6C for the eight likelihood functions during the calibration period (a)-(h)
1087    and the validation period (i)-(p). *The prediction ensembles are generated to consider parametric*
1088    *uncertainty of not only the soil respiration models but also the data models.*
1089



1090