A framework for variational inference and data assimilation of soil biogeochemical models using state space approximations and normalizing flows

Hua Wally Xie^{1,1,1,1,1}, Debora Sujono^{1,1,1,1,1}, Tom Ryder^{2,2,2,2,2}, Erik B. Sudderth^{1,1,1,1,1}, and Steven Allison^{1,1,1,1,1}

¹UC Irvine ²Newcastle University

November 30, 2022

Abstract

Soil biogeochemical models (SBMs) simulate element transfer processes between organic soil pools. These models can be used to specify falsifiable quantitative assertions about soil system dynamics and their responses to global surface temperature warming. To determine whether SBMs are useful for representing and forecasting data-generating processes in soils, it is important to conduct data assimilation and fitting of SBMs conditioned on soil pool and flux measurements to validate model predictive accuracy. SBM data assimilation has previously been carried out in approaches ranging from visual qualitative tuning of model output against data to more statistically rigorous Bayesian inferences that estimate posterior parameter distributions with Markov chain Monte Carlo (MCMC) methods. MCMC inference is better able to account for data and parameter uncertainty, but the computational inefficiency of MCMC methods limits their ability to scale assimilations to larger data sets. With formulation of efficient and statistically rigorous SBM inference frameworks remaining an open problem, we demonstrate the novel application of a variational inference framework that uses a method called normalizing flows to approximate SBMs that have been discretized into state space models. We fit the approximated SBMs to synthetic data sourced from known data-generating processes to identify discrepancies between the inference results and true parameter values and ensure functionality of our method. Our approach trades estimation accuracy for algorithmic efficiency gains that make SBM data assimilation more tractable and achievable under computational time and resource limitations.

A framework for variational inference and data assimilation of soil biogeochemical models using state space approximations and normalizing flows

H. W. Xie^{1,†}, D. Sujono^{2,†}, T. Ryder³, E. B. Sudderth², S. D. Allison^{1,4}

¹Center for Complex Biological Systems, University of California, Irvine, Irvine, CA, United States of

America ²Department of Computer Science, University of California, Irvine, Irvine, CA, United States of America ³School of Mathematics Statistics and Physics, Newcastle University, Newcastle, United Kingdom ⁴Department of Ecology and Evolutionary Biology, University of California, Irvine, Irvine, CA, United

States of America

[†]Authors contributed equally to this work.

Key Points:

1

2

3

4

5

10

11

12

18

13	•	Soil biogeochemical models (SBMs) represent the dynamics of soil systems and
14		predict their responses to global warming.
15	•	SBMs are fit to data sets including soil respiration and organic matter measure-
16		ments to assess model accuracy and estimate parameter values.
17	•	We demonstrate a Bayesian inference workflow in this study that can be used for

efficient SBM data fitting and parameter estimation.

Corresponding author: Hua Xie, xiehw@uci.edu

19 Abstract

Soil biogeochemical models (SBMs) simulate element transfer processes between organic
 soil pools. These models can be used to specify falsifiable quantitative assertions about
 soil system dynamics and their responses to global surface temperature warming.

To determine whether SBMs are useful for representing and forecasting data-generating 23 processes in soils, it is important to conduct data assimilation and fitting of SBMs con-24 ditioned on soil pool and flux measurements to validate model predictive accuracy. SBM 25 data assimilation has previously been carried out in approaches ranging from visual qual-26 27 itative tuning of model output against data to more statistically rigorous Bayesian inferences that estimate posterior parameter distributions with Markov chain Monte Carlo 28 (MCMC) methods. MCMC inference is better able to account for data and parameter 29 uncertainty, but the computational inefficiency of MCMC methods limits their ability 30 to scale assimilations to larger data sets. 31

With formulation of efficient and statistically rigorous SBM inference frameworks 32 remaining an open problem, we demonstrate the novel application of a variational in-33 ference framework that uses a method called normalizing flows to approximate SBMs 34 that have been discretized into state space models. We fit the approximated SBMs to 35 synthetic data sourced from known data-generating processes to identify discrepancies 36 between the inference results and true parameter values and ensure functionality of our 37 method. Our approach trades estimation accuracy for algorithmic efficiency gains that 38 make SBM data assimilation more tractable and achievable under computational time 39 and resource limitations. 40

⁴¹ Plain Language Summary

Soil biogeochemical models (SBMs) simulate soil systems in a quantifiable and fal-42 sifiable manner. Climate researchers rely on SBMs to predict how soil systems could be 43 globally affected by climate change. However, SBMs differ widely in their predictions of 44 changes in soil measurements including rates of soil carbon dioxide emissions. It is un-45 clear which SBMs offer more realistic climate projections, and the establishment of sta-46 tistical techniques to rigorously compare the predictive performance of SBMs is still a 47 work in progress. We make a contribution to SBM comparison efforts by developing a 48 statistical framework to assess SBM accuracy that leverages deep learning for compu-49 tational efficiency gains. Results of our case study demonstrate that we can fit two SBMs 50 to soil observation data and estimate ranges of SBM parameter values compatible with 51 52 those observations. Our modular framework is flexible and stimulates future work to improve on our procedure with modifications of our existing methods. 53

54 1 Introduction

Soil biogeochemical models (SBMs) are differential equation systems that repre-55 sent dynamics of organic matter transfer between soil pools, including the soil organic 56 (SOC), dissolved organic (DOC), and microbial biomass carbon (MBC) pools. The state 57 variables of SBMs typically are densities or masses of elements in those pools (Manzoni 58 & Porporato, 2009), and heterotrophic soil CO_2 emissions can be estimated from those 59 state values and microbial parameters (Allison et al., 2010). As soil microbe communi-60 ties influencing organic mass transfer dynamics evolve and shift under the selection pres-61 sures of terrestrial warming, SBMs have become an important tool for soil scientists and 62 biogeochemists to quantify changes in soil system activity and predict future heterotrophic 63 soil respiration levels (Sulman et al., 2018; Saifuddin et al., 2021). 64

SBMs offer falsifiability of their dynamics through their depiction of biological soil processes as interpretable mathematical equations governed by model parameters θ . How-

ever, the formulation of statistically sound frameworks to assess the dynamical validity 67 and predictive accuracy of SBMs remains an open problem in soil biogeochemistry (Luo 68 et al., 2016; Xie et al., 2020; Bradford et al., 2021; Georgiou et al., 2021; Raczka et al., 69 2021). One approach for assessing SBM utility involves comparing models by their abil-70 ity to assimilate soil observation data with suitable θ values, assuming that models that 71 can more accurately describe the past will also be better at predicting the future (Wieder 72 et al., 2014; Bradford et al., 2021). Past SBM fit evaluations have ranged from visual 73 juxtapositions of manually calibrated model outputs against empirical observations (Sulman 74 et al., 2014; Wieder et al., 2015) to quantitative frequentist comparisons involving cor-75 relation coefficients and root-mean-square errors (K. E. O. Todd-Brown et al., 2013; K. E. Todd-76 Brown et al., 2014; Wieder et al., 2014). In an effort to account for uncertainty in θ val-77 ues and data observations and encode expert domain beliefs, other comparisons have in-78 volved the use of Bayesian Markov chain Monte Carlo (MCMC) inference methods and 79 goodness-of-fit metrics with some success (Hararuk et al., 2014; Hararuk & Luo, 2014; 80 J. Li et al., 2019; Xie et al., 2020; Saifuddin et al., 2021; Wang et al., 2022). 81

MCMC transition sampling methods, such as the Gibbs (Geman & Geman, 1987) 82 and No-U-Turn (NUTS) (Hoffman & Gelman, 2014) samplers deployed in widely-used 83 probabilistic programming platforms like JAGS (Plummer, 2003), Stan (Carpenter et 84 al., 2017), and PyMC (Salvatier et al., 2016), are powerful algorithms for inference, but 85 their relative computational cost presently limits their ability to scale for use on model 86 comparisons involving more complex SBM systems conditioned on larger data sets span-87 ning decades (Kucukelbir et al., 2017). Stochastic gradient optimization variational in-88 ference (VI) is an alternative approach to Bayesian inference and model-fitting that trades 89 asymptotic exactness and the ability to estimate non-parametric posterior distributions 90 for increased computational efficiency and simplicity (Blei et al., 2017). It does so by re-91 framing Bayesian inference from a transition sampling problem to an optimization ob-92 jective of maximizing a metric called the evidence lower bound (ELBO), which corre-93 sponds to minimizing the discrepancy between an approximate parametric posterior and 94 true posterior distribution (Salimans et al., 2015). 95

VI on differential equation models benefits from the use of stochastic differential 96 equation (SDE) over ordinary differential equation (ODE) systems. SDE noise provides 97 a means of adjusting and correcting for proposals of system initial conditions and un-98 derlying dynamics that are inconsistent with the true data-generating process sourcing aq the data observations (Whitaker, 2016; Särkkä & Solin, 2019; Wigvist et al., 2021). Ad-100 ditionally, noise-driven fluctuation and variation in state trajectories can account for out-101 lier data measurements during inference to reduce optimization pressures that can drive 102 rigid deterministic models into unstable θ regimes. SDE noise thereby improves infer-103 ence flexibility, stability, and efficiency through the acommodation and mitigation of dis-104 crepancies between model outputs and data generation or observation. Furthermore, SDEs 105 offer a more realistic and accurate representation of the stochasticity that is inherent to 106 biological processes across all scales (Golightly & Wilkinson, 2011; Abs et al., 2020; Brown-107 ing et al., 2020). The ability to effectively fit SDEs is an advantage of VI over many es-108 tablished MCMC methods; off-the-shelf MCMC implementations are frequently not de-109 signed to tolerate the noisy likelihood estimates of SDEs (Golightly & Wilkinson, 2010; 110 Fuchs, 2013; Chen et al., 2014). 111

With the goal of applying VI to SBMs in mind, we formulated SDE versions of the linear deterministic "conventional" (CON) SBM system (Allison et al., 2010; J. Li et al., 2014) to establish an SCON family of models and leverage the versatility of stochastic optimization. As is the case for CON, SCON models have three state variables representing SOC, DOC, and MBC densities in a soil system. We parameterized two SCON variants, "constant diffusion" SCON (SCON-C) and "state-scaling diffusion" SCON (SCON-SS). Diffusion coefficients are model parameters that govern the noise dynamics of an



Figure 1. In our study, we use normalizing flows to approximate SCON soil biogeochemical model solution trajectories x over time t. The flow operates in a generative direction, mapping a simpler base distribution to a more complex one representing SCON output.

¹¹⁹ SDE system. In SCON-C, diffusion, or noise, was set to be independent of time and states, ¹²⁰ while in SCON-SS, noise was made to depend on and scale with state values.

We used a class of methods called *normalizing flows* to approximate SCON mod-121 els in our inference approach. In simple terms, flows can be thought of as one or more layers of random variable mappings that transform an initial base probability distribu-123 tion to a new distribution (Papamakarios et al., 2021). When we deploy flows to trans-124 form a simpler probability density into a more complex one (Figure 1), as we do in our 125 study, it is classified as a *generative* normalizing flow (Kobyzev et al., 2020). The flow 126 approximation refashions SCON from an SDE that depicts state variable dynamics $\frac{dx}{dt}$ 127 in continuous time to a probabilistic state space model that specifies distributions of state 128 measurements y_t noisily observed from underlying states x_t in discrete time (Särkkä & 129 Solin, 2019). 130

The replacement of differential equation solver integration with state space mod-131 els to approximate dynamical systems offers substantial computational efficiency gains 132 in inference (Ryder et al., 2018; Särkkä & Solin, 2019). At each inference iteration or 133 epoch, rather than sequentially computing state trajectories x one time step at a time 134 with solvers including Euler, Runge-Kutta, and Adams' schemes, as was demonstrated 135 in studies like Xie et al. (2020), we can instead simultaneously sample multiple x in one 136 vectorized draw from a flow-transformed state space distribution object. This increased 137 efficiency allows us to more capably assimilate SBMs with time series data sets spanning 138 longer periods under computing resource limitations, especially when highly paralleliz-139 able graphical processing units (GPUs) can be leveraged. 140

¹⁴¹ Drawing from the methodologies of previous work that test various inference ap-¹⁴² proaches (Golightly & Wilkinson, 2006; Whitaker et al., 2017; Ryder et al., 2018, 2021), ¹⁴³ our study demonstrates functional stochastic VI of flow-approximated SBMs conditioned ¹⁴⁴ on soil observations data y that includes various soil pool and respiration measurements. ¹⁴⁵ To support the notion that our VI approach is operational, we show that it can fit model ¹⁴⁶ output to y sourced from a known data-generating process and estimate model θ pos-¹⁴⁷ teriors in line with the true θ values used by that process.

Hence, to begin our study workflow, we generated synthetic y consisting of SOC, 148 DOC, and MBC state and heterotrophic CO₂ respiration rate observations correspond-149 ing to SCON-C and SCON-SS data-generating processes. The processes used "true" θ 150 values randomly sampled from constrained data-generating distributions that were cho-151 sen to produce faster and more dramatic SOC decay dynamics reminiscent of organic 152 matter decomposition at soil surface, which contrasts with the slower and deeper soil de-153 composition depicted in J. Li et al. (2014) and Xie et al. (2020). Faster decay provided 154 our inference approach with substantive dynamical information in shorter time series to 155 operate and optimize on. We then conditioned our state space model VI on those syn-156



Figure 2. A workflow diagram summarizing the steps involved in our study's stochastic variational Bayesian framework. Our workflow efficiently conducts inference and data assimilation on stochastic differential equation (SDE) soil biogeochemical models (SBMs) with their approximation into state space models (SSMs). Our modular workflow is designed to serve as a basis for building future soil biogeochemical model inferences, as the "black box" inference method used can be modified or substituted. Our "black box" inference method of choice was stochastic gradient descent mean-field variational inference. Within the nodes of the diagram, blue lines and shading correspond to prior means and distributions, while orange lines and shading correspond to posterior means and distributions. Orange dots represent observations upon which the inference is conditioned.

thetic y for estimation of approximate posterior densities $q(\theta)$ that were compared with prior densities $p(\theta)$. Priors were made to be equivalent to our data-generating distributions.

Ultimately, we found that our VI approach allowed us to reasonably fit y. When 160 possible, our fits were checked against solutions from a Kalman smoother algorithm, and 161 we observed that the flow fits were mostly consistent with the Kalman solutions. Cru-162 cially, we were also able to recover some of the true θ values used by our data-generating 163 processes against model identifiability limitations that could not be resolved by the ex-164 tent of information contained in our synthetic data. Model identifiability can be sum-165 marized as the ability to update prior beliefs about θ and align model to true θ based 166 on available data. Our identifiability issues related to the presence of ambiguous SCON 167 equation terms involving the multiplication of more than two parameters. Our work of-168 fers insights and suggestions for improving the identification of θ , which is of interest for 169 experimentalists and biogeochemists who are interested in building effective data sets 170 for SBM inference. 171

¹⁷² 2 Materials and Methods

173

2.1 Inference workflow overview

The general steps constituting our study's SBM data assimilation workflow are outlined in Figure 2. We established SCON-C and SCON-SS to serve as known data-generating processes whose true θ values can be compared with the inferred posteriors to test our flow VI method. Discrepancies between the true θ and posterior means inform on the effectiveness of our selected inference algorithm. Differences between the priors and posterior densities further indicate algorithm efficacy and additionally point to the informativeness of the data y for identifying and constraining posteriors.

True SCON-C and SCON-SS θ were sampled from data-generating distributions 181 truncated between lower and upper support bounds to ensure that data-generating pro-182 cesses would remain in parameter regimes with faster state decay corresponding to soil 183 surface decomposition occurring on the order of thousands of hours, rather than tens or 184 hundreds of thousands. This allows us to generate shorter data sets y that enable reduced 185 computational loads and faster turnaround times for testing our inference algorithm while 186 retaining dynamical richness that can inform the algorithm to estimate more certain pos-187 teriors. We used logit-normal distributions to handle truncation in our data-generating, 188 prior, and posterior distributions, which we will describe in section 2.3. Our inference 189 priors matched our data-generating distributions. 190

Synthetic data y were observed and processed from our data-generating SDE so-191 lution trajectories. We parameterized our SCON models based in time units of hours, 192 so observations were collected every 5 hours by default. State space approximation of 193 SDE output, which we will describe in section 2.4, requires regular time series discretiza-194 tion (Kalman, 1960), so in an empirical setting, all existing, imputed, or missing obser-195 vations must coincide with discrete time steps of the state space model in our approach 196 and cannot transpire in between. Different SDE approximation methods would be needed 197 for irregular time discretization. 198

We selected mean-field stochastic VI as our black box inference method for its math-199 ematical simplicity and efficiency. Mean-field inference makes the simplifying assump-200 tion that model parameters are independently distributed. This aligns with our synthetic 201 data-generating processes, in which our true θ values are sampled from independent logit-202 normal distributions. VI frames Bayesian inference as an optimization goal of finding 203 the set of mean-field posterior distributions that best describes y. The optimization pro-204 cess takes place over a number of training iterations in which θ values are sampled at 205 each iteration and the likelihood of the resultant model output conditioned on y and θ 206 is evaluated in fulfillment of the objective of the VI algorithm to locate θ correspond-207 ing to higher model likelihood. We present an overview of our VI implementation and 208 key algorithm steps in section 2.5. 209

We used normalizing flows to approximate SCON-C and SCON-SS from continuous-210 time SDEs to time-discretized state space models. These state space approximations then 211 served as our bases for VI optimization. A brief treatment on state space models is given 212 in section 2.4. Flow state space approximation increased the computational efficiency 213 214 of sampling SCON solution trajectories (also referred to as *latent variables*, *states*, or *paths* in machine learning literature) such that multiple trajectories would be simultaneously 215 collected from a flow distribution object rather than sequentially simulated from a dif-216 ferential equation solver at each training iteration. The flow is assembled through deep 217 neural network layers that transform simpler random input into more complex approx-218 imated SCON output. The constituent pieces of the machine learning architecture un-219 derlying our flow are detailed in section 2.6. 220

Per equations (3) and (4), SCON-C is a completely linear SDE. Consequently, SCON-C flow-approximated x and its fit of y can be visually benchmarked against output from

an instance of the Kalman smoother algorithm summarized in section 2.7. Given a known 223 data-generating process and observation error, a Kalman smoother exactly solves the true 224 mean latent path x of the SDE data-generating process sourcing y. We successfully com-225 pared SCON-C flow x to the true x solution computed by the smoother, which we de-226 scribe in section 3.1. The smoother algorithm cannot resolve the non-linear diffusion de-227 picted in equation (5), so SCON-SS flow output could not be validated in the same man-228 ner. 229

230

2.2 SCON SDE parameterization and data generation

SDE system equations are frequently written with the state value derivatives dxon the left-hand side, and consist of a drift coefficient vector, frequently notated as α , and a diffusion coefficient matrix, notated as β , on the right-hand side. For biological SDE models, a square-root diffusion structure is frequently used such that these systems follow the form

$$dx_t = \alpha(x_t, t, \theta)dt + \sqrt{\beta(x_t, t, \theta)}dW_t$$
(1)

where dW_t denotes a continuous stochastic Wiener process. Evolution of SDE trajec-231

tories x across a simulation duration T in time increments dt can be interpreted as a se-232 ries of small steps whose values are independently drawn from a normal distribution with 233 mean $\alpha(x_t, t) dt$ and variance $\beta(x_t, t) dt$ (Särkkä & Solin, 2019). 234

Like the CON model introduced in Allison et al. (2010), SCON has three state dimensions made up of soil organic C (SOC), dissolved organic C (DOC), and microbial biomass C (MBC) densities. We notate total state dimensions with \mathcal{D} , so $\mathcal{D} = 3$ for all systems in the SCON family. SOC, DOC, and MBC are respectively notated in the system equations as S, D, and M. Thus, x_t , the solutions of the continuous SCON system at time t, expand to the vector,

$$x_t = \begin{bmatrix} S_t \\ D_t \\ M_t \end{bmatrix}$$
(2)

and observations of the system y_t are similarly three-dimensional. 235

We established two SCON versions for inference and data generation use, SCON-C and SCON-SS. SCON-C and SCON-SS share the same underlying α drift vector, equivalent to the deterministic CON dynamics and following the form:

$$\begin{bmatrix} dS \\ dD \\ dM \end{bmatrix} = \begin{bmatrix} \mathcal{I}_S + a_{DS} \cdot k_D \cdot D + a_M \cdot a_{MSC} \cdot k_M \cdot M - k_S \cdot S \\ \mathcal{I}_D + a_{SD} \cdot k_S \cdot S + a_M \cdot (1 - a_{MSC}) \cdot k_M \cdot M - (u_M + k_D) \cdot D \\ u_M \cdot D - k_M \cdot M \end{bmatrix} dt + \beta^{0.5} \begin{bmatrix} dW_S \\ dW_D \\ dW_M \end{bmatrix}$$
(3)

where β now refers to the diffusion matrix component of the SDE and the W_S , W_D , and 236 W_M elements of the Wiener process vector represent random draws from the distribu-237 tion $\mathcal{N}(0, \sqrt{\mathrm{d}t})$.

238

For simplification purposes, the β diffusion matrices of both systems were made to be diagonal only and free of covariance diffusion terms. SCON-C diffusion dynamics are given by

$$\beta = \begin{bmatrix} c_S & 0 & 0\\ 0 & c_D & 0\\ 0 & 0 & c_M \end{bmatrix}$$
(4)

while SCON-SS diffusion dynamics are

$$\beta = \begin{bmatrix} s_S \cdot S & 0 & 0 \\ 0 & s_D \cdot D & 0 \\ 0 & 0 & s_M \cdot M \end{bmatrix}$$
(5)

²³⁹ Thus, SCON-C diffusion noise is *additive*, meaning it is independent of the values of states

S, D, A and M, A and also stationary, meaning that is not a function of t. Meanwhile, SCON-

SS noise is *multiplicative*, meaning it is dependent on the states. As such, SCON-C is

linear in drift and diffusion, while SCON-SS is linear in drift but non-linear in diffusion.

 \mathcal{I}_S and \mathcal{I}_D respectively represent the exogenous input of C mass in units of mg C g⁻¹ soil h⁻¹ into the SOC and DOC soil pools from litter decay and can be modeled as constants or functions. We used sinusoidal litter input functions with annual periods that assumed litterfall peaking through late summer and early fall in a pattern resembling those observed in tropical forest ecosystems (Giweta, 2020). The functions are given by

$$\mathcal{I}_{S,t} = 0.001 + 0.0005 \cdot \sin\left(\frac{2\pi}{365 \cdot 24}t\right) \tag{6}$$

$$\mathcal{I}_{D,t} = 0.0001 + 0.00005 \cdot \sin\left(\frac{2\pi}{365 \cdot 24}t\right) \tag{7}$$

As was previously instituted for CON (Allison et al., 2010; J. Li et al., 2014), the SCON linear first-order decay parameters $k_{i \in \{S,D,M\}}$ remain dependent on temperature. Temperature sensitivity of the $k_{i \in \{S,D,M\}}$ linear first-order decay parameters is enforced by a function derived from the original Arrhenius equation (Arrhenius, 1889),

$$k_{i,t} = k_{i,\text{ref}} \exp\left[-\frac{Ea_{k_i}}{R} \left(\frac{1}{\text{temp}_t} - \frac{1}{\text{temp}_{\text{ref}}}\right)\right]$$
(8)

243 244 where R is the ideal gas constant 8.314 J K⁻¹ mol⁻¹ and temp_{ref} specifies a "reference" equilibrium temperature which we set at 283 K.

Through changing values of $k_{i \in \{S,D,M\}}$, SCON systems respond to day-night and seasonal temperature cycles through the composite sinusoid forcing function,

$$\operatorname{temp}_{t} = \operatorname{temp}_{\operatorname{ref}} + \frac{5t}{80 \cdot 365 \cdot 24} + 10 \cdot \sin\left(\frac{2\pi}{24}t\right) + 10 \cdot \sin\left(\frac{2\pi}{365 \cdot 24}t\right) \tag{9}$$

The function assumes a gradual linear increase in mean soil surface temperature by 5 °C over 80 years from the start of the simulation, in line with the upper bound of mean surface temperature increases predicted in emissions scenarios by 2100 (O'Neill et al., 2017).

SDE systems rarely admit tractable analytic solutions. To sample state trajectories accurately approximating SCON-C and SCON-SS dynamics and construct our synthetic time series data y, we used the long-established and reliable Euler-Maruyama SDE solver (Maruyama, 1955) to numerically integrate solution paths x corresponding to θ randomly sampled from logit-normal distributions. Our solver step size was set to dt =0.1 hour. We note that we recover the exact SCON-C and SCON-SS processes in continuous time as dt is decreased to 0.

If inference involved conditioning with CO_2 observations in y in addition to state measurements, model CO_2 respiration rate would be computed from the time-corresponding x state values with the equation

$$r_{CO_2,t} = (1 - a_{SD}) \cdot k_{S,t} \cdot S_t + (1 - a_{DS}) \cdot k_{D,t} \cdot D_t + (1 - a_M) \cdot k_{M,t} \cdot M_t$$
(10)

where $r_{CO_2,t}$ is in units of $\mu g g^{-1} \sinh^{-1}$. We then sliced x and CO₂ time series at some regular interval, i.e. every 1 hour or 5 hours, and normally sampled about the sliced values with an observation error vector σ_{obs} in the manner of

$$y_t \sim \mathcal{N}(x_t, \eta_{\text{obs}})$$
 (11)

to arrive at y. We lower bounded y such that $y \in \mathbb{R}_{\geq 0}$ to preclude nonsense negative state measurements. We used constant η_{obs} that was 10% of the overall state mean such that

$$\eta_{\rm obs} = 0.1 \odot \begin{bmatrix} \bar{S} & 0 & 0\\ 0 & \bar{D} & 0\\ 0 & 0 & \bar{M} \end{bmatrix}$$
(12)

where \odot indicates elementwise multiplication. This corresponds to an empirical scenario where measurement instruments and processes introduce a stable level of observation noise.

We generated and conditioned inferences on synthetic y of up to 5000 hours in to-258 tal timespan T. Data-generating θ distribution hyperparameters were chosen to produce 259 stable and informative state dynamics in a shorter span of time and minimize the mem-260 ory footprint of the data set under available computing resources. We used elevated $k_{i,ref}$ means compared to previous literature values (Allison et al., 2010; J. Li et al., 2014; Xie 262 et al., 2020). Sampled θ values and T scale are thereby reminiscent of an organic decay 263 process occurring at the soil surface, rather than a slower subterranean decomposition. 264 θ data-generating distribution hyperparameters, equivalent to the prior distribution $p(\theta)$ 265 hyperparameters, along with the biogeochemical interpretations associated with each θ , 266 are detailed in Table S1. 267

2.3 The generalized univariate logit-normal distribution

268

We used a univariate logit-normal distribution family for our data-generating, in-269 formed prior $p(\theta)$, and mean-field variational posterior $q(\theta|y)$ probability density func-270 tions. To avoid being restricted to the standard [0,1] distribution support that the logit-271 normal density is typically associated with in statistics, we defined a generalized form 272 of the family whose supports could be enclosed between an arbitrary positive [a, b], where 273 $a, b \in \mathbb{R}_{>0}$ and b > a. Generalized logit-normal distributions provide similar utility 274 to truncated normal distributions used previously in SBM inference projects for constrain-275 ing θ values to finite supports (Xie et al., 2020), but are more stable for backpropaga-276 tion, as the inverse cumulative distribution function of the truncated normal distribu-277 tion has inherent stability issues close to support boundaries. 278

We notate logit-normal distribution parameters in order of desired "target" mean μ , standard deviation σ , support lower bound a, and upper bound b akin to

$$\theta \sim \mathscr{LN}(\mu, \sigma, a, b) \tag{13}$$

Via passage through a sigmoid function, logit-normal distributions are transformed from normal distributions $\mathcal{N}(\check{\mu},\check{\sigma})$, where $\check{\mu}$ and $\check{\sigma}$ are respectively the "parent" mean and standard deviation distribution parameters:

$$\check{\theta} \sim \mathcal{N}(\check{\mu}, \check{\sigma}) \tag{14}$$

$$\theta^{\text{mid}} = \frac{1}{1 + \exp(-\check{\theta})} \tag{15}$$

$$\theta = (b-a) \cdot \theta^{\text{mid}} + a \tag{16}$$

The logit-normal distribution has no closed form probability density function and 279 its probability moments are not analytically resolvable, so no formula can be deduced 280 that allows us to make random variable transformations between logit normal and nor-281 mal distributions. Hence, to arrive at a particular logit-normal density with target μ and 282 σ in each VI optimization iteration to sample from, we must first numerically solve for 283 the parent $\check{\mu}$ and $\check{\sigma}$ of a normal distribution that corresponds to the desired logit-normal 284 properties following the transformations from equations (14) to (16). We can do this with 285 root-finding algorithms like the bisection method that search for an appropriate $\check{\mu}$ in the 286 truncated support interval between a and b and $\check{\sigma}$ within a provided range of tolerated 287 standard deviation values (Daunizeau, 2017). 288

2.4 State space model approximation of the SDE

289

Instead of optimizing SCON θ via an iterative SDE solver, we optimized time-discretized state space models approximating the SCON-C and SCON-SS SDEs. State space models describe the distribution of a discrete sequence of observations y sourced from discrete latent states x. They can take the general form

$$x_t \sim p(x_t | x_{t-1}, \theta) \tag{17}$$

$$y_t \sim p(y_t | x_t, \theta) \tag{18}$$

Equation (17) indicates that the transition from x_{t-1} to x_t occurs at a probability density of $p(x_t|x_{t-1},\theta)$ and that subsequent states of a state space model depend on previous ones, thus indicating that x constitutes a Markov chain. Equation (18) specifies that y_t is observed from x_t at a density of $p(y_t|x_t,\theta)$. An initial state x_0 must be nominated to compute x and it can be set as a constant, or informed as a density, $p(x_0)$, which we do in our case.

The state space model θ are the same model parameters as in the SDE counterpart. When accounting for the SDE α drift and β diffusion dynamics, x_t , the latent states of the state space model at time t can be written as

$$x_t = x_{t-1} + \alpha(x_{t-1}, \theta)\Delta t + \epsilon_t \sqrt{\beta(x_{t-1}, \theta)\Delta t}$$
(19)

with the same α and β as in (1). ϵ_t is a random noise vector of length \mathcal{D} independently sampled via $\epsilon_t \sim \mathcal{N}(0, I_{\mathcal{D}})$. $I_{\mathcal{D}}$ is an identity matrix with number of diagonal elements equal to \mathcal{D} . Δt is the state space model time step, not to be confused with SDE solver time step dt. We used $\Delta t = 1.0$ hour for our state space model approximations in contrast to the aforementioned dt = 0.1 for Euler-Maruyama solving of our data generating processes.

There is overlap between SDEs and state space models. Both depict the evolution of state values in a series of steps where future values depend on past ones. Both require the specification of initial conditions x_0 to compute state trajectories.

However, SDEs and state space models treat time differently. A key distinction that makes state space model approximation helpful for inference efficiency is that Δt can be made relatively large versus SDE solver dt. This is helpful for common biological inference and data assimilation situations where y is sparsely observed due to the expense and difficulty of collecting measurements.

Differential equation systems are instead typically numerically integrated and like state space models, are solved in discrete steps, as only smooth analytic solutions can only be obtained from the simplest systems. But, the differential equation integration procedures still assume that states are evolving in continuous time. The integrating solvers almost always require relatively small integration time steps dt that are as close to 0 as possible; the solvers tend to fail at higher dt.

The divergent handling of time in state space models and SDEs renders them more apt for different objectives. State space models are better suited for estimating observations over long T, whereas SDEs are required for precise analyses of accurately simulated system dynamics. With their differing but related purposes, we can ultimately use state space models to represent discrete observations from continuous SDEs.

321 2.5 VI optimization

Under a Bayesian statistics framework, the goal of statistical inference broadly consists of estimation of the θ posterior density function for some model, $p(\theta|y)$, conditioned on some data set y via Bayes' rule,

$$p(\theta|y) = \frac{p(y|\theta)p(\theta)}{p(y)}$$
(20)

 $p(y|\theta)$, also notated as $\ell(\theta|y)$, is the likelihood function, which indicates model goodnessof-fit across various values of individual parameters comprising θ . $p(\theta)$ is the prior probability representing beliefs about θ uncertainty. p(y) is the probability density of the observed data.

The prior density $p(\theta)$ can be specified in an informed fashion, as we did in our workflow with distributions that matched our known data-generating distributions, or with less certainty in the absence of empirical information or domain experience. In most cases, $p(y|\theta)$ is not obtainable in closed analytic form and has to be numerically estimated with methods including Monte Carlo sampling and Laplace approximation (Reid, 2015). Additionally, p(y), sometimes known as the *marginal evidence*, tends to be unresolvable (Gelman et al., 2013; McElreath, 2020). Thus, we take advantage of the proportionality relationship based on (20),

$$p(\theta|y) \propto p(y|\theta)p(\theta) = p(y,\theta) \tag{21}$$

to estimate $p(\theta|y)$ and practically conduct inference.

For Bayesian inference on state space models, we additionally need to account for the transition and observation densities generalized in equations (17) and (18), which influence the θ posterior. Estimation of the posterior of θ in state space model inference must occur along with estimation of the posterior of x, whether in a joint or marginal fashion, in a case such as ours where the transition and observation distributions are not known. We opted for joint estimation. The joint posterior density of θ and x is notated as $p(\theta, x|y)$. We arrive at an expression for $p(\theta, x|y)$ by substituting (17) and (18) into (21):

$$p(\theta, x|y) \propto p(y, \theta, x)$$
 (22)

$$= p(y|\theta, x)p(\theta, x)$$
(23)

$$= p(y|\theta, x)p(\theta)p(x|\theta)$$
(24)

$$= p(\theta) \prod_{i \in \mathfrak{N}} p(y_i | x_i, \theta) \prod_{t=1}^{T} p(x_t | x_{t-1}, \theta)$$
(25)

T denotes the final discretized integer time index of x. Since we set state space model $\Delta t = 1.0$ hour, our final time index matches total synthetic experiment hours and Tcan signify both. We also use T to represent the set of x state space model discretization indices not including the initial state at t = 0. We can then adopt a $\mathfrak{T} \subseteq T$ to indicate the set of y observation time indices not including an initial observation at t =0, which is required in our VI procedure. The total number of x discretizations is N =T + 1 when including the t = 0 index. $\mathfrak{N} = \{0\} \cup \mathfrak{T}$ notates the full set of y indices.

In stochastic VI on state space models, we optimize a parametric density $q(\theta, x)$ to match the true joint posterior $p(\theta, x|y)$ as closely as possible. Per the probability chain rule, $q(\theta, x)$ expands to,

$$q(\theta, x) = q(x|\theta)q(\theta) \tag{26}$$

The density functions we select for our marginalized $q(\theta)$ and $q(x|\theta)$ are known as our *variational families*. As mentioned in section 2.3, we chose a mean-field logit-normal variational family for $q(\theta)$ that assumed independent univariate distributions per θ such that

$$q(\theta) = q(\theta_1, \theta_2, \dots, \theta_{\mathcal{P}}) = \prod_{j=1}^{\mathcal{P}} q_j(\theta_j)$$
(27)

where \mathcal{P} is the total number of individual SBM θ and $\mathcal{P} = 14$ for SCON-C and SCON-SS (Table S1). We chose a class of normalizing flow called a *neural moving average flow* in Ryder et al. (2021) for our $q(x|\theta)$ variational family, which we will describe subsequently in section 2.6.

We can index individual representatives of our joint variational family by the prop-338 erties and hyperparameters of the distribution symbolized in aggregate by $\phi_{(\theta,x)}$ such 339 that an instance of $q(\theta, x)$ is notated as $q(\theta, x; \phi_{(\theta,x)})$. $q(\theta, x; \phi_{(\theta,x)})$ can be decomposed 340 into $q(\theta; \phi_{\theta})q(x|\theta; \phi_x)$ since $\phi_{(\theta,x)} = (\phi_{\theta}, \phi_x)$. ϕ are termed variational parameters, as 341 they represent the distribution settings that can be varied and tuned to adjust the ap-342 proximation. For neural network models like flows, variational parameters would include 343 the hidden neural network parameters and weights. A particular distribution can be re-344 ferred to by its variational parameter index in shorthand. 345

Our VI framework seeks a set of variational parameters ϕ that minimizes discrepancies between $q(\theta, x; \phi_{(\theta,x)})$ and $p(\theta, x|y)$, the approximate and true posteriors. One measure of distance between distributions customarily employed in statistics and machine learning literature is called the *Kullback-Leibler (KL) divergence*, notated as $D_{\text{KL}}[q(\theta, x; \phi_{(\theta,x)})||p(\theta, x|y)]$ (Kullback & Leibler, 1951; Perez-Cruz, 2008; Joyce, 2011). Targeting the KL divergence for minimization, our optimization objective can then be mathematically stated as,

$$q(\theta, x; \phi^*_{(\theta, x)}) = \operatorname{argmin}_{q(\theta, x; \phi_{(\theta, x)})}(D_{\mathrm{KL}}[q(\theta, x; \phi_{(\theta, x)})||p(\theta, x|y)])$$
(28)

where $\phi^*_{(\theta,x)}$ indexes the set of variational parameters that corresponds to the idealized global KL divergence minimum. After several omitted steps that can be referenced in greater detail from Blei et al. (2017), we proceed from (28) to

$$D_{\mathrm{KL}}[q(\theta, x; \phi_{(\theta, x)})||p(\theta, x|y)] = \mathbb{E}_{q(\theta, x; \phi_{(\theta, x)})}[\log q(\theta, x; \phi_{(\theta, x)})] - \mathbb{E}_{q(\theta, x; \phi_{(\theta, x)})}[\log p(y|\theta, x)]$$
(29)
= $\mathbb{E}_{q(\theta, x; \phi_{(\theta, x)})}[\log q(\theta, x; \phi_{(\theta, x)})] - \mathbb{E}_{q(\theta, x; \phi_{(\theta, x)})}[\log p(y, \theta, x)] + \log p(y)$ (30)

where the expectations \mathbb{E} subscripted with $q(\theta, x; \phi_{(\theta,x)})$ are taken with respect to the density $q(\theta, x; \phi_{(\theta,x)})$.

Reviewing the notion that p(y) and in turn, the log marginal evidence, are typically unavailable (Christensen et al., 2010), we then rely on a reduced and rearranged expression that constitutes the ELBO function, notated as \mathcal{L} ,

$$\mathcal{L}[\phi_{(\theta,x)}] = \mathbb{E}_{q(\theta,x;\phi_{(\theta,x)})}[\log p(y,\theta,x)] - \mathbb{E}_{q(\theta,x;\phi_{(\theta,x)})}[\log q(\theta,x;\phi_{(\theta,x)})]$$
(31)

$$= \mathbb{E}_{q(\theta, x; \phi_{(\theta, x)})} [\log p(y, \theta, x) - \log q(\theta, x; \phi_{(\theta, x)})]$$
(32)

$$= \mathbb{E}_{q(\theta, x; \phi_{(\theta, x)})} \left\langle \log p(y, \theta, x) - \log[q(x|\theta; \phi_x)q(\theta; \phi_\theta)] \right\rangle$$
(33)

$$= \mathbb{E}_{q(\theta, x; \phi_{(\theta, x)})} \left\langle \log p(y, \theta, x) - \log q(x|\theta; \phi_x) - \log q(\theta; \phi_\theta) \right\rangle$$
(34)

Substituting in (25) for $p(y, \theta, x)$ results in

$$\mathcal{L}[\phi_{(\theta,x)}] = \mathbb{E}_{q(\theta,x;\phi_{(\theta,x)})} \langle \log \left[p(\theta) \prod_{i \in \mathfrak{N}} p(y_i | x_i, \theta) \prod_{t=1}^T p(x_t | x_{t-1}, \theta) \right] - \log q(x|\theta;\phi_x) - \log q(\theta;\phi_\theta) \rangle \quad (35)$$

which, recalling that the set of total y indices $\mathfrak{N} = \{0\} \cup \mathfrak{T}$, expands into

$$\mathcal{L}[\phi_{(\theta,x)}] = \mathbb{E}_{q(\theta,x;\phi_{(\theta,x)})} \langle \log p(\theta) + \log p(y_0|x_0,\theta) + \sum_{i\in\mathfrak{T}} \log p(y_i|x_i,\theta) + \sum_{t=1}^T \log p(x_t|x_{t-1},\theta) - \log q(x|\theta;\phi_x) - \log q(\theta;\phi_\theta) \rangle \quad (36)$$

We will decompose the marginal variational log-density of x, $\log q(x|\theta; \phi_x)$, in more gran-

³⁴⁹ ular detail when we describe the architecture of the neural moving average flow in sec-

 $_{350}$ tion 2.6.

The ELBO function is called as such because it is the lower bound of the log marginal evidence:

$$\log p(y) = \mathcal{L}[\phi_{(\theta,x)}] + D_{\mathrm{KL}}[q(\theta,x;\phi_{(\theta,x)})||p(\theta,x|y)]$$

$$\geq \mathcal{L}[\phi_{(\theta,x)}]$$
(37)
(38)

Maximizing $\mathcal{L}[\phi_{(\theta,x)}]$, or minimizing the negative ELBO $-\mathcal{L}$, as we need to do in machine learning libraries like PyTorch that implement gradient descent rather than ascent, is commensurate to minimizing $D_{\mathrm{KL}}[q(\theta, x; \phi_{(\theta,x)})||p(\theta, x|y)]$. Hence, $\mathcal{L}[\phi_{(\theta,x)}]$ is our optimization objective function.

For pithier description of the ELBO gradient, $\nabla \mathcal{L}$, used to update $\phi_{(\theta,x)}$ via automatic differentiation, we set $\log p(y, \theta, x) - \log q(\theta, x; \phi_{(\theta,x)})$ in (32) equal to $\mathscr{R}(\theta, x, y, \phi)$, where \mathscr{R} is a log-density ratio function. This reduces the ELBO equation to

$$\mathcal{L}[\phi_{(\theta,x)}] = \mathbb{E}_{q(\theta,x;\phi_{(\theta,x)})} \left[\mathscr{R}(\theta,x,y,\phi_{(\theta,x)}) \right]$$
(39)

and the ELBO gradient is

$$\nabla \mathcal{L}[\phi_{(\theta,x)}] = \nabla_{\phi} \left\langle \mathbb{E}_{q(\theta,x;\phi_{(\theta,x)})} \left[\mathscr{R}(\theta,x,y,\phi_{(\theta,x)}) \right] \right\rangle$$
(40)

$$= \nabla_{\phi} \left[\int_{\theta} \int_{x} q(\theta, x; \phi_{(\theta, x)}) \mathscr{R}(\theta, x, y, \phi_{(\theta, x)}) \mathrm{d}x \mathrm{d}\theta \right]$$
(41)

$$= \int_{\theta} \int_{x} \nabla_{\phi} \left[q(\theta, x; \phi_{(\theta, x)}) \mathscr{R}(\theta, x, y, \phi_{(\theta, x)}) \right] \mathrm{d}x \mathrm{d}\theta \tag{42}$$

which decomposes to

$$\nabla \mathcal{L}[\phi_{(\theta,x)})] = \int_{\theta} \int_{x} q(\theta, x; \phi_{(\theta,x)}) \nabla_{\phi} \left[\mathscr{R}(\theta, x, y, \phi_{(\theta,x)}) \right] \mathrm{d}x \mathrm{d}\theta + \int_{\theta} \int_{x} \mathscr{R}(\theta, x, y, \phi_{(\theta,x)}) \nabla_{\phi} \left[q(\theta, x; \phi_{(\theta,x)}) \right] \mathrm{d}x \mathrm{d}\theta$$
(43)

Note that the gradients ∇_{ϕ} are taken with respect to the variational parameters. This presents a complication, as examining the second term of (43), we are left with the situation that $\nabla_{\phi} \left[q(\theta, x; \phi_{(\theta, x)}) \right]$ is by and large unavailable, as q can be sampled from, but is usually not analytically differentiable. Our joint variational family q is no exception to that pattern; our marginal mean-field $q(\theta; \phi_{\theta})$ has the straightforward analytic form given in (27), but use of the neural moving average flow as the variational family for $q(x|\theta; \phi_x)$ precludes the overall joint density $q(\theta, x; \phi_{(\theta, x)})$ from having an orderly closed form.

To ultimately compute the gradient of an expectation as in (40) in numerical fashion, we thereby turn to the *reparameterization trick* set forth in Salimans and Knowles (2013) and Kingma and Welling (2014). The reparameterization trick involves setting (θ, x) as an output of an invertible, deterministic, and differentiable function $g(z, \phi_{(\theta,x)})$, where z is a random vector sampled from some fixed density q(z). This enables us to tractably take a gradient of a simpler fixed distribution whose probability density function is easier to differentiate and not dependent on the variational parameters ϕ (Ruiz et al., 2016).

In our case, z elements are sampled from standard normal distributions and undergo invertible transformations to proceed to x. θ is still directly sampled from its meanfield logit-normal family described in section 2.5 as part of the operations of g. Hence, \mathcal{L} can be estimated with each VI training iteration with Monte Carlo sampling of z and θ entries starting with the steps

$$z^{(s)} \sim \mathcal{N}(0, I_N) \tag{44}$$

$$\theta^{(s)}, x^{(s)} = g(z^{(s)}, \phi_{(\theta, x)}) \tag{45}$$

where I_N is an identity matrix with number of diagonal entries matching the total x discretization indices N. The superscript (s) indexes an individual draw from a distribution. We can then re-frame (40) from an analytically intractable gradient of an expectation to a numerically assessable expectation of a gradient with

$$\nabla \mathcal{L}[\phi_{(\theta,x)}] = \nabla_{\phi} \left\langle \mathbb{E}_{q(z)} \left[\mathscr{R}(\theta, x, y, \phi_{(\theta,x)}) \right] \right\rangle$$
(46)

$$= \mathbb{E}_{q(z)} \left[\nabla_{\phi} \left\langle \mathscr{R}(\theta, x, y, \phi_{(\theta, x)}) \right\rangle \right]$$
(47)

$$\approx \frac{1}{S} \sum_{s=1}^{S} \nabla_{\phi} \left\langle \mathscr{R}(\theta^{(s)}, x^{(s)}, y, \phi_{(\theta, x)}) \right\rangle$$
(48)

$$\nabla \widehat{\mathcal{L}}[\phi_{(\theta,x)}] = \frac{1}{\mathcal{S}} \sum_{s=1}^{\mathcal{S}} \nabla_{\phi} \left\langle \mathscr{R} \left[g(z^{(s)}, \phi_{(\theta,x)}), y, \phi_{(\theta,x)}) \right] \right\rangle$$
(49)

 \mathcal{S} denotes the total number of independent θ and z samples drawn per training itera-

tion. $\nabla \mathcal{L}[\phi_{(\theta,x)}]$ specifies the Monte Carlo estimate of $\nabla \mathcal{L}[\phi_{(\theta,x)}]$.

371

2.6 Masked neural moving average flow architecture

Delineating a normalizing flow more formally than in section 1, a general flow is a chain of *bijections*, or invertible transformation functions mapping an object in a set one-to-one to an object in another set. Flows can be decomposed into

$$x = g(z) = g_M \circ g_{M-1} \circ \dots \circ g_m \circ \dots \circ g_1(z)$$
(50)

 $_{372}$ where \circ notates function composition operations and M marks the total number of bi-

jections. In the generative direction, our flow takes us from a random object z to a random object x corresponding to a set of approximated SCON state trajectories.

A generative flow is constructed such that computation of $\log q(x|\theta; \phi_x)$ in (36) is available to facilitate the optimization of $q(x|\theta; \phi_x)$. The log-density of x is available through the change of variables formula:

$$\log q(x) = \sum_{t=1}^{T} \varphi(z_t) - \log |\det J|$$
(51)

$$\log q(x) = \sum_{t=1}^{T} \varphi(z_t) - \log \prod_{m=1}^{M} |\det J_m|$$
(52)

$$\log q(x) = \sum_{t=1}^{T} \varphi(z_t) - \sum_{m=1}^{M} \log |\det J_m|$$
(53)

where J is the Jacobian matrix of the overall transformation and J_m is the Jacobian of bijection g_m with respect to the intermediate transformed variable $g_{m-1} \circ g_{m-2} \circ \cdots \circ$ $g_1(z)$. We use $\varphi(z_t)$ to indicate the log-density of each element of z, z_t . We establish that z here is equivalent to the z introduced in section 2.5, so each $\varphi(z_t)$ is then a unit standard normal log-density in our framework. We notate the density function of z with q(z). Since q(z) is the starting distribution before transformations are layered, it is also termed the base distribution.

The particular flow we implemented as the marginal variational family for $q(x|\theta)$ was patterned after the original neural moving average flow introduced in Ryder et al. (2021). Neural moving average flows include *affine* bijections (Dinh et al., 2015; Kingma et al., 2016; Dinh et al., 2017; Papamakarios et al., 2017) among the functions constituting g in which an x^{out} is transformed from an x^{in} in the general form of

$$x^{\text{out}} = \mu + \sigma \odot x^{\text{in}} \tag{54}$$

where \odot represents elementwise multiplication to denote that μ , σ , and x^{in} can be ma-382 trices and vectors in addition to scalars, though our explicit situation involves scalars. 383 μ and σ are respectively known as shift and scale values of the bijection and it is required 384 that $\sigma \in \mathbb{R}_+$. Cumulative μ and σ values of a flow are usually implemented as trained 385 outputs of a neural network and are super- and subscripted to identify the transforma-386 tion layer and input elements they correspond to. They are notated as such by conven-387 tion and not to be confused with the similarly notated target logit-normal mean and stan-388 dard deviation parameters in section 2.3. 389

These linear affine transformations are basic in structure and consequently are individually not so *expressive*, or able to flexibly transition a base distribution into substantially different distributions of varying complexity. However, when layered repeatedly and stacked, their cumulative expressivity increases and with sufficient layers, composite affine functions can come to embody any distribution that is log-concave and bookended by declining density tails (Lee et al., 2021), which represents a large swath of probability distributions.

Neural moving average flows are specifically distinguished from other flows contain-397 ing affine layers through their execution of affine bijections with 1-dimensional convo-398 lutional neural networks (CNNs). To apply 1-dimensional CNNs rather than 2-dimensional 399 CNNs, we note that for systems with $\mathcal{D} > 1$, like SCON family instances, we must be-400 gin with z in a 1-dimensional "melted" form that is $\mathcal{D} \cdot T$ elements in length before re-401 shaping the final transformed x to a $\mathcal{D} \times T$ matrix matching the SDE solution struc-402 ture demonstrated in (2) following the conclusion of g. Thus, in equations (44) and (53), 403 we replace T with $\mathcal{D} \cdot T$ in our implementation. 404

Through masking, in which inputs to the convolution patch are zeroed out through multiplication by weights, the flow is imbued with an *autoregression* property in which the σ_i and μ_i values producing an *i*th element of an output vector x does not depend and convolve on any element $z_{j\geq i}$ in the base input vector. This autoregression is critical for the intent of arranging the computation of $\sum_{m=1}^{M} \log |\det J_m|$ in (53) to be manageable. The autoregression ensures that J is a diagonal matrix whose non-zero elements are the σ scale parameters underlying the overall transformation, which simplifies calculation of det J and $\log q(x)$ to

$$\log q(x) = \sum_{t=1}^{T} \varphi(z_t) - \sum_{m=1}^{M} \sum_{t=1}^{T} \log \sigma_t^m$$
(55)

where σ_t^m is the shift parameter of the bijection producing the *t*th term of the *m*th affine layer output of length *T*.

Figure 3b portrays a schematic of the autoregressive convolutions and affine bijec-407 tions used in our specific neural moving average flow implementation. The operations 408 occur within residual blocks, component pieces of deep learning networks consisting of 409 organizations of layers oriented toward the mitigation of training and approximation er-410 ror that can otherwise snowball with greater network depth. Residual blocks do this with 411 the use of *skip connections*, which preserve and carry over output from previous layers 412 to serve as input to subsequent layers and in doing so prevent noisy degradation of in-413 formation cascading through the network (He et al., 2016). 414

In each residual block, we perform two masked 1-dimensional convolutions, Convolution A and Convolution B, that each have a kernel length of 3 elements and a stride length of 1. To enshrine autoregressiveness of the flow, Convolution A applies a kernel masked as [1, 0, 0] that outputs a shift and scale value pair. The Convolution A operation and associated affine bijection can be generally expressed as

$$(\mu_i, \sigma_i) = f_i^{\mathcal{A}}(x_{i-1}^{\text{in}}) \tag{56}$$

$$x_i^{\text{out}} = \mu_i + \sigma_i \cdot x_i^{\text{in}} \tag{57}$$



Figure 3. Architecture blueprint of the neural moving average flow used as the marginal variational family for $q(x|\theta)$. (a) outlines the sequence of layers and operations. The affine block is a residual block in which the autoregressive convolution operations that distinguish neural moving average flows occur. (b) illustrates the two bijections, Convolution A and Convolution B, that link three hypothetical layers x^{in} , x^{mid} , and x^{out} together in each instance of an affine layer in our particular flow. Convolution A applies a [1,0,0] mask, while Convolution B applies a [1,1,0] mask. The example affine μ and σ parameters are indexed by superscripts and subscripts respectively identifying the layer and element they are associated with.

where μ_i , σ_i , x_i^{in} , and x_i^{out} are scalar elements of vectors and f_i^{A} is the Convolution A operation. The subsequent Convolution B involves a single stride kernel masked as [1, 1, 0] and it can be expressed together with its associated bijection as

$$(\mu_i, \sigma_i) = f_i^{\mathrm{B}}(x_{i-1}^{\mathrm{in}}, x_i^{\mathrm{in}})$$
(58)

$$x_i^{\text{out}} = \mu_i + \sigma_i \cdot x_i^{\text{in}} \tag{59}$$

⁴¹⁵ Combined, the two convolutions in sequence have a total receptive field length of 2.

To be able to produce the μ and σ parameters associated with the affine transformation of vector endpoint elements under autoregressive alignment, both convolutions require zero padding, in which zero elements are added to either end of the vector. As can be gleaned from Figure 3b, without zero padding, the kernels producing $(\mu_1^{\text{mid}}, \sigma_1^{\text{mid}})$ and $(\mu_1^{\text{out}}, \sigma_1^{\text{out}})$ would lack 1 element to convolve on, and the kernels sourcing $(\mu_N^{\text{mid}}, \sigma_N^{\text{mid}})$ and $(\mu_2^{\text{out}}, \sigma_0^{\text{out}})$ would by overhang their vectors by 1. A zero pad of length 1 was thereby sufficient for our purposes.

Simplified from our actual implementation and not pictured in Figure 3b is our ex-423 pansion of input into many *channels*, which are duplicates of the input vector that are 424 stacked on top of each other in a matrix. At each convolution stage, the same kernel is 425 applied in parallel across all the channels. Enlarging channel depth broadens the space of neural network weight values constituting $f_i^{\rm A}$ and $f_i^{\rm B}$ that can be explored per train-427 ing iteration. We set the number of channels at 96 for both convolutions and did not ex-428 periment further with channel depth. Also not pictured in Figure 3b, but implied in Fig-429 ure 3a, is the injection of auxiliary features extracted from y and observation indices \mathfrak{N} 430 in the form of vectors stacked on top of the input channels to inform training of the neu-431 ral network weights associated with the shift and scale values. Further elaboration on 432 the incorporation of auxiliary information is available in the supplement of Ryder et al. 433 (2021).434

In the overall flow procedure, the convolutions and affine bijections in the affine 435 residual block are linked with other transformations that we organize into repeatable sets 436 of layers. The order of transformations for each layer set is outlined in Figure 3a. Pre-437 ceding the affine blocks are *order-reversing permutations*, in which element order of a 438 vector input is flipped such that a vector $[x_1^{\text{in}}, x_2^{\text{in}}, ..., x_N^{\text{in}}]$ becomes $[x_1^{\text{out}}, x_2^{\text{out}}, ..., x_N^{\text{out}}] =$ 439 $[x_N^{\text{in}}, x_{N-1}^{\text{in}}, ..., x_1^{\text{in}}]$. Order-reversing permutations are a method of extending the expres-440 sivity and stability of a flow by enabling more complex dependency structures while pre-441 serving flow autoregression (Papamakarios et al., 2021). We found that adding order re-442 versals allowed us to modestly boost our ELBO learning rates. The permutations can 443 be seamlessly interspersed between other transformations since their absolute Jacobian 444 determinant is valued at 1, so they do not affect the computation of $\log q(x)$. 445

Differing from the neural moving average flow of Ryder et al. (2021), our flow fol-446 lows affine blocks with *batch renormalization* transformations. Batch renormalization 447 is a simple extension of *batch normalization*, which is a means of normalizing and reg-448 ularizing our variational samples such that our optimization is less influenced by ran-449 dom fluctuations in neural network weights and sample characteristics from one train-450 ing iteration to the next (Ioffe & Szegedy, 2015). Similar in intent but not operation to 451 permutations, batch normalization and renormalization are applied to bolster algorithm 452 stability and flexibility with increasing layer depth. They empirically allow VI algorithms 453 to tolerate higher learning rates (Bjorck et al., 2018), poor initialization of variational 454 parameters ϕ (Zhu et al., 2020), and erratic base distribution $z^{(s)}$ draws. 455

Batch normalization and renormalization overlap in the following steps that compute a batch mean μ_{S} and batch standard deviation σ_{S} from input x^{in} samples, not to be confused with the affine bijection and logit-normal μ and σ :

$$\mu_{\mathcal{S}} = \frac{1}{\mathcal{S}} \sum_{s=1}^{\mathcal{S}} x_s^{\rm in} \tag{60}$$

$$\sigma_{\mathcal{S}} = \sqrt{\varepsilon + \frac{1}{\mathcal{S}} \sum_{s=1}^{\mathcal{S}} (x_s^{\rm in} - \mu_{\mathcal{S}})^2}$$
(61)

where ε is a small constant added for stability. $\mu_{\mathcal{S}}$ and $\sigma_{\mathcal{S}}$ are involved in computation of the optimization objective—again, $\mathcal{L}[\phi_{(\theta,x)}]$ for our purposes—during the model training phase. They also update a lagging *running average* $\mu_{\mathcal{R}}$ and *running mean* $\sigma_{\mathcal{R}}$ that are less sensitive to change. $\mu_{\mathcal{R}}$ and $\sigma_{\mathcal{R}}$ are used after training of the model—the joint variational family $q(\theta, x; \phi_{(\theta,x)})$ in this setting—has been halted to estimate the objective metric at the testing stage.

In the testing phase, batch renormalization and normalization are equivalent in transforming input to output:

$$x^{\rm mid} = \frac{x^{\rm in} - \mu_{\mathcal{R}}}{\sigma_{\mathcal{R}}} \tag{62}$$

$$x^{\text{out}} = \gamma \cdot x^{\text{mid}} + \Upsilon \tag{63}$$

The collection of γ_t^m and Υ_t^m parameters in each flow layer set are learned neural network outputs. Batch renormalization diverges from batch normalization during training with the steps

$$r = \frac{\sigma_{\mathcal{S}}}{\sigma_{\mathcal{R}}} \tag{64}$$

$$d = \frac{\mu_{\mathcal{S}} - \mu_{\mathcal{R}}}{\sigma_{\mathcal{R}}} \tag{65}$$

$$x^{\rm mid} = \frac{x^{\rm in} - \mu_{\mathcal{S}}}{\sigma_{\mathcal{S}}} \cdot r + d \tag{66}$$

$$x^{\text{out}} = \gamma \cdot x^{\text{mid}} + \Upsilon \tag{67}$$

where r and d are variable correction factors. r and d are intended to limit the divergence between batch and running sample characteristics. r is clipped between the interval $[1/r_{\max}, r_{\max}]$, where r_{\max} is gradually increased to 3 over the course of inference, and d is clipped between the interval $[-d_{\max}, d_{\max}]$, where d_{\max} is gradually increased to 5. These intervals were established based on guidelines from previous empirical work (Ioffe, 2017). Batch normalization is a special case of batch renormalization where r =1 and d = 0.

Batch renormalization's changes more tightly correlate the batch and running sam-469 ple characteristics and have been documented to minimize discrepancy between train and 470 test objectives (Ioffe, 2017). We observed this with our ELBO results, where consistent 471 gaps remained between the train and test $\mathcal{L}[\phi_{(\theta,x)}]$ until we swapped batch normaliza-472 tion for renormalization. Batch renormalization also improves training on low batch sizes 473 (Ioffe, 2017; Summers & Dinneen, 2020), and in our position where variational path sam-474 ples were limited by GPU video memory constraints, renormalization was helpful for de-475 creasing the total number of training iterations we needed for algorithm convergence. 476

With batch (re)normalization layers, $\log q(x)$ accrues log determinant Jacobian summation terms corresponding to those transformations and develops from (55) to become, in the training phase,

$$\log q(x) = \sum_{t=1}^{T} \varphi(z_t) - \sum_{m=1}^{M} \sum_{t=1}^{T} \left[\log \sigma_t^m - \log r_t^m - \log \gamma_t^m + \log \sigma_{\mathcal{S},t}^m \right]$$
(68)

or in the testing phase,

$$\log q(x) = \sum_{t=1}^{T} \varphi(z_t) - \sum_{m=1}^{M} \sum_{t=1}^{T} \left[\log \sigma_t^m - \log \gamma_t^m + \log \sigma_{\mathcal{R},t}^m \right]$$
(69)

where we now take M to mark the total number of layer sets rather than layers as we did before in (55). This assumes that each layer set always includes 1 single affine block and 1 batch renormalization layer. Substituting (68) or (69) into (36) for $\log q(x|\theta;\phi_x)$ leads respectively to our fully decomposed train or test $\mathcal{L}[\phi_{(\theta,x)}]$ calculation unless an optional single softplus transformation is used to ensure constraint of flow output to $\mathbb{R}_{\geq 0}$. In that case, the resulting train $\log q(x)$ is

$$\log q(x) = \sum_{t=1}^{T} \varphi(z_t) - \sum_{m=1}^{M} \sum_{t=1}^{T} \left[\log \sigma_t^m - \log r_t^m - \log \gamma_t^m + \log \sigma_{\mathcal{S},t}^m \right] - \sum_{t=1}^{T} \log(-e^{-x_t} + 1) \quad (70)$$

where x is our terminally transformed random variable following softplus constraint. Setting

$$\lambda_t = \varphi(z_t) - \log(-\mathrm{e}^{-x_t} + 1) - \sum_{m=1}^M \left(\log \sigma_t^m - \log r_t^m - \log \gamma_t^m + \log \sigma_{\mathcal{S},t}^m\right)$$
(71)

$$\log q(x) = \sum_{t=1}^{T} \lambda_t \tag{72}$$

our fully decomposed train $\mathcal{L}[\phi_{(\theta,x)}]$ calculation that we use in each iteration of VI optimization (Algorithm 1) then consolidates from (36) into

$$\mathcal{L}[\phi_{(\theta,x)}] = \mathbb{E}_{q(\theta,x;\phi_{(\theta,x)})} \langle \log p(\theta) + \log p(y_0|x_0,\theta) - \log q(\theta;\phi_{\theta}) \\ + \sum_{i\in\mathfrak{T}} \log p(y_i|x_i,\theta) + \sum_{t=1}^T \left[\log p(x_t|x_{t-1},\theta) - \lambda_t\right] \rangle \quad (73)$$

with softplus flow termination. The test ELBO equation is equivalent except for use of a different λ_t assignment that lacks the log r_t^m term and swaps $\sigma_{\mathcal{S},t}^m$ for $\sigma_{\mathcal{R},t}^m$.

We note that it is not required for the total permutation layers, affine blocks, and 479 batch renormalization layers constituting a neural moving average flow architecture to 480 match in count; we can choose to omit certain layers in a layer set. To slightly reduce 481 the neural network size, we would frequently use 1 less batch renormalization layer than total affine blocks or permutation layers, omitting batch renormalization in the first layer 483 set since we empirically observed little qualitative difference in visual fit quality between 484 running with 3, 4, or 5 batch renormalizations. If the numbers of affine blocks and batch 485 renormalization layers do not match, then the log Jacobian determinant summations in 486 (68) to (71) need to be adjusted accordingly. 487

It is apparent that each layer set of our neural moving average flow corresponds to a matrix of hidden parameters, including affine and batch renormalization parameters, of dimensions [T, h], where h is the count of hidden parameters per layer set. Thus, when conditioning on long, dense T data that is complex in such a manner that would require many layer sets for flow representation, we note that a different choice of marginal variational family for $q(x|\theta)$ aside from the neural moving average flow may be appropriate for minimizing computational expense. Algorithm 1: Synopsis of the operations occurring in each iteration of our soil biogeochemical state space model VI framework

Data: Time series matrix y of soil pool state and other observations, including CO_2 respiration measurements **Result:** $q(\theta, x; \phi_{(\theta,x)})$ corresponding to the $\mathcal{L}[\phi_{(\theta,x)}]$ value at the stoppage of stochastic gradient optimization Define $q(\theta; \phi_{\theta})$ and $q(x|\theta; \phi_x)$; Initialize (ϕ_{θ}, ϕ_x) ; $n \leftarrow \text{total desired training iterations};$ for $i \leftarrow 1$ to n do for $s \leftarrow 1$ to S do Draw $\theta^{(s)} \sim q(\theta; \phi_{\theta});$ Draw $x^{(s)} \sim q(x|\theta; \phi_x)$ transformed from $z^{(s)}$; end Compute $\mathcal{L}[\phi_{(\theta,x)}]$ (or $-\mathcal{L}[\phi_{(\theta,x)}]$ for gradient descent) as per (73); Compute the gradient $\nabla \hat{\mathcal{L}}[\phi_{(\theta,x)}]$ from (49) with automatic differentiation; Update variational parameters $\phi_{(\theta,x)}$ based on the gradient; end

2.7 Kalman smoother validation

When a state space model is linear in drift and its diffusion is stationary and ad-496 ditive, as is the state space model approximation of SCON-C, the posterior density p(x|y)497 can be determined analytically and precisely in closed form with the Kalman smoother 498 algorithm, provided the algorithm is fed the true θ and observation noise (Kalman, 1960; 499 Rauch et al., 1965). Flow VI in contrast can only numerically estimate p(x|y) through 500 a variational approximation, but has the critical advantage of being capable of function-501 ing without exact knowledge of θ given uninformed prior distributions and is able to es-502 timate the joint density $p(x, \theta|y)$ via variational approximations. Thus, comparing a Kalman-503 derived true p(x|y) to a post-optimization $q(x|\theta; \phi_x)$ can be a revealing means of bench-504 marking flow approximation performance and accuracy before applying an architecture 505 with confidence to approximation, optimization, and θ inference of models like SCON-506 SS that cannot be resolved by the smoother. 507

The Kalman smoother procedure is a two part process consisting of a *forward pass* 508 followed by a backward pass. The forward pass computes a "filtering" posterior $p(x_t|y_{0:t})$. 509 which notates the posterior of x_t given observations up to the time indexed by t, going 510 forward in time from $t = \{0, \ldots, T\}$. The backward pass computes a "smoothing" pos-511 terior $p(x_t|y)$, which notates the posterior of x_t given all observations, going backward 512 in time from $t = \{T, \ldots, 0\}$. Reconciling the "filtering" and "smoothing" posteriors pro-513 duces the true p(x|y). A comprehensive explication of Kalman smoothing is available 514 in Särkkä (2013). 515

516

495

2.8 Flow neural network training tuning choices

We settled on using 5 layer sets of permutation, affine, and batch renormalization layers for our neural moving average flow. This offered qualitatively superior fits over flow architectures with lower layer set counts. For inferences of duration T = 5000 with $\Delta t = 1.0$ with 5 layers, maximum training batch size S at 16 GB of VRAM was 31, so we set S = 31. For T = 1000, we used S = 150, though use of smaller S also appeared functional. For T = 5000 inferences we used 120000 non-warmup ELBO training iterations. For T = 1000 inferences we used 60000 non-warmup ELBO training iterations.

With respect to gradient optimizers including AdaMax (Kingma & Ba, 2015), which 525 was the particular optimizer we selected to carry out gradient descent, the *learning rate* 526 is a hyperparameter that scales the objective gradient and in doing so regulates the ex-527 tent to which neural network weights can updated with each training iteration. The learn-528 ing rate can be adjusted over the course of training based on a schedule. It is frequently 529 decayed over the course of training to promote convergence of our objective function to-530 ward a maximum (for gradient ascent) or minimum (for gradient descent) ((You et al., 531 2019)). We chose a step decline schedule for learning rate decay. For our T = 5000 in-532 ferences, we started with a pre-decay ELBO learning rate of 1×10^{-2} and decayed it 533 by a factor of 0.6 every 10000 iterations. For our T = 1000 inferences, we started with 534 a pre-decay learning rate of 4×10^{-3} and decayed it by a factor of 0.6 every 5000 iter-535 ations. 536

We employed *training warmup*, in which we began optimization with a phase of low learning rate at 1×10^{-6} before increasing the rate to its initial pre-decay levels. As has been demonstrated previously (Goyal et al., 2017), we found warmup allowed us to use higher pre-decay learning rates, experience more stable ELBO loss trajectories, and converge to lower average ELBO values over training (Figure S1). We found 5000 warmup iterations to be sufficient for those purposes.

543

2.9 Software and hardware

544 With respect to the computational software and hardware powering the inference operations, our DGP and inference code was developed for a Python 3.9.7 environment 545 distributed by Anaconda (Anaconda Software Distribution, 2020) and used the Numpy 546 1.20.3 (Harris et al., 2020) and PyTorch 1.10.2 (Paszke et al., 2019) software libraries. 547 PyTorch 1.10.2 was compiled with the Nvidia CUDA 10.2 toolkit. The inferences were 548 run on one Nvidia Tesla V100 GPU at a time updated to CUDA version 11.4.0 with a 549 maximum of 16 GB of video random access memory and two Intel Xeon Gold 6148 CPU 550 cores clocked at 2.40 GHz on the University of California, Irvine HPC3 cluster. Our flow 551 VI framework code modules, data-generating notebooks, and synthetic data are avail-552 able via the address https://doi.org/10.5281/zenodo.6969782. 553

The deterministic CON $p(\theta|y)$ posteriors compared with flow VI $q(\theta; \phi_{\theta})$ in Fig-554 ure 5 were estimated using Stan's NUTS algorithm, which is an extension of the Hamil-555 tonian Monte Carlo inference algorithm (Hoffman & Gelman, 2014). Application of Hamil-556 tonian Monte Carlo for data assimilation and inference of SBMs is further described in 557 Xie et al. (2020) and intuition behind the algorithm can be found in Betancourt (2017). 558 The Stan inference was conducted on a 2017 Intel MacBook Pro in an R 4.0.4 environ-559 ment using Stan 2.29.1 (Carpenter et al., 2017) through the CmdStanR interface (Gabry 560 & Cešnovar, 2021). The NUTS simulation ran with 2 chains of 1000 warmup iterations 561 and 5000 sampling iterations each. In our experience, 1000 warmup iterations were suf-562 ficient for locating the bulk of the posterior density. 563

564 **3 Results**

We generated synthetic y of various lengths, dimensions (i.e. whether CO₂ observations were included in addition to state information), and regular observation densities (i.e. whether we observed measurements from our SCON family data generating processes every 1 or 5 hours). We explored the validity of our state space model VI approach for data assimilation and posterior identification of model θ with inferences conditioned on those y. Below, our results suggest the neural moving average flow framework was functional for approximating the SCON family of SDE systems as state space models, fitting y, constraining posteriors, and recovering some true θ values. We also demonstrate subsequently that stochastic gradient optimization in our case was more stable, efficient, and capable at θ identification than an MCMC procedure involving deterministic ODE models adapted from Xie et al. (2020) conditioned on the same y.

576

3.1 Flow-approximated SCON-C converges to fit synthetic data

Following optimization, an SCON-C state space model approximated by our neural moving average flow implementation reasonably assimilated a T = 5000 hour y produced by an SCON-C data-generating process that included CO₂ observations (Figure 4a). The relatively flat $-\mathcal{L}[\phi_{(\theta,x)}]$ trajectory steadily hovering between -1550 and -1600in the latter half of variational training iterations indicates that our flow VI algorithm converged to a local ELBO minimum (Figure S2).

The mean of the marginal posterior density of latent states $q(x|\theta; \phi_x)$ was estimated from 250 x samples drawn from the joint variational density after ELBO training. The mean latent SOC, DOC, and MBC paths and state-derived CO₂ measurements corresponding to the SCON-C flow sit centrally between the y data points and observation noise across the entire time series (Figure 4a). The latent means are able to adhere to many of the sharp peaks and valleys in the dynamics of the data and the flow CO₂ mean was able to reproduce the rapid oscillatory behavior of the observed CO₂ time series.

Upon closer qualitative inspection and comparison to the true latent distribution 590 computed by a Kalman smoother (Figure 4b), we note the presence of visual discrep-591 ancies between the Kalman and flow means and 95% q(x) diffusion distribution inter-592 vals. Firstly, the extent of SOC diffusion noise is substantially underestimated by the 593 flow, which is line with documentation in literature that a mean-field VI approach tends 594 to underestimate posterior uncertainty compared to more complex full-rank approaches (Kucukelbir et al., 2017). For the other two states, DOC and MBC, the extent of dif-596 fusion noise is more consistent to that which is observed in the Kalman output, but the 597 flow DOC and MBC densities and means appear noisier and more uneven than the Kalman 598 means. 599

Still, the flow encouragingly is generally congruous with the true Kalman solution 600 in dynamics. The flow means fall entirely within the bounds of the 95% Kalman diffu-601 sion interval from t = 0 to 500 as can be seen in Figure 4b and we observed for this par-602 ticular optimization that they almost always remain within those Kalman diffusion bounds 603 through the rest of the time series. Also, we see that the CO_2 mean and distribution calculated from the 250 SCON-C state space model x draws closely matches their Kalman 605 counterparts. The ability of the flow to align with the Kalman smoother in latent state 606 densities improves our confidence in the ability of the neural moving average flow to ap-607 proximate systems that are non-linear in diffusion, like SCON-SS. 608

609 610

3.2 SCON-C flow VI marginal θ posteriors indicate appropriate optimization

Beyond fitting data, we needed to ascertain that proper posterior optimization was 611 612 occurring for confidence in inference algorithm function. In our setting, we would expect our posterior densities to at least always be as informed and certain about θ values as 613 our prior densities, not less. With a mean-field logit-normal variational family for $q(\theta; \phi_{\theta})$, 614 evidence of suitable optimization would come in the form of marginal posterior densi-615 ties being narrower than priors to indicate greater certainty after the introduction of in-616 formation from y along with posterior means separating from prior means and approach-617 ing the true θ used by the data-generating process. 618

Figure 5 indicates that valid posterior optimization indeed occurred in our SCON-C state space model inference to support the notion that our flow VI framework was func-



Figure 4. Marginal posterior $q(x|\theta; \phi_x)$ soil pool state means (orange lines) of the SCON-C state space model approximated by the neural moving average flow following VI optimization. The means are estimated from 250 x samples drawn from the optimized joint density $q(\theta, x; \phi_{(\theta,x)})$. The states are in units of mg C g⁻¹ soil. In (a), the trajectories of flowapproximated state means are compared to the synthetic observations an SCON-C T = 5000hour y backgrounded by the 95% interval of the observation noise (blue dots over blue shading). In (b), we zoom into a subset of the above plot from t = 0 to 500 hour and additionally compare the state means and 95% interval of the diffusion distribution of the optimized model to the true posterior means and 95% diffusion noise computed by a Kalman smoother with knowledge of the true θ values.



Figure 5. SCON-C state space model marginal $q(\theta; \phi_{\theta})$ posterior densities following flow VI optimization (orange) compared to mean-field prior $p(\theta)$ densities (blue) and non-parametric CON ODE marginal $p(\theta|y)$ posterior densities estimated with Stan's NUTS algorithm (green). Flow VI and NUTS were conditioned on the same T = 5000 hour y generated by an SCON-C SDE. The true θ values sampled during data generation are marked by vertical dashed gray lines. Being a deterministic ODE system, CON does not have β diffusion θ , so subplots portraying the marginal $q(\theta; \phi_{\theta})$ densities for the SCON-C state space model c_S , c_D , and $c_M \theta$ were not included in this figure due to a lack of comparison.

tional. Almost all the marginal posterior densities narrowed compared to the priors with the information learned from y by the algorithm. Moreover, many of the marginal $q(\theta; \phi_{\theta})$ means drifted closer to the true θ , including the means of u_M , a_{SD} , and $k_{S,ref}$.

We contrasted the flow VI parametric $q(\theta; \phi_{\theta})$ posterior densities to the non-parametric 624 $p(\theta|y)$ posterior densities estimated with an SBM inference framework conditioned on 625 the same T = 5000 SCON-C y that was previously applied in Xie et al. (2020). This prior 626 framework involves Stan's NUTS algorithm and can only infer θ of deterministic mod-627 els, so the CON system that the SCON family was parameterized from served as the ba-628 sis for inference in this approach. With the flexibility and stability afforded by the abil-629 ity of stochastic optimization to adjust for poor initial condition proposals, noisy state 630 path fluctuations, and outlier observations, the flow VI framework expectedly outper-631 formed the deterministic NUTS workflow. The flow VI marginal $q(\theta; \phi_{\theta})$ densities were 632 all-around better informed and identified, exemplified by the subplots corresponding to 633 the u_M , a_{SD} , a_M , $k_{S,ref}$, $k_{D,ref}$, Ea_S , and $Ea_M \theta$ (Figure 5). Moreover, some NUTS pos-634 terior densities, including those corresponding to the a_{MSC} , a_M , and $Ea_S \theta$, consolidated 635 near their lower or upper support bounds, which points to the deterministic model in-636 ference method compensating for its lack of versatility with more extreme θ proposals. 637

Scrutiny of the posterior for the transfer fraction parameter a_{MSC} brings the issue of θ identifiability limitations to our attention. We see that the SCON-C flow VI marginal a_{MSC} posterior density barely budged from the $a_{MSC} p(\theta)$ density post-optimization (Figure 5). For good posterior identifiability, the a_{MSC} posterior should both narrow substantially to signal reduced uncertainty and shift its density peak toward the true a_{MSC} value.

644 645

3.3 Flow VI can effectively assimilate both full and reduced SCON-SS state space approximations

After visually demonstrating the ability of the flow VI framework to optimize $q(\theta, x; \phi_{(\theta,x)})$ through the fitting of the approximated SCON-C state space model to synthetic SCON-C y and the informing and identification of some marginal $q(\theta; \phi_{\theta})$ densities, we proceeded to test if the flow VI approach could similarly function with a moderately more complex model in SCON-SS that is non-linear in diffusion.

Reviewing the fact that the SCON-SS state space model diffusion is not station-651 ary or additive, it was no longer possible for us to validate SCON-SS $q(x|\theta; \phi_x)$ estimated 652 from post-optimization x samples against a true p(x|y) determined by a Kalman smoother. 653 Nonetheless, we observed that flow VI was able to optimize $q(\theta, x; \phi_{(\theta,x)})$ adequately enough 654 to fit the approximated SCON-SS state space model $q(x|\theta; \phi_x)$ means to T = 5000 y gen-655 erated by an SCON-SS SDE. As was the case for the SCON-C flow, the mean latent SOC. 656 DOC, and MBC trajectories of the trained SCON-SS flow traced a central route through 657 the observed state values and diffusion noise (Figure S3). The trajectories were able to 658 follow the peaks and valleys of the state dynamics recorded in y, and the flow CO₂ mean 659 derived from the sampled states tightly replicated the $y \text{ CO}_2$ oscillations. 660

SCON-SS $q(\theta; \phi_{\theta})$ posterior densities were consistent with proper optimization from 661 information learned in y. Juxtaposed with priors $p(\theta)$, marginal $q(\theta; \phi_{\theta})$ densities mostly 662 narrowed and did not move drastically away from their corresponding true θ to inhibit 663 identifiability (Figure S4). There were clear exceptions for the state-scaling diffusion θ 664 posteriors due to reasonable flow neural network approximation error that prompts an 665 overestimate of diffusion noise and, again, for the a_{MSC} posterior. The modest shift of 666 some θ posterior means away from the true θ is counterbalanced by movement of other 667 related θ , like in the circumstance of the Ea_D posterior mean being counterbalanced by 668 Ea_M to satisfy equation (10). 669



Figure 6. Optimized marginal posterior $q(\theta; \phi_{\theta})$ densities (orange) of a reduced SCON-SS model with all but the $k_{i \in S, D, M, \text{ref}}$ decay and state-scaling diffusion θ fixed compared to mean-field priors $p(\theta)$ (blue).

For one more test to corroborate proper functioning of our flow VI framework, we established a reduced SCON-SS model with all θ fixed in value except for the $k_{i \in S, D, M, \text{ref}}$ linear decay and state-scaling β diffusion parameters. Mirroring above procedures, a synthetic $T = 5000 \ y$ was produced by a reduced SCON-SS data-generating process to condition an SCON-SS state space model optimization. For an appropriately behaving inference algorithm, we would expect that removing degrees of freedom should bolster θ identifiability.

⁶⁷⁷ We verified that identifiability was indeed clarified and improved in the remaining ⁶⁷⁸ drift θ (Figure 6). The marginal $k_{S,\text{ref}} q(\theta; \phi_{\theta})$ posterior density was tightly constrained ⁶⁷⁹ right about the true $k_{S,\text{ref}}$ value. The $k_{D,\text{ref}}$ and $k_{M,\text{ref}}$ posterior means did not align ex-⁶⁸⁰ actly with their true θ , but unambiguously offset each other in a manner that plainly ⁶⁸¹ fulfilled (3) and (10).

We were unable to fix the state-scaling diffusion θ without breaking the flow VI frame-682 work, as it became apparent that the algorithm needed to maintain the ability to over-683 estimate diffusion noise to work. This makes intuitive sense as the flow neural network 684 approximation process will always come with some amount of noisy approximation er-685 ror that adds to the base diffusion of the unapproximated system. The algorithm can 686 no longer work if the flow diffusion noise needs to be fixed at about the same level as it 687 is in the unapproximated system, as it leaves no room for the approximation error to over-688 flow into. So, with the diffusion θ left unfixed during VI training, the algorithm once more 689 overestimated their $q(\theta; \phi_{\theta})$ means, but this is not a cause for concern since the discrep-690 ancy can be explained by neural network approximation error. 691

692 693

3.4 Increasing information in y alters SCON posterior certainty and identifiability

The preceding results all involved y that had CO₂ respiration observations included. Inference conditioned on just state observations is also possible and in our experience was able to fit the data well, but it was much less effective for constraining posteriors and identifying θ (Figure S5). Without CO₂ information, marginal $q(\theta; \phi_{\theta})$ posterior densities tended to be wider and less informed and density means were frequently farther away from true θ , as exemplified by panels corresponding to the a_{SD} and Ea_S SCON-C θ in Figure S5.

We separately observed that increasing the amount of information in y by length-701 ening duration T of the time series greatly benefitted posterior identifiability (Figure S6). 702 Alternatively, θ identifiability was boosted without elongating T by bolstering observa-703 tion density. Individually, the two actions trade off between improvements. In compar-704 ison to densifying observations across T = 1000 such that the set of observation indices 705 \mathfrak{N} matches the set of state space model discretization indices N, extending T to 5000 706 more tightly constrained $q(\theta; \phi_{\theta})$ posterior densities for all θ and concentrated a_{SD} , $k_{S,\text{ref}}$, 707 $k_{D,ref}$, Ea_S , and Ea_M posterior means closer to the true θ . 708

However, increasing observation density for T = 1000 data had the benefit of fur-709 ther constraining posterior densities without also enlarging the divergence in identifica-710 tion of the true SCON-C β diffusion θ , c_S , c_D , and c_M by the means of their correspond-711 ing $q(\theta; \phi_{\theta})$ densities. The enlarged divergence and uncertainty of the diffusion θ in the 712 T = 5000 hour inference compared to the T = 1000 inferences is not unexpected. Cu-713 mulative approximation error of state space x trajectories compounds for the flow with 714 greater T in a manner typical to approximation methods. Larger accrued approxima-715 tion error then corresponds to estimation of greater diffusion noise during inference. 716

717 4 Discussion

We developed a stochastic SBM data assimilation and inference framework that is a versatile, stable, and computationally efficient alternative to MCMC approaches assimilating deterministic ODE systems, especially when GPU hardware is available. The framework involves approximation of SBMs as state space models whose state trajectories can be sampled at reduced computational and temporal cost in comparison to SDEs.

In our demonstration, we carried out state space model approximation with a class 723 of normalizing flows called neural moving average flows that successively transition ran-724 dom variables from simpler to more complex distributions with the stacking of neural 725 network layers. We applied this framework to fit approximated representatives of the SCON 726 family of SBMs to synthetic data. Conditioning with synthetic rather than empirical data 727 allowed us to visualize discrepancies between estimated posterior densities, data-generating 728 densities, and true θ values used by the data-generating process for an assessment of frame-729 work performance. 730

Flow-approximated SCON-C state trajectories were able to effectively track state and CO₂ observations after variational optimization and graphically align with the true latent state distributions determined by a Kalman smoother. Following Kalman validation of our SCON-C inference, we then successfully assimilated synthetic observations and estimated posteriors with SCON-SS, which is non-linear in diffusion and modestly more complex than SCON-C.

737 738

4.1 More data promotes model θ identifiability and constraining of posteriors

In terms of implications for experimental work focused on producing data sets suitable for SBM inference and data assimilation, we firstly recommend that CO_2 respiration measurements be collected and included in y. CO_2 information is highly beneficial for informing the posteriors of SBMs like SCON for which CO_2 efflux equations have been established (Figure S5). Additionally, the collection of supplemental measurements, such as radiolabeled C densities linked to SBM pool transfer fraction θ , should further constrain and identify θ posteriors. Our results indicate that just the CO_2 and state observations were not enough to effectively identify the marginal posterior of the SCON MBCto-SOC transfer θ , a_{MSC} (Figure 5, S4, S5, S6).

With respect to a_{MSC} posterior identifiability or lack thereof, inspection of the SCON 748 system drift in (3) and CO_2 efflux rate equation in (10) suggests that the consistent lack 749 of identifiability is not the consequence of a general algorithm issue but instead stems 750 from a dearth of information in y to further constrain the marginal $a_{MSC} q(\theta; \phi_{\theta})$ den-751 sity. The a_{MSC} parameter appears in two terms of (3) that are each the product of four 752 elements, the $a_M \cdot a_{MSC} \cdot k_M \cdot M$ term in the dS equation denoting C mass transfer 753 from the MBC to SOC soil pool and the $a_M \cdot (1 - a_{MSC}) \cdot k_M \cdot M$ term in the dD equa-754 tion denoting C transfer from the MBC to DOC soil pool. The posterior densities of the 755 a_M and $k_M \theta$ in those terms appear in (10) and are accordingly better constrained and 756 identified with CO_2 measurements in y. This is not the case for a_{MSC} , which is not present 757 in (10). Informing of a_{MSC} can thereby only occur through the state measurements in 758 y, and as only one element in the drift product terms, a_{MSC} can take many values be-759 tween its [0,1] support bounds without greatly affecting the products. 760

Furthermore, our results suggests that raising both study time or data collection 761 frequency would improve posterior estimation accuracy and identifiability in our frame-762 work (Figure S6). But, under budget and personnel limitations, empiricists creating in-763 ference data sets should prioritize one or the other depending on the specific SBMs tar-764 geted for data assimilation and their research objectives. For model comparison of naive 765 stochastic SBM parameterizations where the conceived diffusion θ are less biologically 766 meaningful, accumulating approximation error is less of a concern and prioritizing the 767 maximization of T would be reasonable. In scenarios involving SBMs parameterized with 768 more biogeochemically sophisticated β matrices where accurate estimation of system dif-769 fusion θ takes precedence and falsification of specific dynamics is the goal, it would be 770 more important to minimize approximation error with denser observations. 771

4.2 Future work and research directions

Having demonstrated functional flow VI on more compact synthetic data sets, we
highlight some engineering expansions and modifications to our existing framework that
would facilitate efficient SBM inferences conditioned on empirical data sourced from longterm ecological (LTER) experiments like those documented in Melillo et al. (2017) and
Wood et al. (2019). Efficiently scaling to these data sets is a key priority to assimilate
them into SBMs on the scale of hours rather than weeks, as experienced in Xie et al. (2020),
for statistically rigorous head-to-head model comparison and selection.

The T of data sets sourced from LTER experiments can be on the order of 100,000 780 to 200,000 hours, much larger than the peak T = 5000 hour timespan we explored in our 781 study. With the ability of our framework to leverage GPU hardware, our T = 5000 in-782 ferences typically ran between one to two days to ensure convergence, but even more lim-783 iting than time were GPU memory thresholds preventing adequate variational sample 784 sizes with T much longer than 5000. A way forward for conditioning inferences on y with 785 longer T is to avoid simulating state space model x for the entire T at each training it-786 eration, and this can be done in stochastic gradient optimization with the leveraging of 787 788 the mini-batching technique. Under a mini-batching scheme, y is partitioned into smaller sub-sequences y_i during training, where $i \in 1 : \mathcal{B}$ and \mathcal{B} is the total number of parti-789 tions. In each training iteration, a y_i can then be randomly selected for likelihood eval-790 uation such that the SBM only needs to simulate an x_i subsequence for calculation of 791 the optimization objective. Mini-batching is targeted for future incorporation in our frame-792 work, having been demonstrated in recent flow-related machine learning literature in-793 cluding Papamakarios et al. (2021) and Ryder et al. (2021). 794

LTER data sets tend to have constituent observation vectors whose elements greatly vary in information density and measurement intervals due in part to the varying phys⁷⁹⁷ ical practicality associated with the sampling and measurement of different observations. ⁷⁹⁸ Hence, it would also be helpful to engineer our flow to handle irregular and "ragged" ob-⁷⁹⁹ servations. Moreover, alterations can be made to the flow network architecture to en-⁸⁰⁰ able more efficient conditioning on SBM θ values and allow for feature extraction from ⁸⁰¹ additional auxiliary information, such as the time elapsed between observations.

Beyond flow engineering and architecture, other relevant research priorities include 802 the study of less naive SCON treatments for inference use. SCON family representatives 803 that explicitly and mechanistically model system diffusion as a function of the under-804 lying system reaction stoichiometry can be formulated (Golightly & Wilkinson, 2011; Fuchs, 805 2013) and the stability and predictive accuracy associated with different diffusion covari-806 ance structures can be compared. Moreover, stochastic parameterizations of SBMs that 807 simulate mass transfer with Michaelis-Menten dynamics and would be non-linear in drift 808 should be investigated so that their predictive accuracy can be compared to those of lin-809 ear drift models under our VI framework. This would go toward an existing priority in 810 biogeochemistry to examine whether explicit representation of enzyme catalysis in SBMs 811 improves model performance (J. Li et al., 2014; Sulman et al., 2014; Wieder et al., 2015; 812 J. Li et al., 2019; Xie et al., 2020). 813

Application of our VI approach to head-to-head model comparison and selection 814 begets a need for incorporation of goodness-of-fit quantification into our framework. MCMC 815 has access to metrics like the widely application information criterion (Vehtari et al., 2017), 816 leave-one-out cross-validation, and leave-future-out cross-validation (Bürkner et al., 2020) 817 for Bayesian predictive accuracy quantification, but with their established formulations, 818 these metrics cannot be computed under a VI procedure without prohibitive computa-819 820 tional expense (Dao et al., 2022). The development of Bayesian goodness-of-fit metrics for VI is still an open area (Yao et al., 2018; Giordano et al., 2018), but there has been 821 recent work adapting cross-validation for VI that is promising for integration with a state 822 space model inference pipeline (Magnusson et al., 2019; Dao et al., 2022). 823

4.3 Conclusion

824

Going forward, we recommend that inference approaches involving state space model 825 approximation of stochastic SBMs be used in future biogeochemical data assimilation, 826 fitting, and model comparison research in pursuit of superior computational stability, 827 flexibility, and efficiency. SDE systems are far more robust than ODE systems at accom-828 modating prior density, initial condition, and model structure proposals that are incon-829 sistent with the true data generating process (Whitaker, 2016; Wiqvist et al., 2021). Then, 830 state space approximation greatly reduces the burden of sampling SDE model state tra-831 jectories for likelihood evaluation. Rather than integrating an SDE solver \mathcal{S} times at great 832 computational cost with each algorithm training iteration, we can efficiently sample \mathcal{S} 833 paths from the variational approximation in one pass. Additionally, the discrete nature 834 of state space models integrates well with likelihood estimation conditioned on sparsely 835 observed data sets from long term ecological research where fine-grained knowledge of 836 continuous state dynamics of a model are not necessary or useful for the inference al-837 gorithm. State space model discretization can be handled much more coarsely, which fa-838 cilitates more efficient scaling to larger T. 839

Many of the steps of our data assimilation framework are common to those of other Bayesian inference approaches and hence, a wealth of options exist for modification of this approximated SBM inference framework depending on computational resources and desired posterior estimation accuracy. Different non-variational black box inference methods that are compatible with state space approximation of SDEs can be substituted, such as sequential Monte Carlo algorithms (Golightly & Kypraios, 2018), stochastic gradient Hamiltonian Monte Carlo (Chen et al., 2014), stochastic gradient langevin dynamics (Brosse et al., 2018), and stochastic gradient Markov chain Monte Carlo (Aicher et al., 2019; Nemeth & Fearnhead, 2021).

Researchers using variational Bayesian methods for their black box can opt for $q(\theta)$ variational approximations that are more complex than mean-field representation. These include full-rank multivariate logit-normal families in which θ are not assumed to be independent and covariance is established. The full-rank modeling of covariance mitigates underestimation of $q(\theta)$ uncertainty, which is prevalent in mean-field inference, to correspond to wider marginal $q(\theta)$ densities (Kucukelbir et al., 2017; Sujono et al., 2022).

Additionally, when memory availability inhibits the establishment of larger neu-855 ral networks to train affine shift and scale values for long T or when another method is 856 known to be faster and more convenient for approximating a particular SBM class, dif-857 ferent variational families can be used in place of neural moving average flows for rep-858 resentation and optimization of $q(x|\theta)$. These include multivariate normal distributions 859 with specialized covariance structures (Archer et al., 2015), automatic differention VI 860 (Kucukelbir et al., 2017) with Gauss-Markov distributions (Sujono et al., 2022), neural 861 stochastic differential equations (Tzen & Raginsky, 2019; Jia & Benson, 2019; X. Li et 862 al., 2020), and recurrent neural networks (Krishnan et al., 2017; Ryder et al., 2018), among 863 others. Thus, our framework is flexible and can be repurposed as needed for assimilation of different SBMs or Earth system models that vary in complexity and simulation 865 requirements. 866

⁸⁶⁷ 5 Open Research

A repository containing synthetic time series data of soil pool state and CO₂ observations, Python notebooks containing the code for SCON-C and SCON-SS data-generating processes, and Python modules scaffolding the neural moving average flow VI algorithm are available at https://doi.org/10.5281/zenodo.6969782. Stan code for the deterministic CON inference whose results were compared to in Figure 5 is available at https:// doi.org/10.5281/zenodo.6969769.

874 Acknowledgments

We would like to thank Anton Obukhov (ETH Zurich) for his PyTorch truncated normal distribution class that was used in exploratory work that this study was built on and Andrew Golightly (Durham University) for his advice on ODE-to-SDE conversion and reparameterization.

This research was financially supported by the U.S. National Science Foundation (grant no. DEB-1900885 and IIS-1816365), the U.S. Department of Energy (grant no. DE-SC0014374 and DE-SC0020382), and a UC Irvine ICS Exploration Research Award.

The authors affirm that they have no financial conflict of interest.

References

895

896

897

899

904

- Abs, E., Leman, H., & Ferrière, R. (2020). A multi-scale eco-evolutionary model of cooperation reveals how microbial adaptation influences soil decomposition. *Communications Biology*, 3(1). Retrieved from https://doi.org/10.1038/s42003-020-01198-4
 S42003-020-01198-4 doi: 10.1038/s42003-020-01198-4
- Aicher, C., Ma, Y.-A., Foti, N. J., & Fox, E. B. (2019). Stochastic Gradient
 MCMC for State Space Models. SIAM Journal on Mathematics of Data
 Science, 1(3). Retrieved from https://doi.org/10.1137/18M1214780 doi:
 10.1137/18M1214780
- Allison, S. D., Wallenstein, M. D., & Bradford, M. A. (2010). Soil-carbon response to warming dependent on microbial physiology. *Nature Geoscience*, 3(5). Retrieved from https://doi.org/10.1038/ngeo846 doi: 10.1038/ngeo846
 - Anaconda Software Distribution. (2020). Anaconda Inc. Retrieved 2022-06-17, from https://docs.anaconda.com/
 - Archer, E., Park, I. M., Buesing, L., Cunningham, J., & Paninski, L. (2015). Black box variational inference for state space models. arXiv. Retrieved from https://arxiv.org/abs/1511.07367 doi: 10.48550/ARXIV.1511.07367
- Arrhenius, S. (1889). Über die Dissociationswärme und den Einfluss der Temperatur auf den Dissociationsgrad der Elektrolyte. Zeitschrift für Physikalische Chemie, 4U(1). Retrieved from https://doi.org/10.1515/zpch-1889-0408 doi: 10.1515/zpch-1889-0408
 - Betancourt, M. (2017). A Conceptual Introduction to Hamiltonian Monte Carlo. arXiv. Retrieved from http://arxiv.org/abs/1701.02434
- Bjorck, N., Gomes, C. P., Selman, B., & Weinberger, K. Q. (2018). Understanding Batch Normalization. In S. Bengio, H. Wallach, H. Larochelle,
 K. Grauman, N. Cesa-Bianchi, & R. Garnett (Eds.), Advances in Neural Information Processing Systems (Vol. 31). Curran Associates, Inc.
 Retrieved from https://proceedings.neurips.cc/paper/2018/file/
- 911 36072923bfc3cf47745d704feb489480-Paper.pdf
- Blei, D. M., Kucukelbir, A., & McAuliffe, J. D. (2017). Variational Inference: A
 Review for Statisticians. Journal of the American Statistical Association, 112(518). doi: 10.1080/01621459.2017.1285773
- Bradford, M. A., Wood, S. A., Addicott, E. T., Fenichel, E. P., Fields, N., González-Rivero, J., ... Wieder, W. R. (2021). Quantifying microbial control of soil organic matter dynamics at macrosystem scales. *Biogeochemistry*, 156(1).
 Retrieved from https://doi.org/10.1007/s10533-021-00789-5 doi: 10.1007/s10533-021-00789-5
- 920Brosse, N., Durmus, A., & Moulines, E.(2018).The promises and pitfalls921of Stochastic Gradient Langevin Dynamics.In S. Bengio, H. Wallach,922H. Larochelle, K. Grauman, N. Cesa-Bianchi, & R. Garnett (Eds.), Advances923in Neural Information Processing Systems (Vol. 31).Curran Associates,924Inc. Retrieved from https://proceedings.neurips.cc/paper/2018/file/
- 335cd1b90bfa4ee70b39d08a4ae0cf2d-Paper.pdf
 Browning, A. P., Warne, D. J., Burrage, K., Baker, R. E., & Simpson, M. J. (2020).
 Identifiability analysis for stochastic differential equation models in systems
 biology. Journal of The Royal Society Interface, 17(173). Retrieved from
- 929
 https://royalsocietypublishing.org/doi/abs/10.1098/rsif.2020.0652

 930
 doi: 10.1098/rsif.2020.0652
- Bürkner, P.-C., Gabry, J., & Vehtari, A. (2020). Approximate leave-future-out cross validation for Bayesian time series models. Journal of Statistical Computation
 and Simulation, 90(14). Retrieved from https://doi.org/10.1080/00949655
 .2020.1783262 doi: 10.1080/00949655.2020.1783262
- Carpenter, B., Gelman, A., Hoffman, M. D., Lee, D., Goodrich, B., Betancourt, M.,
 Riddell, A. (2017). Stan: A probabilistic programming language. Journal
 of Statistical Software, 76(1). doi: 10.18637/jss.v076.i01

938	Chen, T., Fox, E., & Guestrin, C. (2014). Stochastic Gradient Hamiltonian Monte
939	Carlo. In E. P. Xing & T. Jebara (Eds.), Proceedings of the 31st International
940	Conference on Machine Learning (Vol. 32). Beijing, China: PMLR. Retrieved
941	from https://proceedings.mlr.press/v32/cheni14.html
942	Christensen, R., Johnson, W., Branscum, A., & Hanson, T. E. (2010). Bayesian
943	Ideas and Data Analysis: An Introduction for Scientists and Statisticians (1st
944	ed.). Taylor & Francis.
045	Dao V H Gunawan D Tran M-N Kohn B Hawkins G E & Brown
945	S D (2022) Efficient selection between hierarchical cognitive mod-
940	ols: Cross validation with variational Bayos — <i>Psychological Mathada</i> — doi:
947	10 1027/mot0000458
948	Daunizary I (2017) Carrier and tical ammenimations to statistical memory of
949	Daumzeau, J. (2017). Semi-analytical approximations to statistical moments of
950	https://orgiu.org/obs/1702_00001_doi: 10.48550/ADXIV.1702.00001
951	Dials I. Karaman D. & Danais V. (2015) NICE New linear Index ender the
952	Dinn, L., Krueger, D., & Bengio, Y. (2015). NICE: Non-linear independent
953	Components Estimation. In Y. Bengio & Y. LeCun (Eds.), 3ra Interna-
954	tional Conference on Learning Representations, ICLR 2015, San Diego,
955	CA, USA, May 7-9, 2015, Workshop Track Proceedings. Retrieved from
956	http://arxiv.org/abs/1410.8516
957	Dinh, L., Sohl-Dickstein, J., & Bengio, S. (2017). Density estimation using Real
958	NVP. In 5th International Conference on Learning Representations, ICLR
959	2017, Toulon, France, April 24-26, 2017, Conference Track Proceedings. Open-
960	Review.net. Retrieved from https://openreview.net/forum?id=HkpbnH91x
961	Fuchs, C. (2013). Inference for Diffusion Processes: With Applications in
962	Life Sciences. Berlin, Heidelberg: Springer Berlin Heidelberg. Retrieved
963	from https://doi.org/10.1007/978-3-642-25969-2 doi: 10.1007/
964	978-3-642-25969-2
965	Gabry, J., & Češnovar, R. (2021). CmdStanR. Retrieved 2022-05-18, from https://
966	mc-stan.org/cmdstanr/
967	Gelman, A., Carlin, J., Stern, H., Dunson, D., Vehtari, A., & Rubin, D. (2013).
968	Bayesian Data Analysis (3rd ed.). New York: Taylor & Francis. Retrieved
969	from https://doi.org/10.1201/b16018
970	Geman, S., & Geman, D. (1987). Stochastic Relaxation, Gibbs Distributions,
971	and the Bayesian Restoration of Images. In M. A. Fischler & O. Firschein
972	(Eds.), <i>Readings in Computer Vision</i> . San Francisco (CA): Morgan Kaufmann.
973	Retrieved from https://www.sciencedirect.com/science/article/pii/
974	B978008051581650057X doi: 10.1016/B978-0-08-051581-6.50057-X
975	Georgiou, K., Malhotra, A., Wieder, W. R., Ennis, J. H., Hartman, M. D., Sul-
976	man, B. N., Jackson, R. B. (2021). Divergent controls of soil organic
977	carbon between observations and process-based models. Biogeochemistry.
978	Retrieved from https://doi.org/10.1007/s10533-021-00819-2 doi:
979	10.1007/s10533-021-00819-2
080	Giordano B Broderick T & Jordan M I (2018) Covariances Bobustness and
0.91	Variational Bayes Journal of Machine Learning Research 19(51) Retrieved
082	from http://imlr.org/papers/v19/17-670.html
092	Ciweta M (2020) Bole of litter production and its decomposition and factors
983	affecting the processes in a tropical forest ecosystem: a review <u>Iournal of</u>
904 085	Ecology and Environment 1/(1) Retrieved from https://doi.org/10.1196/
900	$a_{1610-020-0151-2}$ doi: 10.1186/ $a/1610.020.0151.2$
980	Colightly A is Kympion T (2010) Efficient SMC^2 schemes for stochastic himstic
987	models Statistics and Commuting 09(6) Detrieved from https://doi.org/
988	models. Sumstics and Comparing, $20(0)$. Refrieved from https://doi.org/
989	10.1007/S11222-017-9789-8 (0): 10.1007/S11222-017-9789-8
990	CDE personation and In N. D. Lawrence M. Circlard, M. D. H.
991	SDE parameter estimation. In N. D. Lawrence, M. Girolami, M. Rattray, &
992	G. Sangumetti (Eds.), . Cambridge, MA, USA: MIT Press.

Golightly, A., & Wilkinson, D. J. (2006).Bayesian sequential inference for 993 nonlinear multivariate diffusions. Statistics and Computing, 16(4). Re-994 trieved from https://doi.org/10.1007/s11222-006-9392-x doi: 995 10.1007/s11222-006-9392-x 996 Golightly, A., & Wilkinson, D. J. (2011). Bayesian parameter inference for stochas-997 tic biochemical network models using particle Markov chain Monte Carlo. 998 Interface Focus, 1(6). Retrieved from https://pubmed.ncbi.nlm.nih.gov/ 999 23226583 doi: 10.1098/rsfs.2011.0047 1000 Goval, P., Dollár, P., Girshick, R., Noordhuis, P., Wesolowski, L., Kyrola, A., ... 1001 Accurate, Large Minibatch SGD: Training ImageNet in 1 He, K. (2017).1002 Hour. arXiv. Retrieved from https://arxiv.org/abs/1706.02677 doi: 1003 10.48550/ARXIV.1706.02677 1004 Hararuk, O., & Luo, Y. (2014).Improvement of global litter turnover rate pre-1005 dictions using a Bayesian MCMC approach. Ecosphere, 5(12). Retrieved 1006 from https://esajournals.onlinelibrary.wiley.com/doi/abs/10.1890/ 1007 ES14-00092.1 doi: 10.1890/ES14-00092.1 1008 Hararuk, O., Xia, J., & Luo, Y. (2014).Evaluation and improvement of a global 1009 land model against soil carbon data using a Bayesian Markov chain Monte 1010 Carlo method. Journal of Geophysical Research: Biogeosciences, 119(3). 1011 Retrieved from https://agupubs.onlinelibrary.wiley.com/doi/abs/ 1012 10.1002/2013JG002535 doi: 10.1002/2013JG002535 1013 Harris, C. R., Millman, K. J., van der Walt, S. J., Gommers, R., Virtanen, P., 1014 Array programming with Cournapeau, D., ... Oliphant, T. E. (2020).1015 NumPy. Nature, 585 (7825). Retrieved from https://doi.org/10.1038/ 1016 s41586-020-2649-2 doi: 10.1038/s41586-020-2649-2 1017 He, K., Zhang, X., Ren, S., & Sun, J. (2016).Deep Residual Learning for Image 1018 In 2016 IEEE Conference on Computer Vision and Pattern Recognition. 1019 *Recognition (CVPR).* doi: 10.1109/CVPR.2016.90 1020 Hoffman, M. D., & Gelman, A. (2014). The No-U-Turn Sampler: Adaptively Set-1021 ting Path Lengths in Hamiltonian Monte Carlo. Journal of Machine Learning 1022 Research, 15(47). Retrieved from http://jmlr.org/papers/v15/hoffman14a 1023 .html 1024 Ioffe, S. (2017). Batch Renormalization: Towards Reducing Minibatch Dependence 1025 in Batch-Normalized Models. In Proceedings of the 31st International Confer-1026 ence on Neural Information Processing Systems. Red Hook, NY, USA: Curran 1027 Associates Inc. 1028 Ioffe, S., & Szegedy, C. (2015). Batch Normalization: Accelerating Deep Network 1029 Training by Reducing Internal Covariate Shift. In Proceedings of the 32nd 1030 International Conference on International Conference on Machine Learning -1031 Volume 37. JMLR.org. 1032 Jia, J., & Benson, A. R. (2019).Neural Jump Stochastic Differential Equations. 1033 In Proceedings of the 33rd International Conference on Neural Information 1034 Processing Systems. Red Hook, NY, USA: Curran Associates Inc. 1035 Kullback-Leibler Divergence. In M. Lovric (Ed.), Interna-Joyce, J. M. (2011).1036 tional Encyclopedia of Statistical Science. Berlin, Heidelberg: Springer Berlin 1037 Heidelberg. Retrieved from https://doi.org/10.1007/978-3-642-04898-2 1038 _327 doi: 10.1007/978-3-642-04898-2_327 1039 Kalman, R. E. (1960). A New Approach to Linear Filtering and Prediction Prob-1040 lems. Journal of Basic Engineering, 82(1). Retrieved from https://doi.org/ 1041 10.1115/1.3662552 doi: 10.1115/1.3662552 1042 Kingma, D. P., & Ba, J. (2015).Adam: A Method for Stochastic Optimization. 1043 In Y. Bengio & Y. LeCun (Eds.), 3rd International Conference on Learning 1044 Representations, ICLR 2015, San Diego, CA, USA, May 7-9, 2015, Conference 1045 Track Proceedings. Retrieved from http://arxiv.org/abs/1412.6980

1046

1047	Kingma, D. P., Salimans, T., Jozefowicz, R., Chen, X., Sutskever, I., & Welling,
1048	M. (2016). Improved Variational Inference with Inverse Autoregressive Flow.
1049	In Proceedings of the 30th International Conference on Neural Information
1050	Processing Systems. Red Hook, NY, USA: Curran Associates Inc.
1051	Kingma, D. P., & Welling, M. (2014). Auto-Encoding Variational Bayes. In Y. Ben-
1052	gio & Y. LeCun (Eds.), 2nd international conference on learning represen-
1053	tations. ICLR 2014. banff. ab. canada. april 14-16. 2014. conference track
1054	proceedings. Retrieved from http://arxiv.org/abs/1312.6114
1055	Kobyzev I Prince S & Brubaker M (2020) Normalizing Flows: An Introduc-
1055	tion and Review of Current Methods IEEE Transactions on Pattern Analy-
1050	sis and Machine Intelligence Retrieved from http://dx.doi.org/10.1109/
1057	TDAMI 2020 2002024 doi: 10.1100/trami 2020.2002024
1058	Kwishnan D. C. Shalit II. & Sontag D. (2017). Structured Inference Networks for
1059	Nonlinean State Space Models In <i>Proceedings of the Thirty First AAAL Con</i>
1060	former and Artificial Intelligence AAAI Dress
1061	<i>Jerence on Artificial Intelligence</i> . AAAI Press.
1062	Kucukelbir, A., Iran, D., Ranganath, R., Gelman, A., & Blei, D. M. (2017). Au-
1063	tomatic Differentiation Variational Inference. Journal of Machine Learning Re-
1064	search, 18(1), 430-474.
1065	Kullback, S., & Leibler, R. A. (1951). On Information and Sufficiency. <i>The Annals</i>
1066	of Mathematical Statistics, 22(1). Retrieved from https://doi.org/10.1214/
1067	aoms/1177729694 doi: 10.1214/aoms/1177729694
1068	Lee, H., Pabbaraju, C., Sevekari, A., & Risteski, A. (2021). Universal Approx-
1069	imation for Log-concave Distributions using Well-conditioned Normalizing
1070	<i>Flows.</i> arXiv. Retrieved from https://arxiv.org/abs/2107.02951 doi:
1071	10.48550/ARXIV.2107.02951
1072	Li, J., Wang, G., Allison, S. D., Mayes, M. A., & Luo, Y. (2014). Soil carbon sen-
1073	sitivity to temperature and carbon use efficiency compared across microbial-
1074	ecosystem models of varying complexity. $Biogeochemistry, 119(1-3)$. doi:
1075	10.1007/s10533-013-9948-8
1076	Li, J., Wang, G., Maves, M. A., Allison, S. D., Frev, S. D., Shi, Z., Melillo,
1077	J. M. (2019). Reduced carbon use efficiency and increased microbial
1078	
	turnover with soil warming. Global Change Biology, 25(3). Retrieved from
1079	turnover with soil warming. <i>Global Change Biology</i> , 25(3). Retrieved from https://onlinelibrary.wiley.com/doi/abs/10.1111/gcb.14517 doi:
1079	turnover with soil warming. Global Change Biology, 25(3). Retrieved from https://onlinelibrary.wiley.com/doi/abs/10.1111/gcb.14517 doi: 10.1111/gcb.14517
1079 1080	<pre>turnover with soil warming. Global Change Biology, 25(3). Retrieved from https://onlinelibrary.wiley.com/doi/abs/10.1111/gcb.14517 doi: 10.1111/gcb.14517 Li X Wong T-K L, Chen B, T O, & Duvenaud D, (2020). Scalable Gra-</pre>
1079 1080 1081	 turnover with soil warming. Global Change Biology, 25(3). Retrieved from https://onlinelibrary.wiley.com/doi/abs/10.1111/gcb.14517 doi: 10.1111/gcb.14517 Li, X., Wong, TK. L., Chen, R. T. Q., & Duvenaud, D. (2020). Scalable Gradients for Stochastic Differential Equations In S. Chiappa & R. Calan-
1079 1080 1081 1082	 turnover with soil warming. Global Change Biology, 25(3). Retrieved from https://onlinelibrary.wiley.com/doi/abs/10.1111/gcb.14517 doi: 10.1111/gcb.14517 Li, X., Wong, TK. L., Chen, R. T. Q., & Duvenaud, D. (2020). Scalable Gra- dients for Stochastic Differential Equations. In S. Chiappa & R. Calan- dra (Eds.). Proceedings of the Twenty Third International Conference on
1079 1080 1081 1082 1083	 turnover with soil warming. Global Change Biology, 25(3). Retrieved from https://onlinelibrary.wiley.com/doi/abs/10.1111/gcb.14517 doi: 10.1111/gcb.14517 Li, X., Wong, TK. L., Chen, R. T. Q., & Duvenaud, D. (2020). Scalable Gradients for Stochastic Differential Equations. In S. Chiappa & R. Calandra (Eds.), Proceedings of the Twenty Third International Conference on Artificial Intelligence and Statistica (Vol. 108) PML R. Retrieved from the statistical form.
1079 1080 1081 1082 1083 1084	 turnover with soil warming. Global Change Biology, 25(3). Retrieved from https://onlinelibrary.wiley.com/doi/abs/10.1111/gcb.14517 doi: 10.1111/gcb.14517 Li, X., Wong, TK. L., Chen, R. T. Q., & Duvenaud, D. (2020). Scalable Gradients for Stochastic Differential Equations. In S. Chiappa & R. Calandra (Eds.), Proceedings of the Twenty Third International Conference on Artificial Intelligence and Statistics (Vol. 108). PMLR. Retrieved from https://proceedings.mlr.proce/u108/li20i.html
1079 1080 1081 1082 1083 1084 1085	 turnover with soil warming. Global Change Biology, 25(3). Retrieved from https://onlinelibrary.wiley.com/doi/abs/10.1111/gcb.14517 doi: 10.1111/gcb.14517 Li, X., Wong, TK. L., Chen, R. T. Q., & Duvenaud, D. (2020). Scalable Gra- dients for Stochastic Differential Equations. In S. Chiappa & R. Calan- dra (Eds.), Proceedings of the Twenty Third International Conference on Artificial Intelligence and Statistics (Vol. 108). PMLR. Retrieved from https://proceedings.mlr.press/v108/li20i.html
1079 1080 1081 1082 1083 1084 1085 1086	 turnover with soil warming. Global Change Biology, 25(3). Retrieved from https://onlinelibrary.wiley.com/doi/abs/10.1111/gcb.14517 doi: 10.1111/gcb.14517 Li, X., Wong, TK. L., Chen, R. T. Q., & Duvenaud, D. (2020). Scalable Gra- dients for Stochastic Differential Equations. In S. Chiappa & R. Calan- dra (Eds.), Proceedings of the Twenty Third International Conference on Artificial Intelligence and Statistics (Vol. 108). PMLR. Retrieved from https://proceedings.mlr.press/v108/li20i.html Luo, Y., Ahlström, A., Allison, S. D., Batjes, N. H., Brovkin, V., Carvalhais, N., They T. (2016). Termend means realistic presistions of axil earlier dynamics
1079 1080 1081 1082 1083 1084 1085 1086 1087	 turnover with soil warming. Global Change Biology, 25(3). Retrieved from https://onlinelibrary.wiley.com/doi/abs/10.1111/gcb.14517 doi: 10.1111/gcb.14517 Li, X., Wong, TK. L., Chen, R. T. Q., & Duvenaud, D. (2020). Scalable Gra- dients for Stochastic Differential Equations. In S. Chiappa & R. Calan- dra (Eds.), Proceedings of the Twenty Third International Conference on Artificial Intelligence and Statistics (Vol. 108). PMLR. Retrieved from https://proceedings.mlr.press/v108/li20i.html Luo, Y., Ahlström, A., Allison, S. D., Batjes, N. H., Brovkin, V., Carvalhais, N., Zhou, T. (2016). Toward more realistic projections of soil carbon dynamics https://protection.com/doi/abs/10.1111/gcb.14517
1079 1080 1081 1082 1083 1084 1085 1086 1087 1088	 turnover with soil warming. Global Change Biology, 25(3). Retrieved from https://onlinelibrary.wiley.com/doi/abs/10.1111/gcb.14517 doi: 10.1111/gcb.14517 Li, X., Wong, TK. L., Chen, R. T. Q., & Duvenaud, D. (2020). Scalable Gra- dients for Stochastic Differential Equations. In S. Chiappa & R. Calan- dra (Eds.), Proceedings of the Twenty Third International Conference on Artificial Intelligence and Statistics (Vol. 108). PMLR. Retrieved from https://proceedings.mlr.press/v108/li20i.html Luo, Y., Ahlström, A., Allison, S. D., Batjes, N. H., Brovkin, V., Carvalhais, N., Zhou, T. (2016). Toward more realistic projections of soil carbon dynamics by Earth system models. Global Biogeochemical Cycles, 30(1). Retrieved
1079 1080 1081 1082 1083 1084 1085 1086 1087 1088 1088	 turnover with soil warming. Global Change Biology, 25(3). Retrieved from https://onlinelibrary.wiley.com/doi/abs/10.1111/gcb.14517 doi: 10.1111/gcb.14517 Li, X., Wong, TK. L., Chen, R. T. Q., & Duvenaud, D. (2020). Scalable Gra- dients for Stochastic Differential Equations. In S. Chiappa & R. Calan- dra (Eds.), Proceedings of the Twenty Third International Conference on Artificial Intelligence and Statistics (Vol. 108). PMLR. Retrieved from https://proceedings.mlr.press/v108/li20i.html Luo, Y., Ahlström, A., Allison, S. D., Batjes, N. H., Brovkin, V., Carvalhais, N., Zhou, T. (2016). Toward more realistic projections of soil carbon dynamics by Earth system models. Global Biogeochemical Cycles, 30(1). Retrieved from https://agupubs.onlinelibrary.wiley.com/doi/abs/10.1002/
1079 1080 1081 1082 1083 1084 1085 1086 1087 1088 1089 1089	 turnover with soil warming. Global Change Biology, 25(3). Retrieved from https://onlinelibrary.wiley.com/doi/abs/10.1111/gcb.14517 doi: 10.1111/gcb.14517 Li, X., Wong, TK. L., Chen, R. T. Q., & Duvenaud, D. (2020). Scalable Gra- dients for Stochastic Differential Equations. In S. Chiappa & R. Calan- dra (Eds.), Proceedings of the Twenty Third International Conference on Artificial Intelligence and Statistics (Vol. 108). PMLR. Retrieved from https://proceedings.mlr.press/v108/li20i.html Luo, Y., Ahlström, A., Allison, S. D., Batjes, N. H., Brovkin, V., Carvalhais, N., Zhou, T. (2016). Toward more realistic projections of soil carbon dynamics by Earth system models. Global Biogeochemical Cycles, 30(1). Retrieved from https://agupubs.onlinelibrary.wiley.com/doi/abs/10.1002/ 2015GB005239 doi: 10.1002/2015GB005239
1079 1080 1081 1082 1083 1084 1085 1086 1087 1088 1089 1099	 turnover with soil warming. Global Change Biology, 25(3). Retrieved from https://onlinelibrary.wiley.com/doi/abs/10.1111/gcb.14517 doi: 10.1111/gcb.14517 Li, X., Wong, TK. L., Chen, R. T. Q., & Duvenaud, D. (2020). Scalable Gra- dients for Stochastic Differential Equations. In S. Chiappa & R. Calan- dra (Eds.), Proceedings of the Twenty Third International Conference on Artificial Intelligence and Statistics (Vol. 108). PMLR. Retrieved from https://proceedings.mlr.press/v108/li120i.html Luo, Y., Ahlström, A., Allison, S. D., Batjes, N. H., Brovkin, V., Carvalhais, N., Zhou, T. (2016). Toward more realistic projections of soil carbon dynamics by Earth system models. Global Biogeochemical Cycles, 30(1). Retrieved from https://agupubs.onlinelibrary.wiley.com/doi/abs/10.1002/ 2015GB005239 doi: 10.1002/2015GB005239 Magnusson, M., Andersen, M., Jonasson, J., & Vehtari, A. (2019). Bayesian leave-
1079 1080 1081 1082 1083 1084 1085 1086 1087 1088 1089 1090 1091	 turnover with soil warming. Global Change Biology, 25(3). Retrieved from https://onlinelibrary.wiley.com/doi/abs/10.1111/gcb.14517 doi: 10.1111/gcb.14517 Li, X., Wong, TK. L., Chen, R. T. Q., & Duvenaud, D. (2020). Scalable Gra- dients for Stochastic Differential Equations. In S. Chiappa & R. Calan- dra (Eds.), Proceedings of the Twenty Third International Conference on Artificial Intelligence and Statistics (Vol. 108). PMLR. Retrieved from https://proceedings.mlr.press/v108/li20i.html Luo, Y., Ahlström, A., Allison, S. D., Batjes, N. H., Brovkin, V., Carvalhais, N., Zhou, T. (2016). Toward more realistic projections of soil carbon dynamics by Earth system models. Global Biogeochemical Cycles, 30(1). Retrieved from https://agupubs.onlinelibrary.wiley.com/doi/abs/10.1002/ 2015GB005239 doi: 10.1002/2015GB005239 Magnusson, M., Andersen, M., Jonasson, J., & Vehtari, A. (2019). Bayesian leave- one-out cross-validation for large data. In K. Chaudhuri & R. Salakhutdinov (Ed.). B. difference in the conduction of the store in the store
1079 1080 1081 1082 1083 1084 1085 1086 1087 1088 1089 1090 1091 1092	 turnover with soil warming. Global Change Biology, 25(3). Retrieved from https://onlinelibrary.wiley.com/doi/abs/10.1111/gcb.14517 doi: 10.1111/gcb.14517 Li, X., Wong, TK. L., Chen, R. T. Q., & Duvenaud, D. (2020). Scalable Gra- dients for Stochastic Differential Equations. In S. Chiappa & R. Calan- dra (Eds.), Proceedings of the Twenty Third International Conference on Artificial Intelligence and Statistics (Vol. 108). PMLR. Retrieved from https://proceedings.mlr.press/v108/li20i.html Luo, Y., Ahlström, A., Allison, S. D., Batjes, N. H., Brovkin, V., Carvalhais, N., Zhou, T. (2016). Toward more realistic projections of soil carbon dynamics by Earth system models. Global Biogeochemical Cycles, 30(1). Retrieved from https://agupubs.onlinelibrary.wiley.com/doi/abs/10.1002/ 2015GB005239 doi: 10.1002/2015GB005239 Magnusson, M., Andersen, M., Jonasson, J., & Vehtari, A. (2019). Bayesian leave- one-out cross-validation for large data. In K. Chaudhuri & R. Salakhutdinov (Eds.), Proceedings of the 36th International Conference on Machine Learning
1079 1080 1081 1082 1083 1084 1085 1086 1087 1088 1089 1090 1091 1092 1093 1094	 turnover with soil warming. Global Change Biology, 25(3). Retrieved from https://onlinelibrary.wiley.com/doi/abs/10.1111/gcb.14517 doi: 10.1111/gcb.14517 Li, X., Wong, TK. L., Chen, R. T. Q., & Duvenaud, D. (2020). Scalable Gra- dients for Stochastic Differential Equations. In S. Chiappa & R. Calan- dra (Eds.), Proceedings of the Twenty Third International Conference on Artificial Intelligence and Statistics (Vol. 108). PMLR. Retrieved from https://proceedings.mlr.press/v108/li20i.html Luo, Y., Ahlström, A., Allison, S. D., Batjes, N. H., Brovkin, V., Carvalhais, N., Zhou, T. (2016). Toward more realistic projections of soil carbon dynamics by Earth system models. Global Biogeochemical Cycles, 30(1). Retrieved from https://agupubs.onlinelibrary.wiley.com/doi/abs/10.1002/ 2015GB005239 doi: 10.1002/2015GB005239 Magnusson, M., Andersen, M., Jonasson, J., & Vehtari, A. (2019). Bayesian leave- one-out cross-validation for large data. In K. Chaudhuri & R. Salakhutdinov (Eds.), Proceedings of the 36th International Conference on Machine Learning (Vol. 97). PMLR. Retrieved from https://proceedings.mlr.press/v97/
1079 1080 1081 1082 1083 1084 1085 1086 1087 1088 1089 1090 1091 1092 1093 1094 1095	 turnover with soil warming. Global Change Biology, 25(3). Retrieved from https://onlinelibrary.wiley.com/doi/abs/10.1111/gcb.14517 doi: 10.1111/gcb.14517 Li, X., Wong, TK. L., Chen, R. T. Q., & Duvenaud, D. (2020). Scalable Gra- dients for Stochastic Differential Equations. In S. Chiappa & R. Calan- dra (Eds.), Proceedings of the Twenty Third International Conference on Artificial Intelligence and Statistics (Vol. 108). PMLR. Retrieved from https://proceedings.mlr.press/v108/li20i.html Luo, Y., Ahlström, A., Allison, S. D., Batjes, N. H., Brovkin, V., Carvalhais, N., Zhou, T. (2016). Toward more realistic projections of soil carbon dynamics by Earth system models. Global Biogeochemical Cycles, 30(1). Retrieved from https://agupubs.onlinelibrary.wiley.com/doi/abs/10.1002/ 2015GB005239 doi: 10.1002/2015GB005239 Magnusson, M., Andersen, M., Jonasson, J., & Vehtari, A. (2019). Bayesian leave- one-out cross-validation for large data. In K. Chaudhuri & R. Salakhutdinov (Eds.), Proceedings of the 36th International Conference on Machine Learning (Vol. 97). PMLR. Retrieved from https://proceedings.mlr.press/v97/ magnusson19a.html
1079 1080 1081 1082 1083 1084 1085 1086 1087 1088 1089 1090 1091 1092 1093 1094 1095 1096	 turnover with soil warming. Global Change Biology, 25(3). Retrieved from https://onlinelibrary.wiley.com/doi/abs/10.1111/gcb.14517 doi: 10.1111/gcb.14517 Li, X., Wong, TK. L., Chen, R. T. Q., & Duvenaud, D. (2020). Scalable Gra- dients for Stochastic Differential Equations. In S. Chiappa & R. Calan- dra (Eds.), Proceedings of the Twenty Third International Conference on Artificial Intelligence and Statistics (Vol. 108). PMLR. Retrieved from https://proceedings.mlr.press/v108/li20i.html Luo, Y., Ahlström, A., Allison, S. D., Batjes, N. H., Brovkin, V., Carvalhais, N., Zhou, T. (2016). Toward more realistic projections of soil carbon dynamics by Earth system models. Global Biogeochemical Cycles, 30(1). Retrieved from https://agupubs.onlinelibrary.wiley.com/doi/abs/10.1002/ 2015GB005239 doi: 10.1002/2015GB005239 Magnusson, M., Andersen, M., Jonasson, J., & Vehtari, A. (2019). Bayesian leave- one-out cross-validation for large data. In K. Chaudhuri & R. Salakhutdinov (Eds.), Proceedings of the 36th International Conference on Machine Learning (Vol. 97). PMLR. Retrieved from https://proceedings.mlr.press/v97/ magnusson19a.html Manzoni, S., & Porporato, A. (2009). Soil carbon and nitrogen mineralization: The-
1079 1080 1081 1082 1083 1084 1085 1086 1087 1088 1089 1090 1091 1092 1093 1094 1095 1096 1097	 turnover with soil warming. Global Change Biology, 25(3). Retrieved from https://onlinelibrary.wiley.com/doi/abs/10.1111/gcb.14517 doi: 10.1111/gcb.14517 Li, X., Wong, TK. L., Chen, R. T. Q., & Duvenaud, D. (2020). Scalable Gradients for Stochastic Differential Equations. In S. Chiappa & R. Calandra (Eds.), Proceedings of the Twenty Third International Conference on Artificial Intelligence and Statistics (Vol. 108). PMLR. Retrieved from https://proceedings.mlr.press/v108/li20i.html Luo, Y., Ahlström, A., Allison, S. D., Batjes, N. H., Brovkin, V., Carvalhais, N., Zhou, T. (2016). Toward more realistic projections of soil carbon dynamics by Earth system models. Global Biogeochemical Cycles, 30(1). Retrieved from https://agupubs.onlinelibrary.wiley.com/doi/abs/10.1002/2015GB005239 Magnusson, M., Andersen, M., Jonasson, J., & Vehtari, A. (2019). Bayesian leave-one-out cross-validation for large data. In K. Chaudhuri & R. Salakhutdinov (Eds.), Proceedings of the 36th International Conference on Machine Learning (Vol. 97). PMLR. Retrieved from https://proceedings.mlr.press/v97/magnusson19a.html Manzoni, S., & Porporato, A. (2009). Soil carbon and nitrogen mineralization: Theory and models across scales. Soil Biology and Biochemistry, 41(7). Retrieved
1079 1080 1081 1082 1083 1084 1085 1086 1087 1088 1089 1090 1091 1092 1093 1094 1095 1096	 turnover with soil warming. Global Change Biology, 25(3). Retrieved from https://onlinelibrary.wiley.com/doi/abs/10.1111/gcb.14517 doi: 10.1111/gcb.14517 Li, X., Wong, TK. L., Chen, R. T. Q., & Duvenaud, D. (2020). Scalable Gradients for Stochastic Differential Equations. In S. Chiappa & R. Calandra (Eds.), Proceedings of the Twenty Third International Conference on Artificial Intelligence and Statistics (Vol. 108). PMLR. Retrieved from https://proceedings.mlr.press/v108/1i20i.html Luo, Y., Ahlström, A., Allison, S. D., Batjes, N. H., Brovkin, V., Carvalhais, N., Zhou, T. (2016). Toward more realistic projections of soil carbon dynamics by Earth system models. Global Biogeochemical Cycles, 30(1). Retrieved from https://agupubs.onlinelibrary.wiley.com/doi/abs/10.1002/2015GB005239 Magnusson, M., Andersen, M., Jonasson, J., & Vehtari, A. (2019). Bayesian leave-one-out cross-validation for large data. In K. Chaudhuri & R. Salakhutdinov (Eds.), Proceedings of the 36th International Conference on Machine Learning (Vol. 97). PMLR. Retrieved from https://proceedings.mlr.press/v97/magnusson19a.html Manzoni, S., & Porporato, A. (2009). Soil carbon and nitrogen mineralization: Theory and models across scales. Soil Biology and Biochemistry, 41(7). Retrieved from http://dx.doi.org/10.1016/j.soilbio.2009.02.031
1079 1080 1081 1082 1083 1084 1085 1086 1087 1088 1089 1090 1091 1092 1093 1094 1095 1096 1097 1098	 turnover with soil warming. Global Change Biology, 25(3). Retrieved from https://onlinelibrary.wiley.com/doi/abs/10.1111/gcb.14517 doi: 10.1111/gcb.14517 doi: 10.1111/gcb.14517 Li, X., Wong, TK. L., Chen, R. T. Q., & Duvenaud, D. (2020). Scalable Gradients for Stochastic Differential Equations. In S. Chiappa & R. Calandra (Eds.), Proceedings of the Twenty Third International Conference on Artificial Intelligence and Statistics (Vol. 108). PMLR. Retrieved from https://proceedings.mlr.press/v108/1i20i.html Luo, Y., Ahlström, A., Allison, S. D., Batjes, N. H., Brovkin, V., Carvalhais, N., Zhou, T. (2016). Toward more realistic projections of soil carbon dynamics by Earth system models. Global Biogeochemical Cycles, 30(1). Retrieved from https://agupubs.onlinelibrary.wiley.com/doi/abs/10.1002/2015GB005239 Magnusson, M., Andersen, M., Jonasson, J., & Vehtari, A. (2019). Bayesian leave-one-out cross-validation for large data. In K. Chaudhuri & R. Salakhutdinov (Eds.), Proceedings of the 36th International Conference on Machine Learning (Vol. 97). PMLR. Retrieved from https://proceedings.mlr.press/v97/magnusson19a.html Manzoni, S., & Porporato, A. (2009). Soil carbon and nitrogen mineralization: Theory and models across scales. Soil Biology and Biochemistry, 41(7). Retrieved from http://dx.doi.org/10.1016/j.soilbio.2009.02.031
1079 1080 1081 1082 1083 1084 1085 1086 1087 1098 1090 1091 1092 1093 1094 1095 1096 1097 1098 1099 1099 1100	 turnover with soil warming. Global Change Biology, 25(3). Retrieved from https://onlinelibrary.wiley.com/doi/abs/10.1111/gcb.14517 doi: 10.1111/gcb.14517 doi: 10.1111/gcb.14517 Li, X., Wong, TK. L., Chen, R. T. Q., & Duvenaud, D. (2020). Scalable Gradients for Stochastic Differential Equations. In S. Chiappa & R. Calandra (Eds.), Proceedings of the Twenty Third International Conference on Artificial Intelligence and Statistics (Vol. 108). PMLR. Retrieved from https://proceedings.mlr.press/v108/1i20i.html Luo, Y., Ahlström, A., Allison, S. D., Batjes, N. H., Brovkin, V., Carvalhais, N., Zhou, T. (2016). Toward more realistic projections of soil carbon dynamics by Earth system models. Global Biogeochemical Cycles, 30(1). Retrieved from https://agupubs.onlinelibrary.wiley.com/doi/abs/10.1002/2015GB005239 Magnusson, M., Andersen, M., Jonasson, J., & Vehtari, A. (2019). Bayesian leave-one-out cross-validation for large data. In K. Chaudhuri & R. Salakhutdinov (Eds.), Proceedings of the 36th International Conference on Machine Learning (Vol. 97). PMLR. Retrieved from https://proceedings.mlr.press/v97/magnusson19a.html Manzoni, S., & Porporato, A. (2009). Soil carbon and nitrogen mineralization: Theory and models across scales. Soil Biology and Biochemistry, 41(7). Retrieved from http://dx.doi.org/10.1016/j.soilbio.2009.02.031 Maruyama, G. (1955). Continuous Markov processes and stochastic equations. Ren-

1102	.org/10.1007/BF02846028 doi: 10.1007/BF02846028
1103	McElreath, R. (2020). Statistical Rethinking: A Bayesian Course with Examples
1104	in R and Stan. 2nd Edition (2nd ed.). CRC Press. Retrieved from http://
1105	xcelab.net/rm/statistical-rethinking/
1106	Melillo, J. M., Frey, S. D., DeAngelis, K. M., Werner, W. J., Bernard, M. J., Bowles,
1107	F. P Grandy, A. S. (2017). Long-term pattern and magnitude of soil car-
1108	bon feedback to the climate system in a warming world. Science, 358(6359).
1109	doi: 10.1126/science.aan2874
1110	Nemeth, C., & Fearnhead, P. (2021). Stochastic Gradient Markov Chain Monte
1111	Carlo, Journal of the American Statistical Association, 116(533). Re-
1112	trieved from https://doi.org/10.1080/01621459.2020.1847120 doi:
1113	10.1080/01621459.2020.1847120
1114	O'Neill B C Kriegler E Ebi K L Kemp-Benedict E Biahi K Both-
1114	man D S Solecki W (2017) The roads ahead: Narratives for
1115	shared socioeconomic pathways describing world futures in the 21st cen-
1110	tury Global Environmental Change 42 Retrieved from https://
1117	www.sciencedirect.com/science/article/pii/S0959378015000060 doi:
1110	10 1016/i gloenycha 2015 01 004
1120	Panamakarios G. Nalisnick F. Bezende D. I. Mohamed S. & Lakshmi-
1120	narayanan B (2021) Normalizing Flows for Probabilistic Modeling and
1121	Inference Internal of Machine Learning Research 22(57) Betrieved from
1122	http://imlr.org/papers/u22/19-1028 html
1125	Papamakarios C. Paylakou T. & Murray I. (2017) Masked Autoregressive Flow
1124	for Density Estimation In Proceedings of the 21st International Conference on
1125	Neural Information Processing Systems Red Hook NV USA: Curren Asso
1126	cistes Inc
1127	Dearly A Creas S Massa E Lover A Predbury I Chappen C Chintele
1128	S (2010) DyTorch: An Imporative Style High Performance Deep Learning
1129	Library In H Wallach H Larochalle A Beygelzimer E d'Alché-Buc E Fox
1130	library. In H. Wahach, H. Larochene, A. Deygelanner, F. a Mene-Duc, E. Fox, l. R. Carnett (Eds.) Advances in Neural Information Processing Systems 29
1131	Curran Associates Inc. Betrieved from http://napers.neurins.cc/naper/
1132	9015-nutorch-an-imporative-style-high-performance-deen-learning
1124	-library ndf
1134	Perez-Cruz F (2008) Kullback-Leibler divergence estimation of continuous distri-
1135	butions In 2008 IEEE International Sumposium on Information Theory doi:
1130	10 1109/ISIT 2008 4595271
1157	Plummer M (2003) IACS: A Program for Analysis of Bayesian Craphical Models
1138	using Cibbs Sampling In K Hornik E Loisch & A Zoilois (Eds.) Proceed
1139	ing of the 3rd International Workshop on Distributed Statistical Computing
1140	(DSC 2003) (Vol. 3) Vienna Austria
1141	Baczka B Hoar T I Duarte H F Fox A M Anderson I I. Bowling D B
1142	lz Lin I C (2021) Improving CLM5.0 Biomass and Carbon Exchange
1143	Across the Western United States Using a Data Assimilation System — <i>Journal</i>
1144	of Advances in Modeling Earth Systems 13(7) Betrieved from https://
1145	agunubs onlinelibrary viley com/doi/abs/10 1029/2020MS002421
1140	(e2020MS002421 2020MS002421) doi: 10.1029/2020MS002421
1147	Rauch H E Tung E & Striebel C T (1965) Maximum likelihood estimates of
1148	linear dynamic systems AIAA journal 3(8)
1150	Beid N M (2015) Approximate likelihoode In I Cuo & 7 M Ma (Edg.)
1150	Proceedings of the 8th International Congress on Industrial and Annlied Math
1151	emotics Beijing Chine: Higher Education Pross
1152	Duiz E B Titaing DC ALLER M & Bloi D (2016) The Conoralized Denomena
1153	terization Gradient In D. Lee, M. Sugiverna, H. Luyburg, I. Gunon, & D. Car
1154	nett (Eds.) Advances in Neural Information Processing Systems (Vol. 20)
1155	Curren Associates Inc. Retrieved from https://proceedings.neuring.co/
1120	Curran Abboliates, inc. Iterreven fom https://proceeuings.heurips.cc/

1157	paper/2016/file/f718499c1c8cef6730f9fd03c8125cab-Paper.pdf
1158	Ryder, T., Golightly, A., McGough, A. S., & Prangle, D. (2018). Black-Box Vari-
1159	ational Inference for Stochastic Differential Equations. In J. Dy & A. Krause
1160	(Eds.). Proceedings of the 35th International Conference on Machine Learning
1161	(Vol. 80). PMLR. Retrieved from https://proceedings.mlr.press/v80/
1162	rvder18a.html
1163	Byder T Prangle D Golightly A & Matthews I (2021) The neural moving
1164	average model for scalable variational inference of state space models. In C de
1165	Campos & M H Maathuis (Eds.) Proceedings of the Thirty-Seventh Confer-
1166	ence on Uncertainty in Artificial Intelligence (Vol 161 pp 12–22) PMLB
1167	Retrieved from https://proceedings.mlr.press/v161/rvder21a.html
1169	Saifuddin M Abramoff B Z Davidson E A Dietze M C & Finzi A C
1160	(2021) Identifying Data Needed to Beduce Parameter Uncertainty in a
1170	Coupled Microbial Soil C and N Decomposition Model Journal of Geo-
1170	<i>physical Research: Biogeosciences</i> 126(12) Retrieved from https://
1171	agunubs onlinelibrary wiley com/doi/abs/10 1029/2021 IG006593 doi:
1172	10 1029/2021 IC006593
1173	Solimons T. Kingma D. & Welling M. (2015) Markov Chain Monte Carle and
1174	Variational Informace: Bridging the Cap In F. Bach & D. Bloi (Edg.) Pro
1175	ceedings of the 22nd International Conference on Machine Learning (Vol. 37)
1170	Lille France: PMLR Retrieved from https://proceedings.mlr.press/w37/
1170	calimand 5 html
1178	Salimans T & Knowles D A (2013) Fixed Form Variational Postarior Approx
1179	implies through Stochastic Linear Bogrossion $Bauesian Analysis S(A)$ Bo
1180	triough from https://doi.org/10.1211/13_BA858.doi: 10.1211/13.BA858
1181	Solution I Wight T V & Fornachael C (2016) Drobabilistic programming in
1182	Duthon using DuMC2 Boon I Computer Science & Detriousd from https://
1183	doi org/10 7717/poori-ca 55 doi: 10 7717/poori as 55
1184	Säuldrä S (2012) Paulogian Filtening and Smoothing Combridge, Combridge
1185	University Process – Potnieved from https://www.combridge.com/hocks/
1186	bayogion-filtoring-ond-gmoothing(C272EP21CED0A100E9476C1P22721A67
1187	d_{0} : 10 1017/CBO07811303/4203
1188	Säulta S & Solin A (2010) Applied Stachastic Differential Freeditions Combridge
1189	University Pross. doi: 10.1017/0781108186725
1190	Suiono D. Vie H. W. Allicon S. & Suddowth F. D. (2022). Variational Information
1191	for Soil Diogeochemical Models In <i>ICMI</i> 2000 and AI for Science Workshop
1192	Batriavad from https://apaproviau.nat/farum2id=2 HruarDdAll
1193	Sulman D. N. Maana I. A. Abramaff D. Avanill C. Kindin S. Caanaian K.
1194	Suman, D. N., Moore, J. A., Abramon, R., Averni, C., Kivini, S., Georgiou, K.,
1195	2016). Multiple models and experiments underscore
1196	Targe uncertainty in son carbon dynamics. $Diogeochemistry, 141(2)$. doi: 10.1007/ $_{a}$ 10522.018.0500 $_{a}$
1197	10.1007/S10000-0000-2
1198	Sulman, B. N., Phillips, R. P., Olsni, A. C., Snevilakova, E., & Pacala, S. W. (2014).
1199	Microbe-driven turnover onsets inneral-mediated storage of son carbon under algorithm QQ . Nature Climate Change $1/(12)$ dai: 10.1028/nalimate2426
1200	elevated OO_2 . Nature Cumate Change, 4 (12). doi: 10.1056/fichinate2450
1201	Summers, C., & Dinneen, M. J. (2020). Four Things Everyone Should Know to Im-
1202	totiona Dataianad from https://energy.org/formational Conjetence on Learning Represen-
1203	Tedd Prown K E Dandowon I T Hapling E Arong V Haiima T Level
1204	C Allicon S D (2014) Changes in seil argenie carbon store
1205	$0., \ldots$ Allson, 5. D. (2014). Unalges in soli organic carbon storage predicted
1206	by Earth system models during the 21st century. Biogeosciences, $11(8)$. doi: 10.5104/bg.11.2341.2014
1207	Todd Drown K E O Dondonson I T Deet W M Hefferer E M T
1208	C Schuur E A C & Allicon S D (2012) Courses of vertication in seil series
1209	U., Schuur, E. A. G., & Allison, S. D. (2015). Uauses of variation in soil car-
1210	bon simulations from Owner 5 Earth system models and comparison with obser- vations. Biogeogeoionage $10(3)$. Betrieved from https://he_comparison.com/
1211	values. Diogeosciences, 10(3). Retrieved from fittps://bg.coperfitcus.org/

1212	articles/10/1717/2013/ doi: $10.5194/bg-10-1717-2013$
1213	Tzen, B., & Raginsky, M. (2019). Neural Stochastic Differential Equations: Deep La-
1214	tent Gaussian Models in the Diffusion Limit. arXiv. Retrieved from https://
1215	arxiv.org/abs/1905.09883 doi: 10.48550/ARXIV.1905.09883
1216	Vehtari, A., Gelman, A., & Gabry, J. (2017). Practical Bayesian model evalua-
1217	tion using leave-one-out cross-validation and WAIC. Statistics and Computing,
1218	27(5). doi: 10.1007/s11222-016-9696-4
1219	Wang, S., Luo, Y., & Niu, S. (2022). Reparameterization Required After Model
1220	Structure Changes From Carbon Only to Carbon-Nitrogen Coupling. Jour-
1221	nal of Advances in Modeling Earth Systems, 14(4). Retrieved from https://
1222	agupubs.onlinelibrary.wiley.com/doi/abs/10.1029/2021MS002798 doi:
1223	10.1029/2021MS002798
1224	Whitaker, G. A. (2016). Bayesian inference for stochastic differential mixed-effects
1225	models (Unpublished doctoral dissertation). Newcastle University.
1226	Whitaker, G. A., Golightly, A., Boys, R. J., & Sherlock, C. (2017). Bayesian
1227	Inference for Diffusion-Driven Mixed-Effects Models. Bayesian Analy-
1228	sis, 12(2). Retrieved from https://doi.org/10.1214/16-BA1009 doi:
1229	10.1214/16-BA1009
1230	Wieder, W. R., Boehnert, J., & Bonan, G. B. (2014). Evaluating soil biogeochem-
1231	istry parameterizations in Earth system models with observations. Global Bio-
1232	geochemical Cycles, 28(3). Retrieved from https://agupubs.onlinelibrary
1233	.wiley.com/doi/abs/10.1002/2013GB004665 doi: 10.1002/2013GB004665
1234	Wieder, W. R., Grandy, A. S., Kallenbach, C. M., Taylor, P. G., & Bonan, G. B.
1235	(2015). Representing life in the Earth system with soil microbial functional
1236	traits in the MIMICS model. Geoscientific Model Development, 8(6). doi:
1237	10.5194/gmd-8-1789-2015
1238	Wiqvist, S., Golightly, A., McLean, A. T., & Picchini, U. (2021). Efficient inference
1239	for stochastic differential equation mixed-effects models using correlated par-
1240	ticle pseudo-marginal algorithms. Computational Statistics & Data Analysis,
1241	157. Retrieved from https://www.sciencedirect.com/science/article/
1242	pii/S0167947320302425 doi: 10.1016/j.csda.2020.107151
1243	Wood, T. E., González, G., Silver, W. L., Reed, S. C., & Cavaleri, M. A. (2019). On
1244	the shoulders of giants: Continuing the legacy of large-scale ecosystem manipu-
1245	lation experiments in Puerto Rico. Forests, $10(3)$. doi: $10.3390/f10030210$
1246	Xie, H. W., Romero-Olivares, A. L., Guindani, M., & Allison, S. D. (2020). A
1247	Bayesian approach to evaluation of soil biogeochemical models. <i>Biogeosciences</i> ,
1248	17(15). Retrieved from https://bg.copernicus.org/articles/17/4043/
1249	2020/ doi: 10.5194/bg-17-4043-2020
1250	Yao, Y., Vehtari, A., Simpson, D., & Gelman, A. (2018). Yes, but Did It Work?:
1251	Evaluating Variational Inference. In J. Dy & A. Krause (Eds.), Proceedings of
1252	the 35th International Conference on Machine Learning (Vol. 80). PMLR. Re-
1253	trieved from https://proceedings.mlr.press/v80/yao18a.html
1254	You, K., Long, M., Wang, J., & Jordan, M. I. (2019). How Does Learning Rate De-
1255	cay Help Modern Neural Networks? arXiv. Retrieved from https://arxiv
1256	.org/abs/1908.01878 doi: 10.48550/ARXIV.1908.01878
1257	Zhu, Q., Bi, W., Liu, X., Ma, X., Li, X., & Wu, D. (2020). A Batch Normalized In-
1258	ference Network Keeps the KL Vanishing Away. In Proceedings of the 58th An-
1258 1259	ference Network Keeps the KL Vanishing Away. In Proceedings of the 58th An- nual Meeting of the Association for Computational Linguistics. Online: Asso-
1258 1259 1260	ference Network Keeps the KL Vanishing Away. In <i>Proceedings of the 58th Annual Meeting of the Association for Computational Linguistics</i> . Online: Association for Computational Linguistics. Retrieved from https://aclanthology

Supporting Information for "A framework for variational inference and data assimilation of soil biogeochemical models using state space approximations and normalizing flows"

H. W. Xie^{1,†}, D. Sujono^{2,†}, T. Ryder³, E. Sudderth², S. D. Allison^{1,4}

¹Center for Complex Biological Systems, University of California, Irvine, Irvine, CA, United States of America

²Department of Computer Science, University of California, Irvine, Irvine, CA, United States of America

³School of Mathematics Statistics and Physics, Newcastle University, Newcastle, United Kingdom

⁴Department of Ecology and Evolutionary Biology, University of California, Irvine, Irvine, CA, United States of America

 $^\dagger\mathrm{Authors}$ contributed equally to this work.

Contents of this file

- 1. Figures S1 to S6
- 2. Table S1

Introduction

This document contains figures supporting the validity and functionality of our neural moving average flow VI framework. Figure S1 illustrates the benefit of initiating VI with an ELBO training warmup phase at low learning rates. Figure S2 demonstrates with an example $-\mathcal{L}$ trajectory from an SCON-C approximation inference that our VI algorithm is able to stably converge in ELBO. Figures S3 and S4 indicate that the neural moving

average flow VI approach remains viable for inference on approximated SCON-SS, and by extension, state space models that are linear in drift but non-linear in diffusion. Figure S5 depicts the importance of including CO_2 information in the data y for subtantial

improvement of posterior identifiability and certainty. Figure S6 contrasts the effects of lengthening experiment time span T versus thickening observations in y to better inform and identify posteriors. Finally, Table S1 details the hyperparameters corresponding to our informed and independent univariate logit-normal priors.



 \mathcal{I}_{-}

-1400

-1600

60000

70000

80000

:

Figure S1. Comparison of $-\mathcal{L}$ trajectories from the latter halves of T = 5000 hour SCON-C flow trainings without (blue) and with (orange) warmup indicates that warmup helps stabilize training and speed up convergence. The trajectory corresponding to warmup displays much less prominent instability spiking and has flattened more quickly in contrast to the that of the no warmup counterpart.

90000

Iteration

100000

110000

120000



:

Figure S2. The stabilizing of the $-\mathcal{L}$ trajectory between -1550 and -1600 in the latter half of T = 5000 SCON-C flow VI training indicates convergence to an approximate local minimum $-\mathcal{L}$ and thereby proper algorithm function of the $q(\theta, x; \phi_{\theta,x})$ joint optimization.

Similarly stabilizing $-\mathcal{L}$ trajectories were observed for inferences on SCON-SS state space model approximations.



:

Figure S3. Flow-approximated SCON-SS $q(x|\theta; \phi_x)$ latent state and observed CO₂ means conditioned on T = 5000 SCON-SS data-generating process y estimated from 250 x paths sampled from the optimized joint variational $q(\theta, x; \phi_{(\theta, x)})$ density.



Figure S4. Full SCON-SS state space model marginal $q(\theta; \phi_{\theta})$ posterior densities (orange) conditioned on T = 5000 SCON-SS data-generating process y compared to the prior densities $p(\theta)$ (blue). The true θ values sampled during data generation are marked by vertical dashed gray lines.



Figure S5. Approximate SCON-C state space model marginal $q(\theta; \phi_{\theta})$ posterior densities conditioned with (orange) and without (green) CO_2 information in y produced by the same SCON-C data-generating process compared to mean-field prior densities $p(\theta)$ (blue). The true θ values sampled during data generation are marked by vertical dashed gray lines.



Approximate SCON-C state space model marginal $q(\theta; \phi_{\theta})$ posterior densities Figure S6. conditioned with T = 1000 data observed every 5 hours (blue), T = 5000 data observed every 5 hours (orange), and T = 1000 data observed every hour (green). All three y share the same SCON-C data-generating process and include CO_2 information. The true θ values sampled during data generation are marked by vertical dashed gray lines.

θ	Biogeochemical interpretation	Target hyperparameters	Units
u_M	MBC uptake rate	$\mathscr{LN}(0.0016, 0.0004, 0, 1)$	${ m mgCg^{-1}Ch^{-1}}$
a_{DS}	DOC to SOC transfer fraction	$\mathscr{LN}(0.5, 0.125, 0, 1)$	NA
a_{SD}	SOC to DOC transfer fraction	$\mathscr{LN}(0.5, 0.125, 0, 1)$	NA
a_M	MBC to organic C transfer fraction	$\mathscr{LN}(0.5, 0.125, 0, 1)$	NA
a_{MSC}	MBC to SOC transfer fraction	$\mathscr{LN}(0.5, 0.125, 0, 1)$	NA
$k_{S,\mathrm{ref}}$	SOC decomposition rate	$\mathscr{LN}(0.0005, 0.000125, 0, 0.1)$	${ m mg}{ m C}{ m mg}^{-1}{ m C}{ m h}^{-1}$
$k_{D,\mathrm{ref}}$	DOC decomposition rate	$\mathscr{LN}(0.0008, 0.0002, 0, 0.1)$	$mg C mg^{-1} C h^{-1}$
$k_{M,\mathrm{ref}}$	MBC decomposition rate	$\mathscr{LN}(0.0007, 0.000175, 0, 0.1)$	$mg C mg^{-1} C h^{-1}$
Ea_S	SOC decomposition activation energy	$\mathscr{LN}(20,5,5,80)$	${ m kJmol^{-1}}$
Ea_D	DOC decomposition activation energy	$\mathscr{LN}(20,5,5,80)$	${ m kJmol^{-1}}$
Ea_M	MBC decomposition activation energy	$\mathscr{LN}(20,5,5,80)$	${ m kJmol}^{-1}$
c_S	SCON-C SOC β constant	$\mathscr{LN}(0.1, 0.025, 0, 0.1)$	$\mathrm{mg}\mathrm{C}\mathrm{g}^{-1}\mathrm{soil}$
c_D	SCON-C DOC β constant	$\mathscr{LN}(0.002, 0.0005, 0, 0.1)$	$\mathrm{mg}\mathrm{C}\mathrm{g}^{-1}\mathrm{soil}$
c_M	SCON-C MBC β constant	$\mathscr{LN}(0.002, 0.0005, 0, 0.1)$	$\mathrm{mg}\mathrm{C}\mathrm{g}^{-1}\mathrm{soil}$
s_S	SCON-SS SOC β factor	$\mathscr{LN}(0.0005, 0.000125, 0, 0.1)$	NA
s_D	SCON-SS DOC β factor	$\mathscr{LN}(0.0005, 0.000125, 0, 0.1)$	NA
s_M	SCON-SS MBC β factor	$\mathscr{LN}(0.0005, 0.000125, 0, 0.1)$	NA

:

Table S1. List of SCON-C and SCON-SS θ and their corresponding marginal data-generating and informed prior hyperparameters. The marginal densities are formatted as $\mathscr{LN}(\mu, \sigma, a, b)$, where μ and σ are the desired target density mean and standard deviation and a and b are the truncated distribution support lower and upper bounds.