

Response to Referee #1

We are grateful to the reviewer for their insightful feedback and suggestions. The comments were very helpful for improving the manuscript. Updated responses to the reviewer's comments are presented below.

L86-91: lacks model descriptions; and nitrogen-related increases in complexity has not been addressed in the entire manuscript. I suggest you may discuss in the discussion. Also, figure 1 has not been mentioned in the manuscript

L101-107: We added descriptions of CON and AWB model structures in the methods section of the paper and along with a reference to Figure 1.

L491-503: We discussed data assimilation of nitrogen and phosphorous-cycling models. L504-521 address computational speed limitations that hinder Bayesian evaluation of more complex models and potential strategies for facing those challenges.

L166: Log Pseudomarginal Likelihood (LPML) has popped up without prior introduction

L83: We added a sentence introducing LPML. Additional detail regarding LPML was provided between L 198-210.

L185: The difference in curve shape (Fig. 3a, b)

We corrected references to multiple panels of the same figure to the above format.

L189: Is it different between 95% confidence interval and 95% model response ratio credible interval?

L141-149: We added a clarifying paragraph about the difference between the two terms in the Methods section.

L196: a bit confused as well as missing figure annotation. It would be better to choose clear points to address why CON and AWB are showing differences

L243-254: The paragraph was re-structured to compare model outputs to data at specific time points.

L198: rewording to emphasize how the steady state pool size ratio has been changed based on increasing MIC; the unit should be mg C g⁻¹(uppercase); Please check other lines as well

L255: We re-worded the line in question. Thanks to your comments, we realized that exponents in units were incorrectly displayed in the first *Biogeosciences Discuss* submission. The issue seems to be the result of a bug in Microsoft Word for OSX that prevents proper rendering of exponents in Microsoft Word PDF conversions. We seemed to have found a workaround for this bug by printing to a PDF file instead of using Microsoft Word’s included PDF conversion feature.

L199: need to clarify. By the way, what is the function of the trend lines? Have you tried polynomial function? It seems similar patterns between them.

The intent of the posterior predictive mean fit trend lines in Supp. Fig. 1 was to show the moderate, but consistent effect of increasing pre-warming MIC on the AWB model output slopes. In our revision, we emphasized that the influence of pre-warming MIC on posterior predictive slope magnitude is limited to the AWB model (L255 - 260).

By polynomial function, I assume you mean the polynomial analytic function

$$f(x) = \sum_{n=0}^{\infty} c_n(x)^n = c_0 + c_1x + c_2x^2 + c_3x^3 + \dots$$

We could exactly fit the data using a polynomial analytic function with sufficiently many $c_n x^n$ terms, but the goal of the manuscript was not to find the best-fitting arbitrary model. Instead, we sought to demonstrate the feasibility of rigorous Bayesian model comparisons for more complex dynamical ODE systems with the hope that further conclusions regarding model mechanisms and structure could be made following future model comparison results. Fitting analytic polynomial model would not be able to contribute feedback towards the refinement or rejection of elements of dynamical model structure or formulation. Additionally, the polynomial model would be prone to overfitting and penalization by the LPML, LOO, and WAIC algorithms.

L203-206: Is it possible to replace the supplemental figure 3 to represent SOC loss rather than “SOC fraction remaining”? It is difficult to interpret.

“Fraction remaining” and “fraction change” were confusing phrases. However, we wanted the y-axis of Supplemental Fig. 3 to directly correspond with the y-axis of Supplemental Fig. 10 to show that the change in model SOC stocks at 12.5 years was in a biologically realistic range observed in soil warming experiments across various soil types, so we ultimately did not want to change Supplemental Fig. 3 to display SOC loss. To address the confusion, we drew from our use of flux response ratios and re-named the y-axes of both figures to display “SOC response ratios,” which we feel is a more straightforward term. Corresponding edits were made in the Results and Discussion sections to re-phrase descriptions of SOC losses and changes in SOC stocks in terms of response ratios rather than fractions.

L205: decreased from 16.3 to 11.3 %. Please check other lines as well

We fixed occurrences of redundant percent symbols in the revision.

L218: R2; Annotation for varying SOC: “SOC = 50 -> SOC50”

With a workaround solution found for the Microsoft Word rendering issue, the R^2 should now be displayed correctly. We used your excellent suggestion to standardize the notation for pre-warming steady state densities.

L231: not sure this manuscript compared models through AIC and DIC with WAIC and LOO

Yes, AIC and DIC were not computed and compared for reasons of redundancy, stability, and algorithmic limitations. AIC relies on a maximum likelihood estimate which cannot be calculated from non-uniform priors in a Bayesian setting [4]. AIC is more readily deployed in frequentist model comparisons in which normally distributed prior information is not used. DIC can be computed under Bayesian settings, but is an approximation of WAIC that is calculated from posterior means, whereas WAIC involves integration over the posterior distribution sample [1]. Since WAIC can already be readily calculated, has been demonstrated to be more stable and accurate than DIC in much statistics literature, and is itself an approximation of LOO [4], we felt that calculating and plotting WAIC, LPML, and LOO was sufficient for illustrating the influence of the pre-warming steady state ratios on goodness-of-fit. We initially felt that some of this information was too technical and distracting to discuss in a paper not being submitted to a statistics journal, but realized that we needed to provide more expository background and included more justification for our use of WAIC, LPML, and LOO between L69 - 98.

L255-259: which Figures showing this? Also, the sentences are too complex. Please re-write simpler sentences

The paragraph was moved to L393-399 and re-written.

L273: is 50 mg SOC g⁻¹ soil same in line 218 (SOC =50)? Please use consistent unit

L354: We corrected the notation in this line.

L275: Supplemental Table 3?

L357: The typo was corrected. Now Supplemental Table 5.

L316: Supplemental Fig. 5

We removed erroneous use of “Supplementary” from the manuscript.

L322: Supplemental Fig. 8a, b

We corrected our multiple figure panel references to reflect the above format.

Typing error

1. Please check upper case expression; r-squared is R² (uppercase)
2. L156: mg C g⁻¹ soil

3. Put period after abbreviation of figure, e.g., Fig. xx
4. Please double-check figures and table numbers
5. Supplemental or supplementary?

We addressed the typing errors listed above.

Response to Referee #2

We are grateful to the reviewer for the thoughtful and constructive feedback. We address the specific reviewer comments below.

Why was it that CON generally performed better at most SOC densities? What are the estimated parameter ranges compared to other models/literature values? For example the activation energy for SOM decomposition seems to be higher than the activation energy for other processes like uptake which may have some interesting implications.

We feel that the chief reason that CON performed relatively better than AWB at most SOC densities by information criteria and cross-validation was because of overfitting and parameter count. By R^2 , AWB could fit the data slightly better by absolute residual sums than CON at most pre-warming SOC (Fig. 5), but AWB has more parameters and was penalized for that, and we were reluctant to over-reach on conclusions. However, you made an excellent point about the need for more discussion about the fitted parameter posteriors and their biological implications. Consequently, we included more discussion about interpretation of parameter values in our revised manuscript (L400 - 432).

L62-66: Citation for this discussion of R^2 ? Also there are other metrics for evaluating Bayesian models that are not discussed here you dont need an exhaustive review but ROC/AUC and BIC seem common.

For the discussion of R^2 , we now cite Kvålseth 1985 [5], Spiess and Neumeier 2010 [8], and Gelman et al. 2019 [2].

We initially did not discuss BIC due to its similarity to AIC. BIC is closely related to AIC and its computation is dependent on the pointwise maximum likelihood estimate from frequentist methodology [3]. Hence, BIC, contrary to its name, is not a fully Bayesian metric calculated from the posterior distribution. However, as BIC is indeed used to compare out-of-sample predictive accuracy of groups of models, we agree that it should be mentioned and will revise our manuscript to include BIC in the introduction.

We would like to avoid discussing ROC/AUC because they are indicators of the prediction accuracy of binary classifier models trained and tested on categorical data [7]. In our case, we were specifically looking at metrics that estimated out-of-sample prediction accuracy or goodness-of-fit for models conditional on ordinal data with elements in \mathbb{R} . LOO/WAIC and ROC/AUC correspond to fundamentally different model and data types, so we feel that ROC/AUC would be off-topic.

L125: 0.9995 and 0.001 seems like extreme adaptation and step sizes to me, causing the model to take many small steps. If this is a supported strategy, can you provide a citation or justify further?

To address this question, we start by differentiating between the traditional MCMC and the Hamiltonian Monte Carlo algorithm that Stan uses. The HMC is not a random walk algorithm, and each of its trajectories are deterministically calculated via Hamiltonian dynamics. An MCMC step size parameter is fixed through the duration of the sampling. In contrast, the step size is tuned at each step based on the adapt delta and calculated Hamiltonian trajectory. We now clarify that the step size used was the initial step size (L156). Consequently, starting with a small step size does not mean that the HMC algorithm takes fixed steps of the same size for the rest of the chain (refer to this page from the Stan documentation for more detail).

The adapt delta and initial step sizes were set as such in an effort to reduce the number of divergent transitions during the HMC sampling, which was an issue for the AWB model. As can be seen in the supplement, divergent transitions were still detected for AWB following the implementation of the strategy. This indicates the future need to re-parameterize and re-formulate AWB to obtain smoother and more stable parameter space geometries to be explored. However, the divergences occurred at a reduced rate compared to using Stan's default adapt delta of 0.8, so the strategy proved helpful for obtaining a suitable number of posterior samples.

We did not encounter any divergent transitions with CON, but mirrored the HMC parameters for both models since the HMC parameters ultimately do not alter the overall exploration of the parameter space. Increasing adapt delta results in less sensitive tuning of the step size per iteration. In our case, this does mean our step size will be smaller on average as our initial step size will be less responsive to tuning and slower to increase, but the algorithm will be better at navigating geometrically trickier parameter space to generate fewer trajectory divergences. This trade-off comes at the cost of more computationally expensive and less efficient calculated Hamiltonian trajectories. The smaller steps correspond to a drive for more changes in trajectory direction to cover the same amount of ground. With sufficient chain length and samples (which our Bayesian diagnostics indicates that we obtained), the parameter space should still have been adequately explored [6].

The strategy of using higher adapt delta and lower initial step sizes is less documented in formal literature, but has been previously used by other Stan users to mitigate divergent transitions (see [this link](#) and [this link](#) for further discussion). It was important for us to maximize the amount of samples we obtained for AWB, so the increased computational time per iteration was a worthwhile tradeoff for us.

For our revision, we did re-run our HMC simulations with less extreme HMC parameters, using an adaptation delta of 0.95, an initial step size of 0.1, and maximum tree depth of 12 (L156-157), as we found that doing so moderately increased the number of divergences and reduced effective sample ratios, but was worth it for necessary algorithm speed improvements.

L191: As written, implies that AWB performs better because it has a higher RR in subsequent years after the first year, but the data show that the first year has the highest RR, so CON seems to correspond more closely to this. In the discussion, you can bring up the potential realism of oscillations given the Harvard Forest long term warming experiment.

The revised manuscript now discusses the above points in an expanded description of

qualitative model behavior observations in Section 4.1.

L325-333: This discussion of R^2 and other cost metrics seems repetitive to the introduction.

This paragraph was indeed redundant and has been pruned from the revision.

It seems like the performance metrics would be yet better with a lower SOC density (50 mg SOC/g soil), if it were possible to achieve them without the AWB instability. I think you could fix the instability by changing your decomposition/uptake kinetics. Right now in the uptake equation DOC is in the denominator but its initial concentration is much smaller than MIC. So you can either flip to Reverse M-M for uptake or use ECA where both quantities (SOC and ENZ or MIC and DOC) are in the denominator <- this may be harder to fit because it will be more constrained, but it is also harder to break.

We greatly appreciate your perceptive and insightful modeling suggestions that directed us to some very useful and applicable literature. We added a paragraph in the discussion to describe the importance of exploring the changes to model enzyme kinetics you mentioned on model stability and posterior sampling in subsequent work (L461 - 472). We will apply those AWB model modifications to an in-progress follow-up model comparison project that fits models to a different, larger data set from Harvard Forest. We feel that these reparameterization suggestions, particularly the use of ECA kinetics, could help reduce or eliminate the number of divergent transitions generated during AWB HMC posterior sampling.

References

- [1] BETANCOURT, M. A Unified Treatment of Predictive Model Comparison, Jun 2015. arXiv:1506.02273.
- [2] GELMAN, A., GOODRICH, B., GABRY, J., AND VEHTARI, A. R-squared for Bayesian Regression Models. *The American Statistician* 73, 3 (Jul 2019), 307–309.
- [3] GELMAN, A., HWANG, J., AND VEHTARI, A. Understanding predictive information criteria for Bayesian models. *Statistics and Computing* 24 (Jul 2013).
- [4] GELMAN, A., HWANG, J., AND VEHTARI, A. Understanding predictive information criteria for Bayesian models. *Statistics and Computing* 24, 6 (2014), 997–1016.
- [5] KVÅLSETH, T. O. Cautionary Note about R2. *The American Statistician* 39, 4 (Nov 1985), 279–285.
- [6] LIVINGSTONE, S., BETANCOURT, M., BYRNE, S., AND GIROLAMI, M. On the Geometric Ergodicity of Hamiltonian Monte Carlo, Jan 2016. arXiv:1601.08057.
- [7] SAITO, T., AND REHMSMEIER, M. The precision-recall plot is more informative than the ROC plot when evaluating binary classifiers on imbalanced datasets. *PloS one* 10, 3 (Mar 2015), e0118432–e0118432.
- [8] SPIESS, A. N., AND NEUMEYER, N. An evaluation of R2 as an inadequate measure for nonlinear models in pharmacological and biochemical research: A Monte Carlo approach. *BMC Pharmacology* 10, 1 (Dec 2010), 6.

A Bayesian Approach to Evaluation of Soil Biogeochemical Models

Hua W. Xie¹, Adriana L. Romero-Olivares², Michele Guindani³, and Steven D. Allison⁴

¹Center for Complex Biological Systems, University of California, Irvine, 2620 Biological Sciences III Irvine, California 92697, United States of America

²Department of Natural Resources & the Environment, University of New Hampshire, 114 James Hall, Durham, New Hampshire 03824, United States of America

³Department of Statistics, University of California, Irvine, 2241 Donald Bren Hall, Irvine, California 92697, United States of America

⁴Department of Ecology and Evolutionary Biology, Department of Earth System Science, 321 Steinhaus Hall, University of California, Irvine, California 92697, United States of America

Correspondence to: Hua W. Xie (xieh@uci.edu)

Abstract. To make predictions about the carbon cycling consequences of rising global surface temperatures, Earth system scientists rely on mathematical soil biogeochemical models (SBMs). However, it is not clear which models have better predictive accuracy, and a rigorous quantitative approach for comparing and validating the predictions has yet to be established. In this study, we present a Bayesian approach to SBM comparison that can be incorporated into a statistical model selection framework. We compared the fits of linear and non-linear SBMs to soil respiration data compiled in a recent meta-analysis of soil warming field experiments. Fit quality was quantified using Bayesian goodness-of-fit metrics, including the Widely Applicable information criterion (WAIC) and Leave-one-out cross-validation (LOO). We found that the linear model generally out-performed the non-linear model at fitting the meta-analysis data set. Both WAIC and LOO computed higher overfitting risk and effective numbers of parameters for the non-linear model compared to the linear model, conditional on the data set. Goodness-of-fit for both models generally improved when they were initialized with lower and more realistic steady state soil organic carbon densities. Still, testing whether linear models offer definitively superior predictive performance over non-linear models on a global scale will require comparisons with additional site-specific data sets of suitable size and dimensionality. Such comparisons can build upon the approach defined in this study to make more rigorous statistical determinations about model accuracy while leveraging emerging data sets, such as those from long-term ecological research experiments.

1 Introduction

Coupled Earth system models (ESMs) and constituent soil biogeochemical models (SBMs) are used to simulate global soil organic carbon (SOC) dynamics and storage. As global climate changes, some ESM and SBM simulations suggest that substantial SOC losses could occur, resulting in greater soil CO₂ emissions (Crowther et al., 2016). However, there is vast divergence between model predictions. For instance, one ESM predicts a global SOC loss of 72 Pg C over the 21st century, while another predicts a gain of 253 Pg C (Todd-Brown et al., 2014).

Soil biogeochemical models vary greatly in structure (Manzoni and Porporato, 2009), but can be broadly partitioned into two categories: those that implicitly represent soil C dynamics as first-order linear decay processes and those that explicitly represent microbial control over C dynamics with non-linear Michaelis-Menten functions (Wieder et al., 2015a). Explicit models typically include more parameters than linear models because multiple microbial parameters are needed for each decay process as opposed to a single rate parameter. The additional parameters allow explicit models to represent microbial mechanisms, but at the expense of greater model complexity.

Rigorous statistical approaches should be applied to investigate how explicit representation of microbial processes affects predictive model performance. ESM and SBM comparisons involving empirical soil C data assimilations have been conducted previously (Allison et al., 2010; Li et al., 2014) but few standardized statistical methods for ESM and SBM benchmarking and comparison have been developed that would allow for rigorous model selection. Prior model comparisons have involved graphical qualitative comparisons or use of basic fit metrics such as the coefficient of determination, R², to judge fit quality. However, these simple approaches are

Deleted: effect

Deleted: we

Formatted: Indent: First line: 0"

Deleted: ¶

Deleted: a

Deleted: SBM

Deleted: CO₂ flux

Deleted: two

Deleted: a

Deleted: penalty

Deleted: than

Deleted: Fits

Deleted: ¶

Testing

Field Code Changed

Deleted: (Wieder et al., 2015).

Deleted: (Allison et al., 2010; Li et al., 2014)

insufficient for comparing an increasing number of complex models (Jiang et al., 2015; Luo et al., 2016; Wieder et al., 2015a).

R^2 on its own provides limited information about goodness-of-fit. In unmodified form, it quantifies the extent to which the variation of just one chosen model outcome—for instance the mean outcome for a range of parameter values—corresponds to the variation in the data set (Gelman et al., 2019). R^2 does not capture model complexity, overfitting, or parameter uncertainty, which is a reason why R^2 by itself is not sufficient for model evaluation (Kvålseth, 1985). Without accounting for model complexity and parameter count, focusing on optimizing fit by R^2 values alone can easily lead to overfitting (Spiess and Neumeier, 2010).

Encouragingly, a rich toolset to further inform quantitative model evaluation and comparison can be drawn from Bayesian statistics (Hararuk et al., 2014, 2018; Hararuk and Luo, 2014). These tools include information criteria and approximate cross-validation, goodness-of-fit metrics designed for the simultaneous comparison of multiple structurally diverse models. Like R^2 , information criteria and cross-validation are quantitative measures that estimate the fit quality of a model to a given data set. Differing from R^2 , information criteria and cross-validation are relative rather than absolute measures. These metrics evaluate the extent to which the data set supports particular distributions of parameter values and in turn, the uncertainty of parameter estimates. Consequently, if the distribution of Model A outcomes aligns more closely to the data set than the distribution of Model B outcomes, we regard Model A as being more likely to explain the data compared to Model B. Information criteria and cross-validation metrics also typically include terms penalizing for model complexity and overfitting as part of their computation (Gelman et al., 2014). Hence, information criteria and approximate cross-validation are useful tools for model evaluation because they present a comprehensive summary of model fit to time series data and can estimate model predictive accuracy for unmeasured and out-of-sample data points.

Examples of information criteria popularized by widely used R packages such as lme4 and rjags include the Akaike information criterion (AIC), Bayesian information criterion (BIC), and deviance information criterion (DIC) (Vehtari and Ojanen, 2012). However, these metrics have some limitations. AIC, BIC, and DIC do not use full sampled posterior distributions in their computational processes. AIC and BIC both rely on a pointwise maximum likelihood estimate that cannot be derived from non-uniform Bayesian prior distributions, including normal distributions. AIC and BIC (despite BIC's name) thereby have limited use in Bayesian statistics settings. DIC can accommodate non-uniform priors but is calculated from pointwise simplified posterior means. The compression of full posteriors into pointwise means can prompt DIC to compute an impossible negative effective model parameter count in select situations (Gelman et al., 2014). Consequently, the original forms of AIC, BIC, and DIC are no longer recommended for use in Bayesian model assessment by some statisticians in light of superseding alternatives (Gelman et al., 2014).

Three predictive goodness-of-fit metrics address the limitations and stability issues of AIC, BIC, and DIC by incorporating full, non-uniform posterior distributions in their calculations to better account for overfitting and model size (Christensen et al., 2010; Gelman et al., 2014). These metrics include the Widely Applicable information criterion (WAIC), log pseudomarginal likelihood (LPML), and Pareto-smoothed important sampling leave-one-out cross-validation (PSIS-LOO and hereby referred to as LOO). WAIC, LPML, and LOO can estimate the ability of models to fit unobserved measurements outside of the set of measured data samples (Vehtari et al., 2017). Thus, WAIC, LPML, and LOO can be considered as superior barometers for model predictive accuracy compared to AIC, BIC, and DIC.

The overarching goal of this study was to develop a statistically rigorous and mathematically consistent data assimilation framework for SBM comparison that uses predictive Bayesian goodness-of-fit metrics. We pursued three specific objectives as part of that goal. First, we compared the behaviors of two different SBMs, a linear microbial-implicit model termed the conventional model (CON) and a non-linear microbial-explicit model called the Allison-Wallenstein-Bradford model (AWB) (Fig. 1), following data assimilation with soil respiration data sourced from a meta-analysis of soil warming studies (Romero-Olivares et al., 2017). Second, we characterized the parameter spaces of these models using prior probability distributions of parameter values informed by previous studies and expert judgment. Third, we compared specific Bayesian predictive information criteria, in WAIC, LPML, and LOO, to the coefficient of determination, R^2 , for quantifying goodness-of-fit to data. AIC, BIC, and DIC were not analyzed due to their stability limitations, our usage of non-uniform prior distributions, and redundancy with WAIC.

2 Methods

2.1 Model Structures

Field Code Changed

Deleted: 2015

Deleted: Encouragingly, a rich toolset for quantitative model evaluation and comparison can be drawn from Bayesian statistics. These tools include information criteria and

Formatted: Font color: Auto

Deleted: data

Moved (insertion) [1]

Deleted: In contrast, R^2 provides less information about goodness-of-fit. It quantifies the extent to which the variation of just one model outcome, perhaps the mean outcome for a range of parameter values, corresponds to the variation in the data set. R^2 does not capture model complexity, overfitting, or parameter uncertainty, which is a reason why R^2 by itself is not sufficient for model evaluation. Without accounting for model complexity and parameter count, focusing on optimizing fit by R^2 values alone can easily lead to overfitting. ¶

Well-known examples of information criteria include the Akaike information criterion (AIC) and Deviance information criterion (DIC)

Formatted: Font color: Auto

Formatted: Font color: Text 1

Formatted: Font color: Text 1

Deleted: However, these two metrics have some limitations. Neither AIC nor DIC use full sampled posterior distributions in their computations. Additionally, the original formulations of AIC and DIC are more limited and less stable in their ability to account for overfitting and parameter count

Moved up [1]: (Gelman et al., 2014).

Formatted: Font color: Auto

Moved (insertion) [2]

Formatted: Font color: Text 1

Formatted: Font color: Text 1

Deleted: Two more recently developed metrics, the Widely Applicable information criterion (WAIC) and Leave-one-out cross-validation (LOO), address the stability and parameter count issues and improve upon AIC and DIC by using the full posterior distribution (Gelman et al., 2014; Vehtari et al., 2017). WAIC and LOO also estimate the relative potentials of models for fitting measurements not included within the existing observed data set. Thus, WAIC and LOO can be used as barometers for model predictive accuracy. ¶

Formatted: Font color: Text 1

Deleted: models, one

Deleted: and one

Deleted: ,

Deleted: .

Deleted: , including

154 We compared two SBMs, the CON and AWB models (Allison et al., 2010). The models were selected for
 156 this study due to their relative equation simplicity, their tractable parameter count, and limited biological data input
 158 requirements (Supplemental Index 1). The CON system models three separate C pools as state variables including
 160 SOC, dissolved organic C (DOC), and microbial biomass C (MIC) pools, while AWB includes SOC, DOC, MIC,
 and extracellular enzyme biomass C (ENZ) pools (Fig. 1). Additionally, these models were chosen because they are
 C-only models without nitrogen (N) pools. The increased complexity of N-accounting SBMs will require future
 studies with coupled N data sets (Manzoni and Porporato, 2009).

2.2 Meta-analysis Data

162 The data set for model fitting was compiled from a recent meta-analysis of 27 soil warming studies that
 164 measured CO₂ fluxes (Romero-Olivares et al., 2017). The experiments reported between 1 and 13 years of CO₂ flux
 166 measurements following warming perturbation. The elements of this data set consisted of empirical response ratios
 168 calculated by dividing CO₂ fluxes measured in the warming treatments by time-paired CO₂ fluxes measured in the
 170 control treatments. We calculated an annual mean response ratio for each experiment (if data were available for that
 172 year) after warming treatment began. Using these annual means, we calculated one overall mean response ratio for
 174 each year along with pooled variances and standard deviations. Pooled data points were assumed to be “collected” at
 176 the halfway point of each year. Because the experiments had variable lengths, the sample size for the pooled annual
 178 mean declines with increasing time since warming perturbation. The warming perturbation was 3°C on average
 across all the studies, and this average was used as the magnitude of warming in the model simulations.

Model-outputted response ratios were calculated by dividing simulated CO₂ flux following warming
 perturbation by the CO₂ flux at pre-warming steady state. We fit models to flux response ratios rather than raw flux
 measurements for several reasons (Wieder et al., 2015b). First, we eliminate the need to convert flux measurements
 from different experiments into a common unit. Second, response ratios represent a standardized metric for warming
 response across disparate ecosystem types with varying climate, soil, and vegetation properties. Finally, fitting a
 mean response ratio overcomes data gaps present in individual experiments.

2.3 Hamiltonian Monte Carlo Fitting of Differential Equation Models

CON and AWB ordinary differential equation systems were simulated using the CVODE backward
 differentiation method (Curtiss and Hirschfelder, 1952) from the SUNDIALS library of equation solvers
 (Hindmarsh et al., 2005). Differential equation models contain parameters that affect state variables, and model-
 fitting through Markov chain algorithms involves iterating through parameter space one set of parameters at a time.
 We performed model fitting using a Markov chain algorithm called the Hamiltonian Monte Carlo (HMC), using
 version 2.18.1 of the RStan interface to the Stan statistical software (Carpenter et al., 2017; Guo et al., 2019) and
 version 3.4.1 of R (R Core Team, 2017). HMC is not a random walk algorithm and uses Hamiltonian mechanics to
 determine exploration steps in parameter space. HMC has been theorized to offer more efficient exploration of high-
 dimensional parameter space than traditional Random-Walk Metropolis algorithms (Beskos et al., 2013).

Conditional on the meta-analysis data set, the HMC algorithm computed posterior and posterior predictive
 distributions, from which Bayesian statistical inferences on likely ranges of parameter values were then made.
 Posterior distributions are the distributions of more likely model parameter values conditional on the data. Posterior
 predictive distributions are the distributions of more likely values for unobserved data points from the data-
 generating process conditional on the observations. In the case of this study, the experiments constituting the meta-
 analysis would be the data-generating process.

For the sake of clarity, it is important to distinguish between the frequentist confidence intervals and
 Bayesian posterior predictive intervals and distributions we describe in our study. Confidence intervals are
 calculated from the sample means and standard errors at observed data points and indicate ranges of values that are
 likely to contain the true data values with repeated sample collections using the same methodology. Posterior
 predictive intervals and distributions are computed after estimation of the posterior parameter distributions and
 represent the likely distributions of unobserved data values conditional on observed data values. Bayesian credible
 intervals, which we will also discuss in this study, are ranges of values that parameters are likely to take with some
 probability that are conditional on the observed data. Credible areas indicate the probability densities of parameter
 values across credible intervals.

We ran four chains for 25,000 iterations each for our HMC simulations, with the first 10,000 iterations
 being discarded as burn-in for each chain. Hence, our posterior distributions consisted of 100,000 posterior samples
 per HMC run. In retrospect, because our credible areas displayed sufficient smoothness (Supplemental Fig. 2) and

Deleted: analyzed the fit of

Deleted: (conventional)

Deleted: (Allison-Wallenstein-Bradford)

Deleted: (Allison et al., 2010). CON is a linear ordinary differential equation system, while AWB is a non-linear system (Supplemental Appendix 1). The models were chosen for this study due to their mathematical simplicity and limited data input requirements. Additionally, they

Deleted: based on

Deleted: and were compiled in a recent soil warming meta-analysis

Deleted: Models were fit to

Deleted: and each year

Deleted: ¶

Deleted: Model output response ratios were calculated by dividing simulated CO₂ flux following warming perturbation by the CO₂ flux at steady state.

Deleted: We chose to fit the response ratios rather than raw flux measurements for several reasons. First, there is no need to convert flux measurements from different experiments into a common unit. Second, response ratios represent a standardized metric for warming response across disparate ecosystem types with varying climate, soil, and vegetation properties. Finally, fitting a mean response ratio overcomes data gaps present in individual experiments. ¶

2.3 Markov Chain Monte Carlo Fitting ¶

→ We performed model fitting using a Markov chain Monte Carlo (MCMC) algorithm called the Hamiltonian Monte Carlo (HMC), using version 2.17

Formatted: Indent: First line: 0"

Deleted: to collect posterior distributions and posterior predictive distributions. Posterior distributions are the distributions of more likely model parameter values conditional on the data. Posterior predictive distributions are the distributions of more likely values for unobserved data points from the data-generating process conditional on the observations. In the case of this study, the experiments constituting the meta-analysis would be the data-generating process. ¶

Differential equation models contain parameters that affect state variables, and model-fitting through MCMC involves iterating through parameter space one set of parameters at a time.

Deleted: In the process of fitting and exploring parameter space with MCMCs, we obtained samples from

Deleted: parameter values. Bayesian inference is highly reliant on these distributions, as they provide information... [1]

Deleted: for

Deleted: for a given data set. For each HMC run, we

Deleted: 45

Deleted: 20

Deleted: in

260 Bayesian diagnostics indicated adequate posterior sampling (Supplemental Table 5), we could have reduced
262 simulation time without impairing posterior computation by running shorter chains that consisted of 20,000 to
264 30,000 iterations. To minimize the presence of divergent energy transitions, which indicate issues with exploring the
266 geometry of the parameter space specified by the prior distributions, we set the adaptation δ to 0.95, the initial
268 step size to 0.1, and maximum tree depth to 12. Those parameters determine how the HMC algorithm proposes new
sets of parameters at each step and were set so that the HMC would begin with smaller exploration steps. The
algorithm varies the step size from its initial value throughout posterior sampling to maintain a desired acceptance
rate; the tuning sensitivity of the step size is governed by the adaptation delta value, with higher values indicating
reduced sensitivity.

270 We further constrained our HMC runs to characterize parameter regimes corresponding to higher biological
272 realism. Normal informative priors were used to initiate the runs, and the prior distribution parameters were chosen
based on expert opinion and previous empirical observations (Allison et al., 2010; Li et al., 2014). Prior distributions
had non-infinite supports; supports were truncated to prevent the HMC from exploring parameter space that was
unrealistic (Supplemental Table 2).

274 2.4 Model Steady State Initialization

276 Because we were mainly interested in testing model predictions of soil warming response, the models were
278 initiated at steady state prior to the introduction of warming perturbation to isolate model warming responses from
steady state attraction. We fixed pre-perturbation steady state soil C densities to prevent HMC runs from exploring
parameter regimes corresponding to biologically unrealistic C pool densities and mass ratios.

280 To set pre-warming steady state soil C densities, we first analytically derived steady state solutions of the
282 ordinary differential equations of the models. Then, with the assistance of Mathematica version 12, we re-arranged
the equations by moving the steady state pool sizes to the left-hand side (Supplemental Appendix 2), such that we
284 could determine the value of parameters dependent on pool sizes while allowing the rest of the parameters to vary
for the HMC. Consequently, we could constrain the pre-warming pool sizes from reaching unrealistic values in the
simulations.

286 2.5 Sensitivity Analysis of C Pool Ratios

288 Sensitivity analyses examine how the distributions of model input values influence the distributions of
model outputs. In our study, we considered pre-warming C-pool densities as a model input. We performed a
sensitivity analysis to observe how the choice of pre-warming C pool densities and C-pool ratios would affect the
model fits and posterior predictive distribution of C pool ratios.

290 We compared the model outputs and post-warming response behavior of AWB and CON at equivalent C
292 pool densities and ratios. The ratio of soil microbe biomass C (MIC) density to SOC density has been observed to
294 vary approximately from 0.01 to 0.04 (Anderson and Domsch, 1989; Sparling, 1992), so we used those numbers as
guidelines for establishing the ranges of the C pool densities and density ratios explored in our simulations. One
296 portion of the analysis involved running HMC simulations in which we set the pre-warming MIC density at 2 mg C
g⁻¹ soil and then varied the SOC density from 50 to 200 mg C g⁻¹ soil in increments of 25, stepping from 0.04 to 0.01
with respect to the MIC-to-SOC ratio. A second portion of the analysis involved observing the effect of varying
pre-warming MIC from 1 to 8 mg C g⁻¹ soil while holding pre-warming SOC at 100 mg C g⁻¹ soil.

298 For some combinations of the prior distributions and pre-warming steady state C pool densities
(Supplemental Table 2), AWB HMC runs wandered into unstable parameter regimes that would prevent the
300 algorithm from reliably running to completion. Consequently, we do not compare simulation results for AWB and
302 CON with pre-warming SOC densities below 50 mg C g⁻¹ soil. Other combinations of prior distribution and pre-
warming C pool density choices that were not necessarily biologically realistic allowed stable AWB runs with lower
pre-warming SOC densities.

304 2.6 Information Criteria and Cross-validation

306 In addition to R², we used the WAIC, LPML, and LOO Bayesian predictive goodness-of-fit metrics to
308 evaluate models with the meta-analysis warming response data. LPML is an example of cross validation that is
calculated similarly to LOO (Gelfand et al., 1992; Gelfand and Dey, 1994; Ibrahim et al., 2001) but differs from
LOO in how the importance ratio sampling portion of its computation is handled. For further explanation regarding
importance ratios and their role in evaluating approximate cross-validation metrics, refer to the description of the

Deleted: and

Deleted: HMC

Deleted: respectively to 0.9995 and 0.001. These parameters control

Deleted: (Allison et al., 2010; Li et al., 2014)

Deleted: The fraction of soil microbe biomass C (MIC) density to SOC density has been observed to vary approximately between 0.01 – 0.04 (Anderson and Domsch, 1989; Sparling, 1992)

Deleted: in terms of

Deleted: fraction

Deleted: to

Deleted: In addition to R², we used the WAIC, LOO, and Log Pseudomarginal Likelihood (LPML) Bayesian predictive goodness-of-fit metrics to evaluate models with the meta-analysis warming response data. LPML is also an example of cross validation and is calculated similarly to LOO. However, LPML does not account for over-fitting or penalize for parameter count (Christensen et al., 2011). We used the 'loo' package available for R to calculate our WAIC and LOO values

332 LOO algorithm presented in Vehtari and Ojanen (2012). LOO updates LPML by implementing a smoothing process
334 in which the largest importance ratios are fitted with a Pareto distribution and then replaced by expected values from
336 the distribution, which stabilizes the importance ratio sampling.

334 Algorithmic differences between WAIC and LPML and LOO render them appropriate for different
336 statistical modeling goals and make them complementary metrics. WAIC is suitable for estimating the relative
338 quality of model fits to hypothetical repeated samples collected at existing experimental time points, whereas LOO
340 and LPML are suitable for estimating the quality of fits to hypothetical measurements taken between observed time
342 points (Vehtari et al., 2017).

340 We used version 2.0.0 of the loo package available for R to calculate our WAIC and LOO values (Vehtari
342 et al., 2019). A lower WAIC and LOO and a higher LPML indicate a more likely model for a given data set. LPML
344 can be multiplied by a factor of -2 to occupy a similar scale to LOO.

3 Results

3.1 Parameter Posterior Distributions

344 We obtained distributions of posterior predictive fits to the univariate response ratio data for both AWB
346 and CON across different pre-warming MIC-to-SOC ratios. Posterior samples totaled 100,000 for each simulation.
348 Sampler diagnostics for the HMC runs indicated that the statistical models were valid at all pre-warming steady state
350 values observed (Supplemental Table 6), that model parameter values converged across the four Markov chains
352 (Supplemental Fig. 7), and that the posterior parameter space was effectively sampled and explored (Supplemental
354 Fig. 5) to generate enough independent posterior samples for inference (Supplemental Fig. 6). The ratios of effective
356 posterior parameter samples to total samples for parameters were generally satisfactory: across observed MIC-to-
358 SOC ratios, they were all greater than 0.25 and mostly greater than 0.5 (Supplemental Table 5).

352 We also tracked divergent transitions, which mark points in chains at which the HMC algorithm was
354 inhibited in its exploration and posterior sampling, potentially due to the parameter space becoming geometrically
356 confined and difficult to navigate. Divergent transitions occurred in the AWB HMC runs (Supplemental Fig. 9),
358 though the ratios of divergent transitions to sampled iterations was relatively low for all runs. The highest divergent
360 transition ratio observed was 0.0217, corresponding to the simulation initiated with pre-warming SOC = 200 mg C
362 g⁻¹ soil. There were no divergent transitions in the CON runs.

3.2 Model Behaviors

360 The CON curve monotonically decreases in response ratio over time, whereas the AWB curve displays
362 changes in slope sign (Fig. 2). The difference in curve shape (Fig. 3a, b) is in line with CON's linear status and
364 AWB's non-linear formulation with more parameters (Allison et al., 2010). By 50 years after warming, mean fit
366 curves for AWB and CON return to 1.0 after their initial increase (Fig. 3c, d), consistent with prior observations and
368 expectations at steady state (van Gestel et al., 2018; Romero-Olivares et al., 2017).

364 From a cursory visual evaluation, neither of the models clearly out-performs the other across all
366 prewarming steady states. The 95% confidence interval of the first data point at t = 0.5 years does not include the
368 AWB SOC100 posterior predictive mean as it does for the CON SOC100 mean (Fig. 2), which most likely impaired
370 AWB's quantitative goodness-of-fit metrics. However, the 95% response ratio posterior predictive interval suggests
372 that AWB is able to replicate the response ratio increase in the data from 1.5 to 3.5 years following the warming
374 perturbation, which CON does not. The shape of the AWB posterior predictive interval also fits the data points and
376 confidence intervals occurring eight years or more after the perturbation more closely than that of CON (Fig. 3a, b).

372 For both AWB and CON, increasing the pre-warming SOC to higher densities from SOC = 50 to 200 mg C
374 g⁻¹ soil (hereby labeled from SOC50 to SOC200) while holding pre-warming MIC at 2 mg C g⁻¹ soil, DOC at 0.2 mg
376 C g⁻¹ soil, and ENZ at 0.1 mg C g⁻¹ soil, corresponded to lower initial mean response ratios in the first year at the t =
378 0.5 year time point, which certainly inhibited the quantitative goodness-of-fit (Fig. 3a, b). For CON, increasing pre-
380 warming SOC also reduced the magnitude of the mean fit slope. For AWB, increasing pre-warming SOC had no
clear effect on the curve slope, but the model needed more time to achieve peak mean response ratio from a lower
start, with the peak being reached at t = 1.5 years in the SOC50 case and t = 3.5 years in the SOC200 case (Fig. 3b).
At higher pre-warming SOC, CON's reduced slope magnitude and AWB's lagging response ratio peak caused both
models to exhibit slower returns to the steady state response ratio of 1.0 (Fig. 3c, d). On their trajectories back to
steady state, the mean SOC200 CON curve substantially overshoots the data means after t = 7.5 years (Fig. 3a),

Moved up [2]: (Vehtari et al., 2017).

Formatted: Font color: Text 1

Deleted: A lower WAIC and LOO and a higher LPML indicate a more likely model for a given data set.

Formatted: Font color: Text 1

Deleted: posterior parameter

Deleted: and

Deleted: generally

Deleted: convergence for

Deleted: Markov chains and usable posteriors

Deleted: Fig 5 – 7).

Deleted: that indicate the presence of regions of parameter space that are too

Deleted: explore by the HMC.

Deleted: ,

Deleted: none exceeding 0.025.

Deleted: Effective sample proportion for parameters was generally satisfactory and greater than 0.3 for parameters.. [2]

Deleted: system

Deleted: (Allison et al., 2010).

Deleted: -

Deleted: the

Deleted: mean

Deleted: ,

Deleted: could negatively impact

Deleted: and information criteria

Deleted: model

Deleted: credible

Deleted: trend of

Deleted: 1-

Deleted: mean

Deleted: fit

Deleted: matches

Deleted: after

Deleted: CON. Visually, though, it is not clear

Deleted: model provides the better

Deleted: .

422 whereas the SOC200 AWB curve exceeds the data means at a more moderate extent through the $t = 8.5, 9.5, 10.5$
423 and 11.5 year time points (Fig. 3b).

424 Changing the pre-warming MIC-to-SOC steady state pool size ratio by increasing pre-warming MIC from
425 1 to 8 mg C g⁻¹ soil (hereby labeled from MIC1 to MIC8) while holding pre-warming SOC at 100 mg C g⁻¹ soil had
426 marginal to moderate qualitative effects on the mean response ratio curves for CON and AWB. The CON MIC1 and
427 MIC8 curves are visually indistinguishable (Supplemental Fig. 1a, b), while the AWB MIC1 and MIC8 curves differ
428 with the MIC8 curve displaying more gradual changes in slope and lower slope magnitudes (Supplemental Fig. 1c,
d).

3.3 Sensitivity Analysis of Parameter Distributions to Pre-warming C Pool Densities and Density Ratios

430 In addition to response ratio fits, we observed the influence of pre-warming MIC-to-SOC ratios on model
431 SOC stock response ratios in AWB and CON simulations following warming. Similar to the model flux response
432 ratios, SOC response ratios were calculated by dividing evolved post-warming SOC densities by pre-warming
433 densities. The SOC response ratios at 12.5 years for CON and AWB increased as pre-warming SOC was raised (and
434 hence, the MIC-to-SOC ratio decreased) with other pre-warming C densities held constant, indicating reduced
435 proportional SOC loss when SOC stocks were initiated at higher pre-warming densities (Supplemental Fig. 3a). For
436 CON, SOC loss decreased from 27.1% at SOC50 to 9.2% at SOC200. In a similar trend for AWB, SOC loss
437 decreased from 17.2% at SOC50 to 8.1% at SOC200. In contrast, raising pre-warming MIC densities (and hence,
438 increasing the MIC-to-SOC ratio) with other pre-warming C densities held constant did not produce a shared trend
439 for CON and AWB (Supplemental Fig. 3b). CON SOC loss decreased from 18.8% at MIC1 to 17.4% at MIC8
440 while AWB SOC loss increased from 11.3% at MIC1 to 16.3% at MIC8.

441 Truncation of prior supports, or distribution domains, generally did not prevent posterior densities from
442 retaining normal distribution shapes. Deformation away from Gaussian shapes for the densities of E_{a_s} from CON
443 was observed at SOC50 and SOC75. For AWB, deformation was observed for the densities of E_{a_v} , E_{a_k} , and $E_{C_{ref}}$.
444 All CON and AWB parameter posterior densities were otherwise observed to be Gaussian from SOC100 to
445 SOC200. Example posterior densities and means for select model parameters at pre-warming SOC100 are presented
446 in Fig. 4 and Supplemental Fig. 2. Parameter posterior means corresponding to other pre-warming C pool densities
and ratios are presented in Supplemental Table 3.

3.4 Sensitivity Analysis of Quantitative Fit Metrics to Pre-warming C Pool Densities and Density Ratios

450 For both CON and AWB, LOO, WAIC, LPML, and R² all worsened as pre-warming steady state SOC
451 density was increased from SOC50 to the less biologically realistic SOC200 (Fig. 5). CON's LOO and WAIC values
452 increased respectively from -15.704 and -15.818 at SOC50 to -6.891 and -6.966 at SOC200, while AWB's LOO and
453 WAIC values increased respectively from -11.028 and -11.379 at SOC50 to -5.97 and -6.579 at SOC200
454 (Supplemental Table 4a, b). Compared to AWB's metrics, CON's goodness-of-fit metrics deteriorated at a faster
455 rate with the increase of pre-warming SOC. Nonetheless, CON outperformed AWB in LOO, WAIC, and LPML
456 across all observed pre-warming SOC densities. The Bayesian metrics accounted for AWB's larger model size and
457 increased propensity for overfitting as demonstrated by the consistently higher effective parameter counts associated
with AWB (Supplemental Fig. 8a, b).

458 Varying pre-warming steady state MIC from MIC1 to MIC8 modestly impaired goodness-of-fit across the
459 various metrics (Supplemental Fig. 4). CON's LOO and WAIC values increased respectively from -11.963 and -
460 12.035 at MIC1 to -11.731 and -11.802 at MIC8, while AWB's LOO and WAIC values increased respectively from
461 -8.63 and -9.302 at MIC1 to -8.181 and -8.711 at MIC8 (Supplemental Table 4c, d). CON did not deteriorate in
462 goodness-of-fit at a faster rate than AWB with respect to increasing pre-warming MIC. Increasing pre-warming
463 MIC has the opposite effect on MIC-to-SOC ratio compared to increasing pre-warming SOC, but both changes
464 worsened goodness-of-fit across all metrics, indicating that changes to pre-warming MIC-to-SOC ratio did not
produce consistent trends.

4 Discussion

465 Our study develops a quantitative, data-driven framework for model comparison that could be applied
466 across different research questions, ecosystems, and scales. We demonstrated the novel deployment of WAIC and
467 LOO, two more recently developed Bayesian goodness-of-fit metrics that estimate model predictive accuracy, to

Deleted: → For both AWB and CON, higher pre-warming SOC corresponds to lower initial response ratio (Fig 3a-b). For CON, higher initial SOC reduces the magnitude of the mean fit slope and slows the return of the response curve to 1.0. For AWB, more time is needed to reach the peak response ratio and return to pre-warming response ratios. Changing the pre-warming MIC-to-SOC steady state pool size ratio by increasing MIC has a subtle effect on the fit curve; the magnitude and severity of slope changes decreases from MIC = 1 to MIC = 8 mg C g⁻¹ soil (Supplemental Fig 1). Increasing MIC did not have an appreciable qualitative effect on CON fit.¶

Deleted: fit

Deleted: fractional

Deleted: loss for

Deleted: fractional

Deleted: loss

Deleted: decreased

Deleted: increased

Deleted:).

Deleted: ranged

Deleted: SOC = 50

Deleted: SOC = 200 (Supplemental Fig 3). For

Deleted: it ranged

Deleted: SOC = 50 to 8.1% at SOC = 200. Similarly, AWB's

Deleted: 16.3%

Deleted: 11.3% as MIC was reduced from 8 to 1. In contrast,

Deleted: 17.4% to 18.8% when MIC was reduced from 8 to 1.¶

Deleted: ¶

Deleted: was observed at SOC = 50 mg C g⁻¹ soil and SOC

Deleted: E_{as} for

Deleted: E_{av}, E_{ak}, and E_{Cref} for AWB.

Deleted: SOC = 100 mg C g⁻¹ soil

Deleted: SOC = 200 mg C g⁻¹ soil.

Deleted: SOC = 100 mg C g⁻¹

Deleted: plotted

Deleted: → ¶ ... [7]

Deleted: appeared

Deleted: slightly reduce

Deleted: quality

Deleted: as MIC ranged from 1 to 8 mg C g⁻¹ soil

Deleted:), though the trend was

Deleted: consistent

Deleted: LOO and WAIC. Since

Deleted: SOC, these results indicate no consistent effect for

536 evaluate SBMs using data from longitudinal soil warming experiments. WAIC and LOO improve upon older and
538 more frequently used metrics, such as AIC and DIC, by accounting for model complexity and overfitting of data in a
540 more comprehensive, stable, and accurate fashion. The quantitative agreement between WAIC, LOO, and LPML
reinforces the reliability and validity of information criteria and cross-validation metrics to complement use of
frequentist R^2 .

542 We constrained the fitting of AWB and CON to biologically reasonable parameter space by fixing pre-
warming steady state C pool densities and establishing prior distributions informed by expert judgment
(Supplemental Table 2). We observed that, despite the qualitative difference in the shapes of their mean posterior
544 predictive fit curves, CON and AWB could both potentially account for the soil warming response in the meta-
analysis data set. For both models, posterior predictive fit distributions overlapped with the confidence intervals of
546 the data points (Fig. 2). However, with respect to the Bayesian goodness-of-fit metrics, CON quantitatively
outperformed AWB across all pre-warming SOC and MIC densities observed (Fig. 5 and Supplemental Fig. 4)
548 because the Bayesian metrics adjusted for AWB's larger model size and consistently higher effective parameter
count (Supplemental Fig. 8). For both models, lower pre-warming SOC densities corresponded to better warming
550 response fits (Fig. 5).

4.1 Model Responses to Warming over Time

552 After fitting, the response ratio curves of CON and AWB both trended toward the pre-warming steady state
response ratio of 1.0 following the soil warming perturbation (Fig. 3). The settling of the curves to the pre-warming
554 model steady states aligns with previous literature which demonstrated that the magnitude of CO_2 flux tends to fall
after reaching a post-warming maximum (Crowther et al., 2016; Romero-Olivares et al., 2017). In the meta-analysis
556 data set, this peak is reached immediately at the first data point at $t = 0.5$ years (Fig. 2). CON matched this data
pattern in all of our observed simulations in outputting maximum response ratios at the first time point after
558 warming (Fig. 3a, c and Supplemental Fig. 1a, b). AWB was unable to output maximum response ratios at the first
time point (Fig. 3b, d) and was therefore penalized in quantitative goodness-of-fit. Examining AWB's system of
560 equations (Supplemental Appendix 1b), we surmise that one reason for the later peak was due to the slower growth
of MIC in the biologically truncated parameter space that AWB was limited to. MIC is a driving force for the
562 increase of CO_2 flux as a numerator term in the AWB flux equation (Supplemental Appendix 1b, Equation A10).
Unlike MIC biomass in CON (Supplemental Appendix 1a, Equation A3), MIC biomass growth in AWB has two
564 loss terms in its differential equation (Supplemental Appendix 1b, Equation A8).

566 This is not to say that CON was clearly superior from a qualitative standpoint. CON's mean posterior
predictive curves were not able to match a subsequent local data maximum in the meta-analysis data set at $t = 3.5$
568 years, a trend which AWB's curves were able to replicate. The mean CON curves also substantially overshoot the
data at later time points following $t = 7.5$ years (Fig. 2a, Fig. 3a, c, and Supplemental Fig. 1a, b) because of the
inability of first order linear models such as CON to display oscillatory dynamics (Hale and LaSalle, 1963).

570 In contrast, AWB displays damped oscillations in its response ratios following warming due to its non-
linear dynamics (Fig. 2 and Fig. 3). AWB was able to match the points after $t = 7.5$ years more closely than CON.
572 The presence of respiration oscillations has been observed in long-term warming experiments, such as the one taking
place at Harvard Forest (Melillo et al., 2017). It is possible AWB would be quantitatively rewarded in goodness-of-
574 fit metrics over CON for its ability to replicate biologically realistic oscillations in larger, site-specific data sets such
as those from Harvard Forest.

4.2 Sensitivity Analyses of C Pool Densities and Density Ratios

578 We performed a goodness-of-fit sensitivity analysis to check whether the response ratio trends stayed
consistent, biologically realistic, and interpretable across a range of pre-warming, steady state soil C densities and
580 pool-to-pool density ratios. For instance, we imposed constraints to reflect that MIC-to-SOC density ratios range
between 0.01 and 0.04 across various soil types (Anderson and Domsch, 1989; Sparling, 1992). CON and AWB
582 response ratio curves exhibited realistic values and qualitatively consistent shapes across all pre-warming SOC and
MIC steady state densities, even at less realistic SOC densities above 100 mg C g^{-1} soil (Fig. 3). There was enough
584 uncertainty in the data that the 95% posterior predictive intervals for the model output always overlapped with the
95% confidence intervals of each fitted data point (Fig. 2). In most cases, the posterior mean response ratio curve
586 also fell within the 95% data confidence interval.

588 We were unable to initiate our pre-warming SOC steady state density below SOC50 with the priors and
MIC-to-SOC ratios used for AWB. Under SOC50, AWB HMC runs would not reliably run to conclusion and would

Deleted: can

Deleted: explain

Deleted: to warming

Deleted:) and that certain

Deleted: soil C densities and density ratios for

Deleted: correspond

Deleted: CON and AWB both displayed similar general trends in

Deleted: progression of their

Deleted: 2

Deleted: return

Deleted: their

Deleted: demonstrates

Deleted: falls following

Deleted: peak

Deleted: . ¶

AWB, unlike CON,

Deleted: . However, it is unclear whether oscillations quantitatively aid AWB with its fit to our response ratio data set

Deleted: For an additional check on model realism, we tallied SOC loss percentages from pre-warming SOC stocks after 12.5 years for AWB and CON. SOC losses ranged from 8.14% to 27.1% across both models (Supplemental Fig 3). These results aligned with a recent comprehensive meta-analysis of 143 soil warming studies (Supplemental Fig 10). The largest loss of 27.1%, occurring in CON at SOC = 50, is sizable, but the van Gestel et al.

Moved down [3]: meta-analysis included 7 studies measuring losses greater than 20%, with the maximum loss observed at 54.4% (van Gestel et al., 2018). ¶

Deleted: For both AWB and CON, increasing pre-warming SOC reduced C loss fraction following the perturbation. Varying pre-warming MIC more prominently affected the fraction of SOC lost from AWB compared to CON, with soil C loss increasing as MIC increased. In CON's case, there was a minimal decline in SOC loss as MIC was increased. The larger effect of increasing MIC on the fraction of SOC lost in AWB is likely due to MIC influence on SOC-to-DOC turnover, which is not a feedback included in the CON model. ¶

Deleted: Analysis

Deleted: We performed a sensitivity analysis to check whether the response ratio trends stayed consistent, biologically realistic, and interpretable across a range of pre-warming, steady state soil C densities and pool-to-pool density ratios. For instance, we imposed constraints to reflect that MIC-to-SOC density ratios range between 0.01 and 0.04

Deleted: 50 mg SOC g^{-1} soil

Deleted: 50 mg SOC g^{-1} soil

658 terminate due to ODE instabilities. Even at SOC50, we saw a reduction in independent and effective samples for
660 certain parameters, namely E_{a_K} and $E_{C_{ref}}$ (Supplementary Table 5). We did not drop under SOC50 for CON, as we
662 sought to compare AWB and CON at similar MIC-to-SOC ranges. Our experience underscores the challenge of
664 choosing realistic steady state soil C densities, density ratios, and prior distributions to obtain valid model
666 comparisons limited to biologically realistic regimes.

662 The information criteria and cross-validation fit metrics generally indicated higher relative probability and
664 predictive performance at lower pre-warming SOC values for AWB and CON (Fig. 5). The fit results suggest that
666 SOC density of the soil at the sites included in the meta-analysis was likely closer to the lower end of the SOC
668 density ranges examined in our sensitivity analysis. A less pronounced trend toward better fits was observed as pre-
warming MIC density was decreased while pre-warming SOC density was held constant (Supplemental Fig. 4). No
clear relationship was observed between MIC-to-SOC ratio and goodness-of-fit in the AWB and CON models.

668 The worsening IC and CV results at higher SOC densities support the notion that pre-warming steady state
670 SOC densities should not be initialized over SOC100 in AWB and CON when fitting to this meta-analysis data set.
672 Pre-warming SOC density was not observed to exceed 50 mg SOC g⁻¹ soil at sites included in the meta-analysis,
674 reaching a maximum of 45 mg SOC g⁻¹ soil for the top 20 cm in one study with alpine wetland soil (Zhang et al.,
2014). The majority of the CO₂ respired by soil microbes is sourced from surface soil (Fang and Moncrieff, 2005),
and it is well-documented that SOC densities increase toward the soil surface (Jobbágy and Jackson, 2000). ¹⁴C
676 measurements of CO₂ fluxes suggest that SOC densities representing the source of most heterotrophic respiration
range between 40 to 80 mg SOC g⁻¹ soil (Trumbore, 2000), so the effective SOC densities associated with soil
respiration at some meta-analysis sites may have been in this range.

678 Overall, the Bayesian metrics from the goodness-of-fit sensitivity analysis suggest that CON is superior to
AWB at explaining the meta-analysis data set when accounting for model parsimony, particularly when the models
680 are initiated in more realistic ranges of pre-warming SOC densities under SOC100. However, we caution against
using these results to conclude that CON is a comprehensively superior predictive model over AWB without
682 comparisons involving other longitudinal soil warming data sets. And other data aside, we observe that AWB has a
useful advantage over CON conditional on the meta-analysis data set alone: AWB was more tolerant of changes in
684 pre-warming conditions, displaying less IC and CV than CON as pre-warming SOC is increased (Fig. 5a – c).
AWB's compensatory ability stemming from its larger model size could be more quantitatively rewarding in
686 goodness-of-fit sensitivity analyses conducted on data assimilations with larger data sets.

686 For an additional check on the biological realism and plausibility of our simulations, we conducted a
688 sensitivity analysis examining changes in model SOC stocks following warming. The response ratios of post-
warming SOC stocks after 12.5 years, evaluated as the ratio of post-warming to pre-warming SOC densities, was
690 computed from observed CON and AWB simulations at the posterior parameter means. SOC losses indicated by the
response ratios ranged from 8.13 to 27.1% across both models (Supplemental Fig. 3). These results aligned with a
692 recent comprehensive meta-analysis of 143 soil warming studies (Supplemental Fig. 10). The largest loss of 27.1%,
occurring in CON at SOC50, is sizable, but the meta-analysis included 7 studies measuring losses greater than 20%,
with the maximum loss observed at 54.4% (van Gestel et al., 2018).

694 Raising pre-warming SOC reduced SOC loss after 12.5 years of warming for both models (Supplemental
Fig. 3a). For CON, SOC loss decreased from 27.1% at SOC50 to 9.2% at SOC200. For AWB, SOC loss decreased
696 from 17.2% at SOC50 to 8.13% at SOC200. Varying pre-warming MIC affected the SOC response ratio more
substantially for AWB than CON (Supplemental Fig. 3b). For AWB, SOC loss increased from 11.4% at MIC1 to
698 16.3% at MIC8, while SOC loss decreased from 18.8% at MIC1 to 17.4% at MIC8 for CON. The larger effect of
increasing MIC on the SOC response ratio in AWB is likely due to MIC influence on SOC-to-DOC turnover, which
700 is not a feedback accounted for in the equations of the CON model (Supplemental Appendix 1a).

702 The posterior means for the Arrhenius activation energy parameters E_a of CON and AWB returned by the
HMC simulations across the observed pre-warming C densities (Supplemental Table 3) differed somewhat from the
704 parameter values used in Allison et al. (2010) and Li et al. (2014), which were in turn tuned based on activation
energies estimated in a prior empirical analysis of enzyme-catalyzed soil organic matter decomposition processes
(Trasar-Cepeda et al., 2007). In Allison et al. (2010), CON parameters $E_{a_{S_2}}$, $E_{a_{D_2}}$, and E_{a_M} were respectively set at
706 47, 40, and 40 kJ mol⁻¹ and AWB parameters E_{a_V} and $E_{a_{VU}}$ were both set at 47 kJ mol⁻¹. The AWB Michaelis-
Menten K_M terms were not parameterized to have Arrhenius temperature dependence in Allison et al. (2010). In Li
708 et al. (2014), CON parameters $E_{a_{S_2}}$, $E_{a_{D_2}}$, and E_{a_M} were set at 47, 47, and 20 kJ mol⁻¹ and AWB parameters $E_{a_{V_2}}$,
 $E_{a_{VU_2}}$, $E_{a_{K_2}}$, and $E_{a_{KU}}$ were set at 47, 47, 30, and 30 kJ mol⁻¹. These values were in line with the activation energies
710 calculated in Trasar-Cepeda et al. (2007), which ranged from 17.0 to 57.7 kJ mol⁻¹, with the energies corresponding

Deleted: 50 mg SOC g⁻¹ soil

Deleted: E_{a_V}

Deleted: E_{a_K}

Deleted: 13

Deleted: 50 mg SOC g⁻¹ soil

Deleted: Similarly, we were unable to drop our pre-warming MIC steady state below 1 mg SOC g⁻¹ soil.

Deleted: for the data

Deleted: soil C densities should not be initialized over 100 mg C g⁻¹ soil in AWB and CON when fitting to this meta-analysis data set.

Moved down [4]: The majority of the CO₂ respired by soil microbes is sourced from surface soil (Fang and Moncrieff, 2005), and it is well-documented that

Deleted: the highest SOC densities are in the top 20 centimeters of soil (Jobbágy and Jackson, 2000).

Moved (insertion) [4]

Deleted: ¹⁴C measurements of CO₂ fluxes suggest that SOC densities representing the source of most heterotrophic respiration in topsoil range between 40 to 80 mg SOC g⁻¹ soil (Trumbore, 2000).

Deleted: 4.3

Moved (insertion) [3]

732 to the decomposition of plant litter and protected organic matter being on the higher end and the energies
733 corresponding to microbial biomass degradation being on the lower.

734 Our HMC simulations arrived at higher E_a values, with the posterior means of E_{a_S} , E_{a_D} , and E_{a_M}
735 respectively ranging from 51.3 to 77.6 kJ mol⁻¹, 50.1 to 50.3 kJ mol⁻¹, and 51.8 to 52.6 kJ mol⁻¹ in the pre-warming
736 SOC-varied simulations for CON, and the posterior means of E_{a_V} , $E_{a_{VU}}$, E_{a_K} , and $E_{a_{KU}}$ respectively ranging
737 from 58.5 to 74.8 kJ mol⁻¹, 50.2 to 51.1 kJ mol⁻¹, 25.8 to 42.4 kJ mol⁻¹, and 49.0 to 49.8 kJ mol⁻¹ for AWB.
738 However, these values are still within the ranges of organic matter decomposition activation energies, which have
739 been empirically estimated to exceed 100 kJ mol⁻¹ at their highest in the A-horizons of temperate soils (Steinweg et
740 al., 2013), suggesting that the E_a posterior means, aided by prior truncation, effectively remained within
741 biologically realistic space across all observed pre-warming C densities. The presence of higher E_{a_S} posterior means
742 also agreed with the empirical trends of higher activation energies for the degradation of SOC-related organic
743 compounds and lower activation energies for the degradation of material associated with microorganisms.

744 We found it less useful to compare the posterior means of other fitted parameters including the C pool
745 transfer coefficients, C use efficiency E_{c_u} and V_{max} to empirical estimates for biological benchmarking purposes.
746 Unitless parameters like transfer coefficients and E_c defy straightforward interpretation, measurement, and
747 estimation from experiments (Bradford and Crowther, 2013). Very different values can be found based on whether
748 substrate-specific or substrate-nonspecific assumptions and methods are used (Geyer et al., 2019; Hagerty et al.,
749 2018). V_{max} parameters are not unitless but display even higher variance than the bounded C transfer and efficiency
750 coefficients. The V_{max} parameter corresponding to a specific enzyme can vary over orders of magnitude when the
751 sensitivity of the enzyme to an interval of temperatures is considered (Nottingham et al., 2016). The process of
752 consolidating experimental substrate-specific and substrate-nonspecific measurements into a single number to
753 correspond to a model V_{max} value introduces further complications and uncertainty, rendering comparisons of
754 potentially drastically different V_{max} values less informative regarding model biological realism.

4.3 HMC Parameter Space Exploration

756 Truncating prior and posterior parameter distributions proved useful for establishing biological constraints
757 and only modestly deformed posterior densities for AWB and CON. From SOC100 to SOC200, CON and AWB
758 posterior densities showed little or no deformation from typical normal distribution shapes. Moderate posterior
759 density deformation was observed for some parameters in both models at SOC50 and SOC75, namely E_a for CON
760 and $E_{c_{ref}}$ for AWB (Supplemental Fig. 11). Even so, most of the other parameter posterior densities still remained
761 undeformed at those SOC values. Thus, prior truncation generally did not prevent posterior means from falling
762 within biologically realistic intervals, suggesting that priors were appropriately informed and chosen.

763 A small frequency of divergent transitions was detected in the AWB HMC simulations. Divergent
764 transitions can be thought of as algorithm trajectory errors arising during the HMC's exploration of a convoluted
765 region of parameter space: a more thorough description of the theory, computation, and implications of divergent
766 transitions can be found in literature focusing on the Hamiltonian Monte Carlo algorithm (Betancourt, 2016, 2017).
767 The number of divergent transitions generally increased as the pre-warming MIC-to-SOC steady state ratio was
768 reduced (Supplemental Fig. 9). Prior truncation and the fixing of select parameters to constrain the pre-warming
769 steady state mass values for biological realism could have played a combined role in generating the Markov chain
770 divergences by hindering the smooth exploration of parameter space. We were unable to eliminate divergent
771 transitions by adjusting HMC parameter proposal step size, suggesting that other methods, such as modification of
772 the HMC algorithm itself or introduction of auxiliary parameters to AWB that reduce correlation between existing
773 model parameters may be more applicable in reducing divergent transitions in our case (Betancourt and Girolami,
774 2015). Additionally, the interaction between the ranges of values used for the prior distributions and the limited
775 number of observations in the data set could have contributed to the shaping of geometric inefficiencies (Betancourt,
776 2017).

777 It is possible that the instability that prevented consistent solving and HMC exploration of AWB under
778 SOC50 could be traced to the forward Michaelis-Menten formulation of decomposition and uptake kinetics used in
779 the present version of the AWB model (Supplemental Appendix 1 Equations A7, A8). We initialized the system
780 with a small DOC density lower than that of MIC at 0.1 mg C g⁻¹ soil. Since DOC was in the denominator of these
781 decomposition and uptake expressions, those expressions could become larger than tolerable for the system in
782 certain parameter regimes.

783 Some suggestions for the re-parameterization of AWB to improve model stability have been proposed that
784 could reduce or even eliminate divergent transitions by facilitating a smoother and steadier parameter space

Deleted: pre-warming SOC = 100

Deleted: SOC = 200

Deleted: SOC = 50

Deleted: 75 ($E_{c_{ref}}$)

Deleted: and E_{a_S} for CON

Deleted: for the AWB HMC runs. A

Deleted: (Betancourt, 2016, 2017).

Deleted: (Betancourt, 2017)

Formatted: Font: Not Bold

Deleted: 4.4 Applying and Interpreting Bayesian Predictive Fit Metrics

With respect to the IC and CV metrics, in both Fig 5 and Supplementary Fig 5, there is disagreement between LOO and WAIC versus LPML. LPML displays more consistent trends for CON and AWB across the range of pre-warming SOC values with a unidirectional change in slope. LPML is calculated similarly to LOO but does not account for overfitting and parameter count (Gelfand and Dey, 1994; Gelman et al., 2014). The computational difference accounts for the divergence between the results of LPML and those of LOO and WAIC. The effective parameter count and penalty for overfitting in both the WAIC and LOO calculations generally increases as pre-warming SOC is reduced (Supplemental Fig 8a and 8b). Thus, while the LPML results appear clearer, we do not recommend use of LPML by itself to quantitatively compare model fits because it does not fully account for the impacts of differing model structure, parameterization, and parameter count on overfitting for a data set.

General agreement between WAIC, LOO, and LPML reinforces the usage of IC and CV metrics alongside usage of R^2 . R^2 is not suitable as sole quantitative metric for model evaluation and selection. The traditional unadjusted R^2 calculation does not have a cost function for parameter counts. R^2 estimates the strength of the relationship between a linear model and a dependent variable and is calculated from the variance in data and residuals separating model outputs from observations. The metric cannot be applied to nonlinear models. Model selection involves a relative comparison of models, but the value of R^2 can result in misleading conclusions regarding absolute goodness of fit of a model to data. For instance, a model appropriate for a data set can correspond to a low R^2 calculation, while a flawed model can correspond to a high R^2 (Spiess and Neumeier, 2010). Adjusted R^2 accounts for model parameter count, but still shares other pitfalls with non-adjusted R^2 .

4.5 Conclusion and Future Directions

Recent SBM comparisons have been unable to demonstrate the superiority of one model over another because the uncertainty boundaries of the data were not sufficient for distinguishing model outcomes (Sulman et al., 2018; Wieder et al., 2018). Similar to Sulman et al., our results indicate that more data is needed to constrain model outputs and to verify the strengths and limitations of linear versus nonlinear SBMs in Earth system modeling.

840 [conductive for HMC exploration. One intermediate possibility would be to modify AWB to use reverse Michaelis-](#)
842 [Menten kinetics, which would replace the DOC term in the denominators of the decomposition and uptake](#)
844 [expressions with the larger MIC term. The use of reverse instead of Michaelis-Menten dynamics has been used to](#)
846 [stabilize and constrain other SBMs \(Sulman et al., 2014; Wieder et al., 2015b\). A more extensive re-formulation](#)
848 [involves the replacement of Michaelis-Menten expressions with equilibrium chemistry approximation \(ECA\)](#)
[kinetics, which would increase the number of denominator terms in decomposition expressions for further stability.](#)
[ECA equations have been shown to be more consistent in behavior and robust to parameter regime variation than](#)
[their Michaelis-Menten counterparts, and thus have been encouraged as a wholesale replacement for Michaelis-](#)
[Menten formulations \(Tang, 2015; Wang and Allison, 2019\). These re-parameterizations should be implemented](#)
[and examined in future work that involves sampling and computation of AWB posteriors.](#)

4.4 Outlook and Conclusions

850 [Recent SBM comparisons have been unable to demonstrate the superiority of one model over another](#)
852 [because the uncertainty boundaries of the data were not sufficient for distinguishing model outcomes \(Sulman et al.,](#)
854 [2018; Wieder et al., 2014, 2015b, 2018\). Similar to these previous studies, our results indicate that more data is](#)
856 [needed to constrain and differentiate between model posterior predictive distributions. Conditional on the meta-](#)
[analysis data set, CON demonstrates superior quantitative goodness-of-fit over AWB, but we are not confident that](#)
[the relative model parsimony of CON and other linear first-order models makes them universally more suitable for](#)
[predictive use.](#)

858 [Consequently, future SBM comparisons would benefit from additional data collection efforts sourced from](#)
860 [long-term ecological research experiments to globally verify the strengths and limitations of linear versus non-linear](#)
862 [SBMs, including CON and AWB, in Earth system modeling. The limited number of longitudinal soil warming](#)
864 [studies presents a challenge for facilitating site-specific model comparisons. We addressed this issue by using meta-](#)
866 [analysis data to aggregate warming responses across sites, but this approach does not provide site-specific](#)
[parameters. Additional data from ongoing and future field warming studies in the vein of the Harvard Forest and](#)
[Tropical Responses to Altered Climate experiments that demonstrate more varied flux dynamics over time than the](#)
[meta-analysis data set will be of critical importance for model testing \(Melillo et al., 2017; Wood et al., 2019\).](#)
[Model parameters could also be better constrained through the use of multivariate data sets, for example microbial](#)
[biomass dynamics in addition to soil respiration.](#)

868 [Our approach can be expanded to compare the predictive accuracies of linear microbial-implicit models to](#)
870 [those of recently developed non-linear microbial-explicit SBMs that are much larger than AWB, such as CORPSE](#)
872 [\(Sulman et al., 2014\) and MIMICS \(Wieder et al., 2014\). Such comparisons will help broadly determine if inclusion](#)
874 [of more detailed microbial dynamics in models offers predictive advantages that can overcome the overfitting](#)
876 [burdens associated with an increase in parameter count. With the appropriate data sets, our approach can also be](#)
878 [applied to consider the predictive performance of SBMs that describe the cycling of nitrogen \(N\), phosphorus \(P\),](#)
[and other limiting nutrients in addition to C dynamics. Models that represent N and P mineralization have yet to see](#)
[extensive head-to-head statistical benchmarking against C-only models with respect to predictive use \(Manzoni and](#)
[Porporato, 2009\). With models growing ever larger in size and specificity, there is a need to verify whether detailed](#)
[representation of microbial processes and the cycling of limiting nutrients are worth the increase in variable,](#)
[parameter, and equation counts. After all, “the tendency of more recent models towards more sophisticated \(and](#)
[generally more mathematically complex\) approaches is not always paralleled by improved model performance or](#)
[ability to interpret observed patterns” \(Manzoni and Porporato, 2009\).](#)

880 [The data assimilation and posterior sampling of complex models in future work comes with computing](#)
882 [performance challenges. Markov chain Monte Carlo algorithms are effective for exploring multidimensional](#)
884 [parameter space but are limited by temporal and computational expense, particularly when it comes to fitting non-](#)
886 [linear differential equation models \(Calderhead et al., 2009; Nemeth and Fearnhead, 2019\). Time per Markov chain](#)
888 [iteration drastically increases with number of parameters and data points. In fact, the present speed limitations of the](#)
890 [family of HMC algorithms make it necessary to use a hybrid approach utilizing Monte Carlo and deep learning](#)
892 [algorithms for parameter estimation at a global scale; Monte Carlo fitting is used to constrain parameter estimates at](#)
[a site-based scale before those estimates are tuned globally by deep learning using spatial information derived from](#)
[satellite maps \(Tao et al., 2020\). However, Monte Carlo algorithms are still the optimal methods for posterior](#)
[computation \(Duan et al., 2018\) and are necessary for Bayesian model comparisons conditional on site-based data.](#)
[Consequently, recent Monte Carlo algorithm innovations and developments that offer theoretical speed](#)
[improvements by trading thorough posterior sampling for numerical efficiency have been encouraging and are ripe](#)
[to be tested in future SBM comparisons involving more complex models and larger data sets. These developments](#)

Deleted: .

Formatted: English (US)

Deleted: Our approach can also be used to compare the predictive accuracy of linear models that only implicitly represent microbial activity to that of more complex non-linear SBMs that explicitly represent the Michaelis-Menten dynamics of soil microbial processes, such as CORPSE (Sulman et al., 2014) and MIMICS (Wieder et al., 2015). Such comparisons will help determine if inclusion of more detailed microbial dynamics in models offers predictive advantages that can overcome the overfitting burdens associated with an increase in parameter count.¶
Despite limited data availability, the

906 include stochastic gradient Monte Carlo sampling methods, a class of techniques in which a posterior is
907 approximated by fitting to a small subset of data at each iteration rather than estimated through exhaustive sampling
908 (Ma et al., 2015), and Gaussian process acceleration, in which a smooth distribution of likely solutions for a
909 differential equation system is specified and sampled in place of explicitly solving for the state variables during
910 every Markov chain iteration (Dondelinger et al., 2013; Wang and Barber, 2014).

911 Alongside advances in Monte Carlo algorithms, additional developments in Bayesian cross-validation and
912 information criteria measures are also available for practical trialing in soil biogeochemical data assimilation.
913 Gelman et al. have proposed a stable Bayesian counterpart of frequentist R^2 defined as “the variance of the predicted
914 values divided by the variance of predicted values plus the expected variance of the errors” that allows for more
915 intuitive and direct comparison to R^2 (Gelman et al., 2019). A Bayesian R^2 distribution provides a signal about the
916 absolute rather than relative goodness-of-fit of an associated posterior predictive distribution to the data. Bürkner et
917 al. (2019) have proposed a leave-future-out (LFO) cross-validation metric which is formulated to estimate relative
918 model predictive accuracy for hypothetical time series data occurring after existing experiment observations. LFO
919 and LOO are computed similarly, and LOO can also be used for time series data, as we demonstrated in this study.
920 However, the algorithmic differences between LFO and LOO make them better suited for different goals. LOO does
921 not inform about the quality of model fits for hypothetical samples collected after final reported measurements and
922 is more appropriate for estimating out-of-sample model predictive accuracy for hypothetical data samples taken
923 between the interval of observed measurement times (Vehtari et al., 2017).

924 The development of our formalized, statistically rigorous approach for model comparison and evaluation is
925 a critical step toward the goal of projecting global SOC levels and soil emissions throughout the 21st century. Our
926 initial results indicate promise in continued refinement and expansion of our approach to evaluate the predictive
927 performance of linear and non-linear SBMs. The future integration of updated Markov chain algorithms and
928 Bayesian predictive accuracy metrics into our framework will expand the ability to efficiently and thoroughly
929 compare differential equation models, even if they vary widely in structure and complexity.

Code and Data Availability

930 The R scripts, Stan code, and respiration data set used for HMC model fitting along with the original soil respiration
931 meta-analysis data set (Romero-Olivares et al., 2017) are available from the directory located at
932 https://osf.io/7mey8/?view_only=af1d54f858c34e41ab4854551d015896 (Xie et al., 2020).

Author contribution

934 SDA and HWX designed the study with assistance from MG. HWX and ALR performed the data cleaning and
935 analysis. HWX wrote the necessary code for the study with assistance from SDA. SDA and HWX prepared the
936 paper with suggestions from MG.

Competing interests

938 The authors declare they have no conflict of interest.

Acknowledgments

940 We would like to thank Stan development team members Aki Vehtari (Aalto University), Michael Betancourt, Bob
941 Carpenter (Flatiron Institute), Ben Bales (Columbia University), Charles Margossian (Columbia University), and
942 Sebastian Weber (Novartis) for their patient help with Stan code implementation and troubleshooting. We would
943 also like to thank both anonymous reviewers for their valuable and constructive comments, which not only aided in
944 the revision of the manuscript but also provided valuable insights to guide future work.

Financial support

946 This research was supported by funds from the National Science Foundation under grant DEB-1900885, the U.S.
947 Department of Energy Office of Science BER-TES program under grant DESC0014374, and the National Institutes
948 of Health T32 Training Program under grant EB009418.

Deleted: improving the forecasting of

Deleted: through the rest of

Deleted: development

Deleted: better

Deleted: a range of

Deleted: that

Deleted: parameter count

Deleted: (Xie et al., 2019).

Deleted:)

Deleted: Michael Betancourt

Deleted: The

960 **References**

- 962 Allison, S. D., Wallenstein, M. D. and Bradford, M. A.: Soil-carbon response to warming dependent on microbial
physiology, *Nat. Geosci.*, 3(5), 336–340, doi:10.1038/ngeo846, 2010.
- Anderson, T.-H. and Domsch, K. H.: Ratios of microbial biomass carbon to total organic carbon in arable soils, *Soil
964 Biol. Biochem.*, 21(4), 471–479, doi:10.1016/0038-0717(89)90117-X, 1989.
- Beskos, A., Pillai, N., Roberts, G., Sanz-Serna, J. M. and Stuart, A.: Optimal tuning of the hybrid Monte Carlo
966 algorithm, *Bernoulli*, 19(5 A), 1501–1534, doi:10.3150/12-BEJ414, 2013.
- Betancourt, M.: Diagnosing Suboptimal Cotangent Disintegrations in Hamiltonian Monte Carlo, [arXiv e-prints,
968 arXiv:1604.00695](#) [online] Available from: <http://arxiv.org/abs/1604.00695>, 2016.
- Betancourt, M.: A Conceptual Introduction to Hamiltonian Monte Carlo, [arXiv e-prints, arXiv:1701.02434](#) [online]
970 Available from: <http://arxiv.org/abs/1701.02434>, 2017.
- Betancourt, M. and Girolami, M.: Hamiltonian Monte Carlo for Hierarchical Models, *Curr. Trends Bayesian
972 Methodol. with Appl.*, 79–101, doi:10.1201/b18502-5, 2015.
- [Bradford, M. A. and Crowther, T. W.: Carbon use efficiency and storage in terrestrial ecosystems, *New Phytol.*,
974 199\(1\), 7–9, doi:10.1111/nph.12334, 2013.](#)
- [Bürkner, P.-C., Gabry, J. and Vehtari, A.: Approximate leave-future-out cross-validation for Bayesian time series
976 models, arXiv e-prints, arXiv:1902.06281](#) [online] Available from:
<https://ui.adsabs.harvard.edu/abs/2019arXiv190206281B>, 2019.
- [Calderhead, B., Girolami, M. and Lawrence, N. D.: Accelerating Bayesian Inference over Nonlinear Differential
978 Equations with Gaussian Processes, in *Advances in Neural Information Processing Systems 21*, edited by D. Koller,
D. Schuurmans, Y. Bengio, and L. Bottou, pp. 217–224, Curran Associates, Inc. \[online\] Available from:
980 \[http://papers.nips.cc/paper/3497-accelerating-bayesian-inference-over-nonlinear-differential-equations-with-
982 gaussian-processes.pdf\]\(http://papers.nips.cc/paper/3497-accelerating-bayesian-inference-over-nonlinear-differential-equations-with-gaussian-processes.pdf\), 2009.](#)
- Carpenter, B., Gelman, A., Hoffman, M. D., Lee, D., Goodrich, B., Betancourt, M., Brubaker, M. A., Guo, J., Li, P.
984 and Riddell, A.: Stan: A probabilistic programming language, *J. Stat. Softw.*, 76(1), doi:10.18637/jss.v076.i01,
2017.
- Christensen, R., Johnson, W., Branscum, A. and Hanson, T. E.: *Bayesian Ideas and Data Analysis: An Introduction
986 for Scientists and Statisticians*, 1st ed., CRC Press., 2010.
- 988 Crowther, T. W., Todd-Brown, K. E. O., Rowe, C. W., Wieder, W. R., Carey, J. C., MacHmuller, M. B., Snoek, B.
L., Fang, S., Zhou, G., Allison, S. D., Blair, J. M., Bridgham, S. D., Burton, A. J., Carrillo, Y., Reich, P. B., Clark, J.
990 S., Classen, A. T., Dijkstra, F. A., Elberling, B., Emmett, B. A., Estiarte, M., Frey, S. D., Guo, J., Harte, J., Jiang, L.,
Johnson, B. R., Kroël-Dulay, G., Larsen, K. S., Laudon, H., Lavallee, J. M., Luo, Y., Lupascu, M., Ma, L. N.,
992 Marhan, S., Michelsen, A., Mohan, J., Niu, S., Pendall, E., Peñuelas, J., Pfeifer-Meister, L., Poll, C., Reinsch, S.,

Deleted: [online] Available from:
papers2://publication/uuid/90F0F273-FB35-4093-B8BB-
67E2429E94BF...

Deleted: ., 2011

998 Reynolds, L. L., Schmidt, I. K., Sistla, S., Sokol, N. W., Templer, P. H., Treseder, K. K., Welker, J. M. and
Bradford, M. A.: Quantifying global soil carbon losses in response to warming, *Nature*, 540(7631), 104–108,
doi:10.1038/nature20150, 2016.

1000 [Curtiss, C. F. and Hirschfelder, J. O.: Integration of Stiff Equations, *Proc. Natl. Acad. Sci. U. S. A.*, 38\(3\), 235–243, doi:10.1073/pnas.38.3.235, 1952.](#)

1002 [Dondelinger, F., Husmeier, D., Rogers, S. and Filippone, M.: ODE parameter inference using adaptive gradient
1004 matching with Gaussian processes, in *Proceedings of the Sixteenth International Conference on Artificial
Intelligence and Statistics*, vol. 31, edited by C. M. Carvalho and P. Ravikumar, pp. 216–228, PMLR, Scottsdale,
Arizona, USA. \[online\] Available from: <http://proceedings.mlr.press/v31/dondelinger13a.html>, 2013.](#)

1006 [Duan, L. L., Johndrow, J. E. and Dunson, D. B.: Scaling up Data Augmentation MCMC via Calibration, *J. Mach.
Learn. Res.*, 19\(1\), 2575–2608, 2018.](#)

1008 Fang, C. and Moncrieff, J. B.: The variation of soil microbial respiration with depth in relation to soil carbon
composition, *Plant Soil*, 268(1), 243–253, doi:10.1007/s11104-004-0278-4, 2005.

1010 Gelfand, A. E. and Dey, D. K.: Bayesian Model Choice : Asymptotics and Exact Calculations, *J. R. Stat. Soc. Ser.
B*, 56(3), 501–514, 1994.

1012 [Gelfand, A. E., Dey, D. K. and Chang, H.: Model determination using predictive distributions, with implementation
1014 via sampling-based methods \(with discussion\), in *Bayesian Statistics 4*, edited by J. M. Bernardo, J. O. Berger, A. P.
Dawid, and A. F. . Smith, pp. 147–167, Oxford University Press., 1992.](#)

Gelman, A., Hwang, J. and Vehtari, A.: Understanding predictive information criteria for Bayesian models, *Stat.
1016 Comput.*, 24(6), 997–1016, doi:10.1007/s11222-013-9416-2, 2014.

1018 [Gelman, A., Goodrich, B., Gabry, J. and Vehtari, A.: R-squared for Bayesian Regression Models, *Am. Stat.*, 73\(3\),
307–309, doi:10.1080/00031305.2018.1549100, 2019.](#)

van Gestel, N., Shi, Z., van Groenigen, K. J., Osenberg, C. W., Andresen, L. C., Dukes, J. S., Hovenden, M. J., Luo,
1020 Y., Michelsen, A., Pendall, E., Reich, P. B., Schuur, E. A. G. and Hungate, B. A.: Predicting soil carbon loss with
warming, *Nature*, 554(7693), E4–E5, doi:10.1038/nature25745, 2018.

1022 [Geyer, K. M., Dijkstra, P., Sinsabaugh, R. and Frey, S. D.: Clarifying the interpretation of carbon use efficiency in
1024 soil through methods comparison, *Soil Biol. Biochem.*, 128, 79–88,
doi:<https://doi.org/10.1016/j.soilbio.2018.09.036>, 2019.](#)

Guo, J., Gabry, J. and Goodrich, B.: RStan: the R interface to Stan, 2019.

1026 [Hagerty, S. B., Allison, S. D. and Schimel, J. P.: Evaluating soil microbial carbon use efficiency explicitly as a
1028 function of cellular processes: implications for measurements and models, *Biogeochemistry*, 140\(3\), 269–283,
doi:10.1007/s10533-018-0489-z, 2018.](#)

[Hale, J. K. and LaSalle, J. P.: Differential Equations: Linearity vs. Nonlinearity, *SIAM Rev.*, 5\(3\), 249–272,](#)

Deleted: Author (s): A . E . Gelfand and D . K . Dey
Published by : Blackwell Publishing for the Royal Statistical
Society Stable URL : <http://www.jstor.org/stable/2346123>, J.

Deleted: Gelman, A.: Conservative prior distributions for
variance parameters in hierarchical models, *Bayesian Anal.*,
1(3), 515–533, doi:10.1002/ejs.5550340302, 2006.¶

1036 [doi:10.1137/1005068.1963](https://doi.org/10.1137/1005068.1963).

1038 [Hararuk, O. and Luo, Y.: Improvement of global litter turnover rate predictions using a Bayesian MCMC approach, *Ecosphere*, 5\(12\), art163, doi:10.1890/ES14-00092.1, 2014.](#)

1040 [Hararuk, O., Xia, J. and Luo, Y.: Evaluation and improvement of a global land model against soil carbon data using a Bayesian Markov chain Monte Carlo method, *J. Geophys. Res. Biogeosciences*, 119\(3\), 403–417, doi:10.1002/2013JG002535, 2014.](#)

1042 [Hararuk, O., Zwart, J. A., Jones, S. E., Prairie, Y. and Solomon, C. T.: Model-Data Fusion to Test Hypothesized Drivers of Lake Carbon Cycling Reveals Importance of Physical Controls, *J. Geophys. Res. Biogeosciences*, 123\(3\), 1130–1142, doi:10.1002/2017JG004084, 2018.](#)

1046 [Hindmarsh, A. C., Brown, P. N., Grant, K. E., Lee, S. L., Serban, R., Shumaker, D. E. and Woodward, C. S.: SUNDIALS: Suite of nonlinear and differential/algebraic equation solvers, *ACM Trans. Math. Softw.*, 31\(3\), 363–396, doi:10.1145/1089014.1089020, 2005.](#)

1048 [Ibrahim, J. G., Chen, M.-H. and Sinha, D.: Bayesian Survival Analysis, 1st ed., Springer-Verlag New York, New York City, New York., 2001.](#)

1050 Jiang, L., Yan, Y., Hararuk, O., Mikle, N., Xia, J., Shi, Z., Tjiputra, J., Wu, T. and Luo, Y.: Scale-dependent performance of CMIP5 earth system models in simulating terrestrial vegetation carbon, *J. Clim.*, 28(13), 5217–5232, doi:10.1175/JCLI-D-14-00270.1, 2015.

1052 Jobbágy, E. and Jackson, R. B.: The Vertical Distribution of Soil Organic Carbon and Its Relation to Climate and Vegetation, *Ecol. Appl.*, 10(April), 423–436, doi:Doi 10.2307/2641104, 2000.

1054 [Kvålseth, T. O.: Cautionary Note about R2, *Am. Stat.*, 39\(4\), 279–285, doi:10.1080/00031305.1985.10479448.1985.](#)

1056 Li, J., Wang, G., Allison, S. D., Mayes, M. A. and Luo, Y.: Soil carbon sensitivity to temperature and carbon use efficiency compared across microbial-ecosystem models of varying complexity, *Biogeochemistry*, 119(1–3), 67–84, doi:10.1007/s10533-013-9948-8, 2014.

1060 Luo, Y., [Ahlström, A.](#), Allison, S. D., Batjes, N. H., Brovkin, V., Carvalhais, N., Chappell, A., Ciais, P., Davidson, E. A., Finzi, A., Georgiou, K., Guenet, B., Hararuk, O., Harden, J. W., He, Y., Hopkins, F., Jiang, L., Koven, C., Jackson, R. B., Jones, C. D., Lara, M. J., Liang, J., McGuire, A. D., Parton, W., Peng, C., Randerson, J. T., Salazar, A., Sierra, C. A., Smith, M. J., Tian, H., Todd-Brown, K. E. O., Torn, M., van Groenigen, K. J., Wang, Y. P., West, T. O., Wei, Y., Wieder, W. R., Xia, J., Xu, X., Xu, X. and Zhou, T.: [Toward](#) more realistic projections of soil carbon dynamics by [Earth](#) system models, *Global Biogeochem. Cycles*, 30(1), 40–56, doi:10.1002/2015GB005239, 2016.

1062 [Ma, Y.-A., Chen, T. and Fox, E. B.: A Complete Recipe for Stochastic Gradient MCMC, in *Proceedings of the 28th International Conference on Neural Information Processing Systems - Volume 2*, pp. 2917–2925, MIT Press, Cambridge, MA, USA., 2015.](#)

1068

Deleted: Ahlström

Deleted: . T

Deleted: . W

Deleted: Towards

Deleted: earth

Deleted: , (February

Deleted: .Received

1076 Manzoni, S. and Porporato, A.: Soil carbon and nitrogen mineralization: Theory and models across scales, *Soil Biol. Biochem.*, 41(7), 1355–1379, doi:10.1016/j.soilbio.2009.02.031, 2009.

1078 Melillo, J. M., Frey, S. D., DeAngelis, K. M., Werner, W. J., Bernard, M. J., Bowles, F. P., Pold, G., Knorr, M. A. and Grandy, A. S.: Long-term pattern and magnitude of soil carbon feedback to the climate system in a warming world, *Science* (80-.), 358(6359), 101–105, doi:10.1126/science.aan2874, 2017.

1080 [Nemeth, C. and Fearnhead, P.: Stochastic gradient Markov chain Monte Carlo, arXiv e-prints, arXiv:1907.06986 \[online\] Available from: <https://ui.adsabs.harvard.edu/abs/2019arXiv190706986N>, 2019.](#)

1082 [Nottingham, A. T., Turner, B. L., Whitaker, J., Ostle, N., Bardgett, R. D., McNamara, N. P., Salinas, N. and Meir, P.: Temperature sensitivity of soil enzymes along an elevation gradient in the Peruvian Andes, *Biogeochemistry*, 127\(2\), 217–230, doi:10.1007/s10533-015-0176-2, 2016.](#)

1084 [R Core Team: R: A Language and Environment for Statistical Computing, \[online\] Available from: <http://www.r-project.org>, 2017.](#)

1086 Romero-Olivares, A. L., Allison, S. D. and Treseder, K. K.: Soil microbes and their response to experimental warming over time: A meta-analysis of field studies, *Soil Biol. Biochem.*, 107, 32–40, doi:10.1016/j.soilbio.2016.12.026, 2017.

1088 Sparling, G. P.: Ratio of microbial biomass carbon to soil organic carbon as a sensitive indicator of changes in soil organic matter, *Aust. J. Soil Res.*, 30(2), 195–207, doi:10.1071/SR9920195, 1992.

1092 Spiess, A. N. and Neumeyer, N.: An evaluation of R^2 as an inadequate measure for nonlinear models in pharmacological and biochemical research: A Monte Carlo approach, *BMC Pharmacol.*, 10(1), 6, doi:10.1186/1471-2210-10-6, 2010.

1094 [Steinweg, J. M., Jagadamma, S., Frerichs, J. and Mayes, M. A.: Activation Energy of Extracellular Enzymes in Soils from Different Biomes, *PLoS One*, 8\(3\), 1–7, doi:10.1371/journal.pone.0059943, 2013.](#)

1096 Sulman, B. N., Phillips, R. P., Oishi, A. C., Shevliakova, E. and Pacala, S. W.: Microbe-driven turnover offsets mineral-mediated storage of soil carbon under elevated CO_2 , *Nat. Clim. Chang.*, 4(12), 1099–1102, doi:10.1038/nclimate2436, 2014.

1098 Sulman, B. N., Moore, J. A. M., Abramoff, R., Averill, C., Kivlin, S., Georgiou, K., Sridhar, B., Hartman, M. D., Wang, G., Wieder, W. R., Bradford, M. A., Luo, Y., Mayes, M. A., Morrison, E., Riley, W. J., Salazar, A., Schimel, J. P., Tang, J. and Classen, A. T.: Multiple models and experiments underscore large uncertainty in soil carbon dynamics, *Biogeochemistry*, 141(2), 109–123, doi:10.1007/s10533-018-0509-z, 2018.

1104 [Tang, J. Y.: On the relationships between the Michaelis–Menten kinetics, reverse Michaelis–Menten kinetics, equilibrium chemistry approximation kinetics, and quadratic kinetics, *Geosci. Model Dev.*, 8\(12\), 3823–3835, doi:10.5194/gmd-8-3823-2015, 2015.](#)

1106 [Tao, F., Zhou, Z., Huang, Y., Li, Q., Lu, X., Ma, S., Huang, X., Liang, Y., Hugelius, G., Jiang, L., Doughty, R.,](#)

1108

Deleted: R2as

Deleted: ,

Deleted: -11

Deleted: CO 2

1114 [Ren, Z. and Luo, Y.: Deep Learning Optimizes Data-Driven Representation of Soil Organic Carbon in Earth System Model Over the Conterminous United States, *Front. Big Data*, 3, 1–17, doi:10.3389/fdata.2020.00017, 2020.](#)

1116 Todd-Brown, K. E. O., Randerson, J. T., Hopkins, F., Arora, V., Hajima, T., Jones, C., Shevliakova, E., Tjiputra, J., Volodin, E., Wu, T., Zhang, Q. and Allison, S. D.: Changes in soil organic carbon storage predicted by Earth system models during the 21st century, *Biogeosciences*, 11(8), 2341–2356, doi:10.5194/bg-11-2341-2014, 2014.

1118 [Trasar-Cepeda, C., Gil-Sotres, F. and Leirós, M. C.: Thermodynamic parameters of enzymes in grassland soils from Galicia, NW Spain, *Soil Biol. Biochem.*, 39\(1\), 311–319, doi:https://doi.org/10.1016/j.soilbio.2006.08.002, 2007.](#)

1120 Trumbore, S.: Age of soil organic matter and soil respiration: Radiocarbon constraints on belowground C dynamics, *Ecol. Appl.*, 10(2), 399–411, doi:10.1890/1051-0761(2000)010[0399:AOSOMA]2.0.CO;2, 2000.

1122 [Vehuri, A. and Ojanen, J.: A survey of Bayesian predictive methods for model assessment, selection and comparison, *Stat. Surv.*, 6\(1\), 142–228, doi:10.1214/12-ss102, 2012.](#)

1124 [Vehuri, A., Gelman, A. and Gabry, J.: Practical Bayesian model evaluation using leave-one-out cross-validation and WAIC, *Stat. Comput.*, 27\(5\), 1413–1432, doi:10.1007/s11222-016-9696-4, 2017.](#)

1126 [Vehuri, A., Gabry, J., Magnusson, M., Yao, Y. and Gelman, A.: loo: Efficient leave-one-out cross-validation and WAIC for Bayesian models, \[online\] Available from: <https://mc-stan.org/loo>, 2019.](#)

1128 [Wang, B. and Allison, S. D.: Emergent properties of organic matter decomposition by soil enzymes, *Soil Biol. Biochem.*, 136, 107522, doi:https://doi.org/10.1016/j.soilbio.2019.107522, 2019.](#)

1130 [Wang, Y. and Barber, D.: Gaussian Processes for Bayesian Estimation in Ordinary Differential Equations, in *Proceedings of the 31st International Conference on International Conference on Machine Learning - Volume 32*, pp. 1485–1493, JMLR.org, Beijing, China., 2014.](#)

1132 [Wieder, W. R., Grandy, A. S., Kallenbach, C. M. and Bonan, G. B.: Integrating microbial physiology and physio-chemical principles in soils with the Microbial-Mineral Carbon Stabilization \(MIMICS\) model, *Biogeosciences*, 11\(14\), 3899–3917, doi:10.5194/bg-11-3899-2014, 2014.](#)

1134 Wieder, W. R., Allison, S. D., Davidson, E. A., Georgiou, K., Hararuk, O., He, Y., Hopkins, F., Luo, Y., Smith, M. J., Sulman, B., Todd-Brown, K., Wang, Y. P., Xia, J. and Xu, X.: Explicitly representing soil microbial processes in Earth system models, *Global Biogeochem. Cycles*, 29(10), 1782–1800, doi:10.1002/2015GB005188, [2015a](#).

1138 [Wieder, W. R., Grandy, A. S., Kallenbach, C. M., Taylor, P. G. and Bonan, G. B.: Representing life in the Earth system with soil microbial functional traits in the MIMICS model, *Geosci. Model Dev.*, 8\(6\), 1789–1808, doi:10.5194/gmd-8-1789-2015, \[2015b\]\(#\).](#)

1140 Wieder, W. R., Hartman, M. D., Sulman, B. N., Wang, Y. P., Koven, C. D. and Bonan, G. B.: Carbon cycle confidence and uncertainty: Exploring variation among soil biogeochemical models, *Glob. Chang. Biol.*, 24(4), 1563–1579, doi:10.1111/gcb.13979, 2018.

1142 Wood, T. E., González, G., Silver, W. L., Reed, S. C. and Cavaleri, M. A.: On the shoulders of giants: Continuing

Deleted: 2015

1148 the legacy of large-scale ecosystem manipulation experiments in Puerto Rico, *Forests*, 10(3), 1–18,
doi:10.3390/f10030210, 2019.

1150 Xie, H. W., Romero-Olivares, A. L., Treseder, K. K. and Allison, S. D.: A Bayesian Approach to Evaluation of Soil
Biogeochemical Models R [And Stan Code, \[online\] Available from:](#)
https://osf.io/7mey8/?view_only=af1d54f858c34c41ab4854551d015896. 2020.

1152 Zhang, B., Chen, S., He, X., Liu, W., Zhao, Q., Zhao, L. and Tian, C.: Responses of soil microbial communities to
experimental warming in alpine grasslands on the Qinghai-Tibet Plateau, *PLoS One*, 9(8),
1154 doi:10.1371/journal.pone.0103859, 2014.

1156

1158

1160

1162

1164

1166

1168

1170

1172

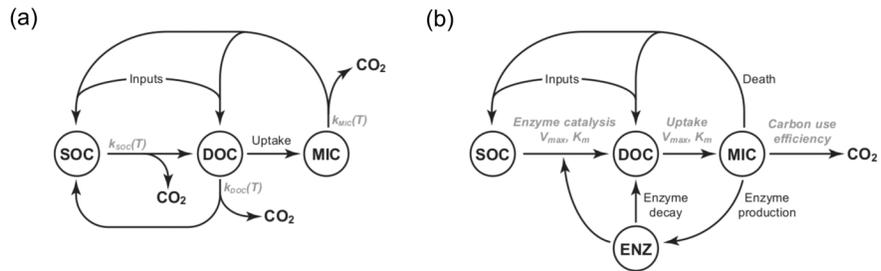
1174

Deleted: and Stan Code, doi:10.17605/OSF.IO/7MEY8, 2019...

Formatted: Don't adjust right indent when grid is defined, Space Before: 0 pt, After: 6 pt, Line spacing: 1.5 lines, No widow/orphan control, Don't adjust space between Latin and Asian text, Don't adjust space between Asian text and numbers

Formatted: Font: Not Bold

Formatted: None, Don't adjust right indent when grid is defined, Space Before: 0 pt, After: 6 pt, Line spacing: 1.5 lines, No widow/orphan control, Don't adjust space between Latin and Asian text, Don't adjust space between Asian text and numbers



1178 **Figure 1:** Diagrams of the pool structures of the (a) CON model; and (b) AWB model, drawn from Allison et al.,
 1180 (2010). Pools are shown within circles including soil organic carbon (SOC), dissolved organic carbon (DOC), and
 1182 microbial (MIC) pools. AWB has SOC, DOC, and MIC pools as in CON, but also an extra enzymatic (ENZ) pool.
 AWB additionally differs from CON in its non-linear feedbacks and assumption that MIC can influence SOC-to-
 DOC turnover through the ENZ pool.

Deleted: .

1184

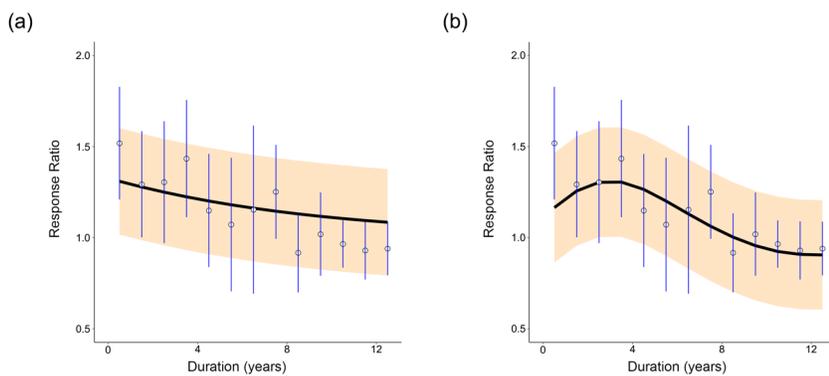
1186

1188

1190

1192

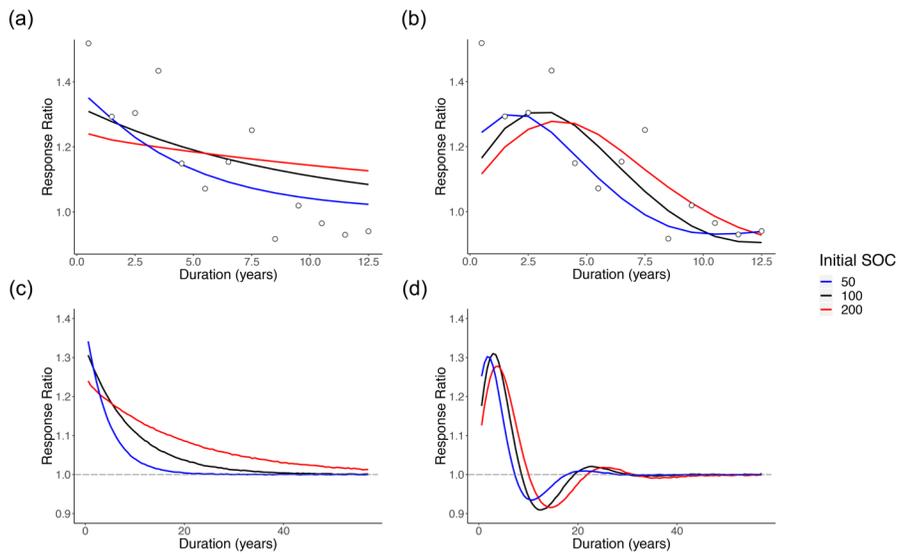
1194



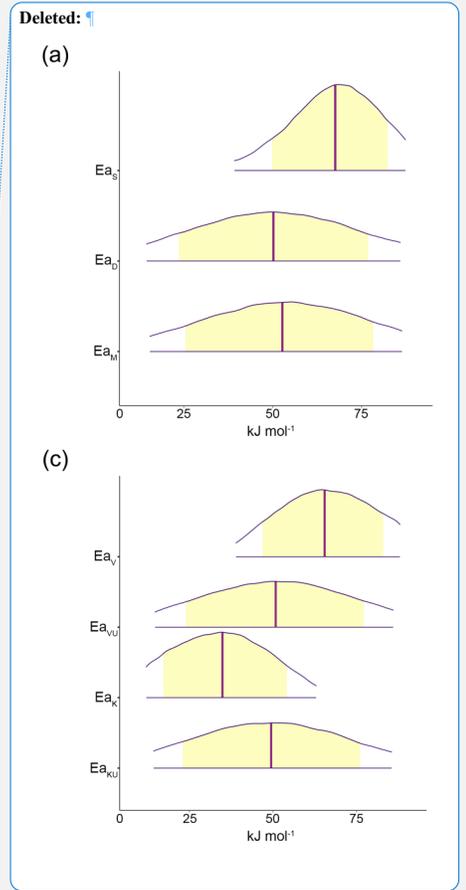
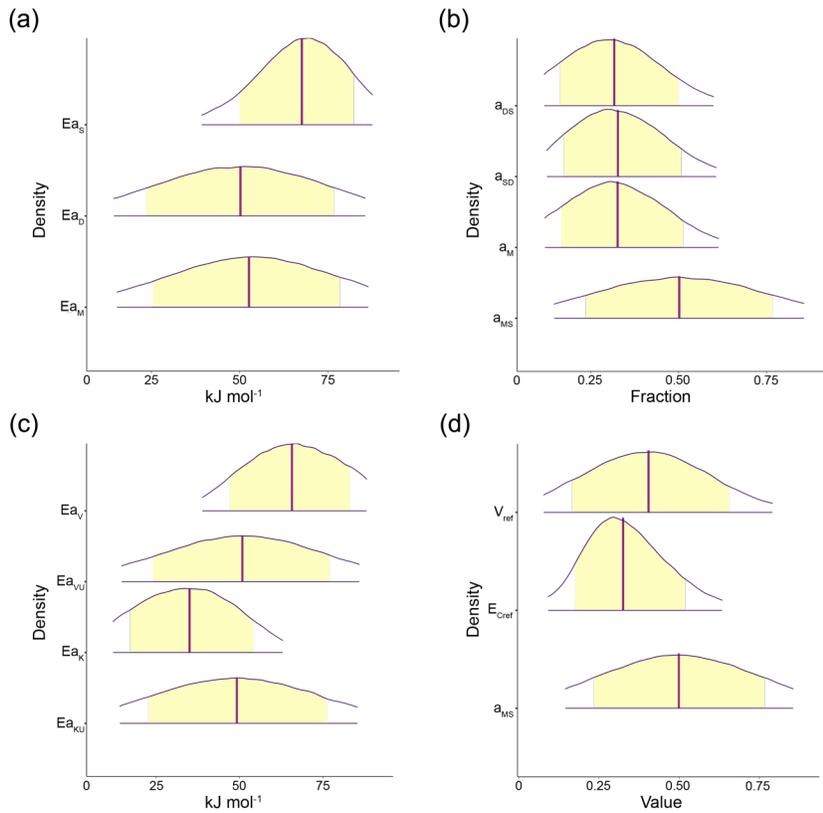
1196
 1198
 1200
 1202
 1204
 1206
 1208

Figure 2: Distribution of fits of **(a)** CON; and **(b)** AWB to the meta-analysis data from Romero-Olivares et al., (2017). Open circles show the meta-analysis data points. Blue vertical lines mark the 95% confidence interval for each data point calculated from the pooled standard deviation. The black line indicates the mean model response ratio fit. The orange shading marks the 95% posterior predictive interval for the fit. For **(a)**, pre-warming steady state soil C densities were set at SOC = 100 mg C g⁻¹ soil, MIC = 2 mg C g⁻¹ soil, DOC = 0.2 mg C g⁻¹ soil. For **(b)**, pre-warming steady state soil C densities were set at SOC = 100 mg C g⁻¹ soil, MIC = 2 mg C g⁻¹ soil, DOC = 0.2 mg C g⁻¹ soil, and ENZ = 0.1 mg C g⁻¹ soil.

Deleted: .
 Deleted: (and median)



1212 **Figure 3:** Intra-model comparisons of mean posterior predictive response ratio fits for AWB and CON across
 1214 different MIC-to-SOC ratios. Open circles show the meta-analysis data points for reference. The blue, black, and red
 1216 lines indicate model mean fits corresponding to different pre-warming-perturbation steady state SOC values of 50
 mg C g⁻¹ soil, 100 mg C g⁻¹ soil, and 200 mg C g⁻¹ soil. The dashed gray line indicates the steady state expectation at
 the response ratio of 1.0. Mean fits are plotted in order of (a) CON; and (b) AWB over the time span of the data and
 (c) CON; and (d) AWB over 57 years.



1218
1220
1222
1224
1226
1228

Figure 4: 95% probability density credible areas for model parameters corresponding to pre-warming steady state SOC = 100 mg C g⁻¹ soil, DOC = 0.2 mg C g⁻¹ soil, MIC = 2 mg C g⁻¹ soil, and (for AWB) ENZ = 0.1 mg C g⁻¹ soil. Yellow shaded regions represent 80% credible areas and vertical purple lines indicate distribution mean. **(a)** CON activation energy parameters E_{a_s} , E_{a_D} , and E_{a_M} ; **(b)** CON C pool partition fraction parameters a_{DS} , a_{SD} , a_M , and a_{MS} ; **(c)** AWB activation energy parameters E_{a_v} , $E_{a_{vU}}$, E_{a_K} , and $E_{a_{KU}}$; **(d)** AWB parameters V_{ref} , $E_{C_{ref}}$, and a_{MS} . V_{ref} is the SOC V_{max} at the reference temperature 283.15 K, $E_{C_{ref}}$ is the carbon use efficiency fraction at the reference temperature, and like its CON counterpart, the AWB a_{MS} parameter is the fraction parameter representing the proportion of dead microbial biomass C transferred to the SOC pool. Parameter units are displayed in Supplemental Table 1. Credible areas for AWB parameters $V_{U_{ref}}$ and m_t are shown in Supplemental Fig. 2 because of differing horizontal axes scales.

- Deleted:** E_{a_s} , E_{a_D} , E_{a_M} ;
- Deleted:** a_{DS} , a_{SD} , a_M , and a_{MS} ;
- Deleted:** E_{a_v} , $E_{a_{vU}}$, E_{a_K} , $E_{a_{KU}}$;
- Deleted:** V_{ref} , $E_{C_{ref}}$, and a_{MS} . V_{ref}
- Deleted:** $E_{C_{ref}}$
- Deleted:** a_{MS}
- Formatted:** Font color: Black, Not Superscript/ Subscript
- Deleted:** $V_{U_{ref}}$
- Deleted:** m_t
- Formatted:** Font color: Black

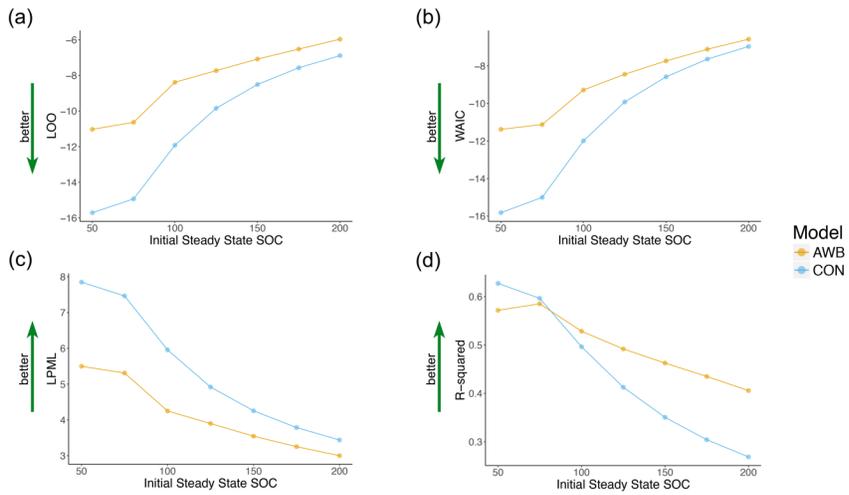


Figure 5: Goodness-of-fit metrics plotted against initial steady state SOC for AWB and CON models for (a) LOO; (b) WAIC cross-validation; (c) LPML; and (d) R^2 values. Pre-perturbation steady state MIC, DOC, and ENZ (for AWB) is held constant as pre-perturbation SOC is varied.

Deleted:
 Deleted: Fit metric versus
 Deleted:),

1240
1242

Page 3: [1] Deleted	Updated	7/2/20 8:29:00 AM
Page 5: [2] Deleted	Updated	7/2/20 8:29:00 AM
Page 6: [3] Deleted	Updated	7/2/20 8:29:00 AM
Page 6: [4] Deleted	Updated	7/2/20 8:29:00 AM
Page 6: [5] Deleted	Updated	7/2/20 8:29:00 AM
Page 6: [6] Deleted	Updated	7/2/20 8:29:00 AM
Page 6: [7] Deleted	Updated	7/2/20 8:29:00 AM
Page 6: [8] Deleted	Updated	7/2/20 8:29:00 AM
Page 7: [9] Deleted	Updated	7/2/20 8:29:00 AM