Training physics-based machine-learning parameterizations with gradient-free ensemble Kalman methods

Ignacio Lopez-Gomez^{1,1}, Costa D Christopoulos^{1,1}, Haakon Ludvig Langeland Ervik^{1,1}, Oliver Dunbar^{1,1}, Yair Cohen^{1,1}, and Tapio Schneider^{1,1}

¹California Institute of Technology

November 30, 2022

Abstract

Most machine learning applications in Earth system modeling currently rely on gradient-based supervised learning. This imposes stringent constraints on the nature of the data used for training (typically, residual time tendencies are needed), and it complicates learning about the interactions between machine-learned parameterizations and other components of an Earth system model. Approaching learning about process-based parameterizations as an inverse problem resolves many of these issues, since it allows parameterizations to be trained with partial observations or statistics that directly relate to quantities of interest in long-term climate projections. Here we demonstrate the effectiveness of Kalman inversion methods in treating learning about parameterizations as an inverse problem. We consider two different algorithms: unscented and ensemble Kalman inversion. Both methods involve highly parallelizable forward model evaluations, converge exponentially fast, and do not require gradient computations. In addition, unscented Kalman inversion provides a measure of parameter uncertainty. We illustrate how training parameterizations can be posed as a regularized inverse problem and solved by ensemble Kalman methods through the calibration of an eddy-diffusivity mass-flux scheme for subgrid-scale turbulence and convection, using data generated by large-eddy simulations. We find the algorithms amenable to batching strategies, robust to noise and model failures, and efficient in the calibration of hybrid parameterizations that can include empirical closures and neural networks.

Training physics-based machine-learning parameterizations with gradient-free ensemble Kalman methods

Ignacio Lopez-Gomez¹, Costa Christopoulos¹, Haakon Ludvig Langeland Ervik¹, Oliver R. A. Dunbar¹, Yair Cohen¹, Tapio Schneider^{1,2}

¹Department of Environmental Science and Engineering, California Institute of Technology, Pasadena, CA, USA. ²Jet Propulsion Laboratory, California Institute of Technology, Pasadena, CA, USA.

Key Points:

1

2

3

4

5

6

7

8

9

14

10	•	Ensemble Kalman methods can be used to train parameterizations regardless of
11		their architecture.
12	•	They enable learning from partial observations or statistics in the presence of noise
13	•	Their effectiveness is demonstrated by calibrating an atmospheric turbulence and

convection scheme.

 $Corresponding \ author: \ Ignacio \ Lopez-Gomez, \ \texttt{ilopezgo@caltech.edu}$

15 Abstract

Most machine learning applications in Earth system modeling currently rely on gradient-16 based supervised learning. This imposes stringent constraints on the nature of the data 17 used for training (typically, residual time tendencies are needed), and it complicates learn-18 ing about the interactions between machine-learned parameterizations and other com-19 ponents of an Earth system model. Approaching learning about process-based param-20 eterizations as an inverse problem resolves many of these issues, since it allows param-21 eterizations to be trained with partial observations or statistics that directly relate to 22 quantities of interest in long-term climate projections. Here we demonstrate the effec-23 tiveness of Kalman inversion methods in treating learning about parameterizations as 24 an inverse problem. We consider two different algorithms: unscented and ensemble Kalman 25 inversion. Both methods involve highly parallelizable forward model evaluations, con-26 verge exponentially fast, and do not require gradient computations. In addition, unscented 27 Kalman inversion provides a measure of parameter uncertainty. We illustrate how train-28 ing parameterizations can be posed as a regularized inverse problem and solved by en-29 semble Kalman methods through the calibration of an eddy-diffusivity mass-flux scheme 30 for subgrid-scale turbulence and convection, using data generated by large-eddy simu-31 lations. We find the algorithms amenable to batching strategies, robust to noise and model 32 failures, and efficient in the calibration of hybrid parameterizations that can include em-33 pirical closures and neural networks. 34

35 Plain Language Summary

Artificial intelligence represents an exciting opportunity in Earth system model-36 ing, but its application brings its own set of challenges. One of these challenges is to train 37 machine learning systems within Earth system models from partial or indirect data. Here 38 we present algorithms, known as ensemble Kalman methods, that can be used to train 39 such systems. We demonstrate their use in situations where the data used for training 40 are noisy, only indirectly informative about the model to be trained, and may only be-41 come available sequentially. As an example, we present training results for a state-of-42 the-art model for turbulence, convection, and clouds for use within Earth system mod-43 els. This model is shown to learn efficiently from data in a variety of configurations, in-44 cluding situations where the model contains neural networks. 45

46 1 Introduction

The remarkable achievements of machine learning over the past decade have led to renewed interest in informing Earth system models with data (Schneider et al., 2017; Reichstein et al., 2019). The spotlight is often on creating or improving models of processes that are deemed important for the correct representation of the Earth system as a whole. Examples of these processes include moist convection (Brenowitz et al., 2020), cloud microphysical and radiative effects (Seifert & Rasp, 2020; Villefranque et al., 2021; Meyer et al., 2022), and evapotranspiration (Zhao et al., 2019), among others.

Processes governed by poorly understood dynamics, such as biological processes, are obvious candidates for representation by purely data-driven models. On the other end of the spectrum are fluid transport processes, which are governed by the Navier-Stokes equations. Uncertain representation of these processes comes from a lack of resolution, not lack of knowledge about the underlying dynamics. Hybrid modeling approaches that incorporate domain knowledge and augment it by learning from data are attractive for such processes, because they reduce what needs to be learned from data.

For processes with known dynamics, data-informed models fall into three broad categories according to their leverage of domain knowledge. In the first category are models that try to learn the entire dynamics using a sufficiently expressive hypothesis set,

such as deep neural networks. This approach has proved successful for predicting pre-64 cipitation over short time horizons (Ravuri et al., 2021), and it has been explored for medium-65 range weather forecasting (Rasp & Thuerey, 2021; Pathak et al., 2022; Lopez-Gomez et 66 al., 2022). An advantage of these models is that they are typically easy to implement 67 and cheap to evaluate. They can afford very large time steps (Weyn et al., 2021), or they 68 may learn directly mappings from the initial state to a probability distribution of final 69 states with no need of time marching or ensemble forecasting (Sønderby et al., 2020). 70 A deficiency of these models is that they often require an extreme amount of data to con-71 strain the many (often $> 10^6$) parameters in them and to achieve acceptable performance. 72

Methods in the second and third categories employ models of subgrid processes to 73 solve the closure problem that arises when coarse-graining the known dynamics, which 74 are retained. Retaining the coarse-grained equations of motion ensures conservation of 75 mass, momentum, and energy, which is more difficult when using models in the first cat-76 egory (Beucler et al., 2021; Brenowitz et al., 2020). The second category encompasses 77 methods that try to learn the functional form of these closures avoiding the use of em-78 pirical laws. For example, Zanna and Bolton (2020) use relevance vector machines to prune 79 a library of functions, resulting in a closed form expression of mesoscale eddy fluxes in 80 ocean simulations; Ling et al. (2016) learn a neural network closure of the Reynolds stress 81 anisotropy tensor while explicitly encoding rotational invariance in the context of k-82 ϵ models of turbulence. 83

Finally, the third category refers to methods that seek to learn the parameters that 84 arise in empirical closures of subgrid processes. In general, models in the third category 85 are more restrictive, and they may be expected to underperform with respect to those 86 in the second category given sufficient data on the target distributions. However, the lim-87 ited parametric complexity of these closures makes them amenable to physical interpre-88 tation, robust to overfitting, and better suited for learning in the low-data regime. This 89 may be attractive for Earth system models, for which online learning from limited high-90 resolution data may be a useful strategy to assimilate computationally generated data 91 of the changing climate (Schneider et al., 2017). 92

A barrier delimiting data-driven and empirical subgrid-scale closures is the access 93 to practical calibration tools. Neural network parameterizations are easily calibrated us-94 ing stochastic gradient descent through backpropagation, which limits datasets to those 95 including output labels, and models to those that afford automatic differentiation with respect to their parameters. Empirical closures, which may depend on time-evolving terms 97 with memory (e.g., Lopez-Gomez et al., 2020) or yield unobservable outputs (e.g., tur-98 bulent versus dynamical entrainment in Cohen et al., 2020) cannot be trained using this 99 approach. Traditional Bayesian inference techniques, like random walk Metropolis (Metropolis 100 et al., 1953) or sequential Monte Carlo (Moral et al., 2006), can be used in this context 101 if the number of parameters is small and the model to be trained is cheap to evaluate. 102 Such methods additionally provide uncertainty quantification, but they become intractable 103 for expensive models with many parameters (e.g., Cotter et al., 2013; Souza et al., 2020). 104 Model-agnostic tools that enable fast calibration of subgrid-scale closures from diverse 105 data are a necessary step toward the development of hybrid closures that leverage the 106 strengths of all modeling approaches. 107

With this goal in mind, we present calibration strategies for models of subgrid pro-108 cesses, formulating the learning task as an inverse problem (Kovachki & Stuart, 2019). 109 Solutions to the inverse problem are sought using the ensemble and unscented Kalman 110 inversion algorithms (Iglesias et al., 2013; Huang, Schneider, & Stuart, 2022). Empha-111 112 sis is given to practical aspects of this specific inverse problem, which have not previously been explored in the literature. These include the construction of a domain-agnostic 113 loss function from high-dimensional observations, a heuristic a priori estimate of model 114 error, systematic handling of model failures during the training process, and the use of 115

the Kalman inversion algorithms when only noisy evaluations of the loss function are avail-able.

The strategies presented here are designed to have several attractive properties com-118 pared to other learning algorithms. First, framing learning as an inverse problem enables 119 the use of partial observations or statistically aggregated data. Second, calibration is per-120 formed using gradient-free methods, which are well suited for stochastic models and/or 121 models whose derivatives do not exist or are difficult to obtain. Finally, the strategies 122 presented are amenable to parallelization and the use of high-dimensional correlated ob-123 servations. The last two properties draw heavily on the use of Kalman inversion algo-124 rithms to tackle the inverse problem, which themselves build on the success of the en-125 semble Kalman filter (EnKF) for data assimilation (Evensen, 1994; Houtekamer & Mitchell, 126 1998; Burgers et al., 1998) and are closely related to iterative EnKF (Chen & Oliver, 2012; 127 Emerick & Reynolds, 2013; Bocquet & Sakov, 2013). The methods presented here are 128 applicable to models of subgrid-scale processes, within the second and third categories 129 described above. They provide an alternative to learning algorithms that impose strin-130 gent requirements on either the model architecture, its computational cost, or the na-131 ture of the training data. 132

The article is organized as follows. Section 2 casts learning about parameteriza-133 tions as an inverse problem, which can be solved through the minimization of a regu-134 larized low-dimensional encoding of the data-model mismatch. Section 3 reviews the ap-135 plication of the ensemble and unscented Kalman inversion algorithms to inverse prob-136 lems and proposes modifications to their update equations that enable training models 137 that may experience failures. Section 4 then applies these ensemble Kalman algorithms 138 to the calibration of closures within an eddy-diffusivity mass-flux (EDMF) scheme of tur-139 bulence and convection, using data generated from large-eddy simulations (LES). The 140 robustness of these learning strategies is demonstrated by calibrating the EDMF scheme 141 using noisy loss evaluations and partial information, and their flexibility is emphasized 142 by learning the parameters in a hybrid model containing both empirical and neural net-143 work closures. Finally, Section 5 ends with a discussion of the findings and concluding 144 remarks. 145

¹⁴⁶ 2 Learning about parameterizations as an inverse problem

We consider the problem of learning the parameters ϕ of a dynamical model $\Psi(\phi)$, 147 using noisy observations y of the true dynamical system ζ that $\Psi(\phi)$ seeks to represent. 148 In the context of subgrid parameterizations, $\Psi(\phi)$ represents a closed version of the coarse-149 grained dynamical system (e.g., the filtered Navier-Stokes equations), where closures are 150 parameterized by ϕ . The model $\Psi(\phi)$ maps a user-defined initial state φ_0 and a forcing 151 $F_{\varphi}(t)$ to a state trajectory $\hat{\varphi}(t)$. Thus, our definition of $\Psi(\phi)$ can be interpreted as the 152 iterative application of the resolvent operator on the initial field φ_0 (Brajard et al., 2021). 153 In the following, we denote any set of initial and forcing conditions collectively as the 154 configuration $x_c = \{\varphi_0, F_{\varphi}\}_c$; the definition of all symbols is summarized in the appendix. 155

For each configuration x_c , the dynamical model can be related to the observations y_c by the observational map \mathcal{H}_c , which encapsulates all averaging and post-processing operations necessary to yield the model predictions associated with the observations. More precisely, the relationship between the observations y_c , the true dynamics ζ , and the dynamical model $\Psi(\phi)$ for a given configuration may be expressed as

$$y_c = \mathcal{H}_c \circ \zeta(x_c) + \eta_c = \mathcal{H}_c \circ \Psi(\phi; x_c) + \delta(x_c) + \eta_c, \tag{1}$$

where $\phi \in \mathbb{R}^p$ is the vector of learnable parameters, η_c is the observational noise asso-

ciated with y_c , and $\delta(x_c)$ is the model or representation error, which we define as the mis-

match between the denoised observations $\mathcal{H}_c \circ \zeta(x_c)$ and the output of a best-fitting model

 $\mathcal{H}_c \circ \Psi(\phi^*; x_c)$, following Kennedy and O'Hagan (2001). Thus, the model error is ap-

proximated as additive (Cohn, 1997; van Leeuwen, 2015) and defined with respect to the observational map \mathcal{H}_c and the optimal parameters ϕ^* that minimize its contribution to the data-model relation (1).

Observations are taken to come from finite spatial and temporal averages of fields 163 such as temperature. Learning from averages can help prevent overfitting to trajecto-164 ries in chaotic systems by focusing on the statistics of the dynamics (Morzfeld et al., 2018). 165 It also improves numerical stability when coupling to a parent model (Brenowitz & Brether-166 ton, 2018). Under this definition of observations, it is reasonable to assume the noise η_c 167 to be additive and Gaussian. In the following, we will further consider $\delta(\cdot)$ to be a cen-168 tered Gaussian, although this constitutes a significantly stronger assumption (e.g., that 169 the model is unbiased) and may not be appropriate for a detailed characterization of pos-170 terior uncertainty (van Leeuwen, 2015; Brynjarsdóttir & O'Hagan, 2014). The construc-171 tion of more precise error models remains a challenge beyond the scope of this work. These 172 assumptions enable us to write $\delta(x_c) + \eta_c \sim \mathcal{N}(0, \Gamma_c)$. 173

In general, we are interested in minimizing the mismatch between y_c and the model output for a wide range of configurations $C = \{x_c, c = 1, ..., |C|\}$ that are representative of the conditions in which the model will operate. This defines the global datamodel relation

$$y = \mathcal{H} \circ \Psi(\phi) + \delta + \eta, \tag{2}$$

where $y = [y_1, \ldots, y_{|C|}]^T \in \mathbb{R}^d$, $\delta = [\delta(x_1), \ldots, \delta(x_{|C|})]^T$, $\eta = [\eta_1, \ldots, \eta_{|C|}]^T$, $\mathcal{H} \circ \Psi(\phi) = [\mathcal{H}_1 \circ \Psi(\phi; x_1), \ldots, \mathcal{H}_{|C|} \circ \Psi(\phi; x_{|C|})]^T$ and $\delta + \eta \sim \mathcal{N}(0, \Gamma)$. In addition, implicit in the definition of the dynamical model $\Psi(\phi)$ is a discrete resolution Δ . This dependence may be lifted if the closures are designed to be scale-aware or scale-independent, in which case the relation (2) should be augmented by stacking copies of y and evaluating $\mathcal{H} \circ \Psi(\phi, \Delta_i)$ for different discretizations Δ_i .

In practice, the parameters ϕ are often defined over some subspace $U \subset \mathbb{R}^p$ outside of which the model trajectories are unphysical or numerically unstable. Examples of these are parameters controlling the diffusion or turbulent dissipation of a scalar field, for which negative values are not physically valid. On the other hand, many algorithms designed to solve inverse problems assume $\phi \in \mathbb{R}^p$. This obstacle may be circumvented by defining a transformation $\mathcal{T} : U \to \mathbb{R}^p$, such that the global data-model relation (2) can be defined in an unconstrained parameter space,

$$y = \mathcal{G}(\theta) + \delta + \eta, \tag{3}$$

where

$$\mathcal{G} \coloneqq \mathcal{H} \circ \Psi \circ \mathcal{T}^{-1}, \qquad \phi = \mathcal{T}^{-1}(\theta). \tag{4}$$

In expressions (3) and (4), $\theta \in \mathbb{R}^p$ is the parameter vector in unconstrained space and $\mathcal{G} : \mathbb{R}^p \to \mathbb{R}^d$ is the map from transformed parameters to model predictions, which represents the forward model. The task of learning a set of model parameters θ under relation (3) can be cast as the Bayesian inverse problem of finding the posterior (Kaipio & Somersalo, 2006; Tarantola, 2005; Huang, Huang, et al., 2022)

$$\rho(\theta|y,\Gamma) = \frac{e^{-\mathcal{L}(\theta;y)}}{Z(y|\Gamma)}\rho_{\text{prior}}(\theta), \qquad \mathcal{L}(\theta;y) = \frac{1}{2}||y - \mathcal{G}(\theta)||_{\Gamma}^{2}, \tag{5}$$

where $Z(y|\Gamma)$ is a normalizing constant, $\|\cdot\|_{\Gamma}^2$ denotes the Mahalanobis norm $\langle\cdot, \Gamma^{-1}\cdot\rangle$, 180 \mathcal{L} is the loss or negative log-likelihood, and $\rho_{\text{prior}}(\theta)$ is the prior density. We stress that 181 the posterior $\rho(\theta|y,\Gamma)$ is conditioned on our approximation of the noise $\delta + \eta$; see Kennedy 182 and O'Hagan (2001) for a discussion on the usefulness and caveats of such an approach. 183 Given the inverse problem (3)-(5), we may be interested in finding the maximum a pos-184 teriori (MAP), approximations of the density $\rho(\theta|y,\Gamma)$ around the MAP for uncertainty 185 quantification, or simply the maximum likelihood estimator (MLE) if we have no prior 186 information about θ . Algorithms to perform these tasks are described in Section 3. 187

The error covariance Γ_c appearing in each model-data relation (1), and ultimately defining the inverse problem (3)–(5), is yet to be defined. In Section 2.1, we suggest an estimate of Γ_c relevant to the calibration of models with an unknown error structure $\delta(\cdot)$. In addition, the choice of observational map \mathcal{H}_c may not be evident when training dynamical models that aim to represent complex dynamical systems ζ with many observable fields. Section 2.2 suggests a model-agnostic definition of \mathcal{H}_c that can be used to construct a regularized inverse problem.

2.1 Estimate of noise covariances

Since the structure of the representation or model error δ is unknown a priori, we must either parameterize it and calibrate it as well (Brynjarsdóttir & O'Hagan, 2014), or use a heuristic to capture its magnitude. Here, we follow the second route and offer a heuristic that has worked well for us in practice. If we take $y_c = y_c(t)$ to be an observation of the true system in configuration x_c aggregated over a time interval $[t, t + \tau]$, we can write equation (1) as

$$y_c(t) - y_c(0) = \mathcal{H}_c \circ \Psi(\phi; x_c, t) - y_c(0) + \delta(x_c; t) + \eta_c(t).$$
(6)

If we further consider a model with no predictive power of the first kind (Lorenz, 1975; Schneider & Griffies, 1999), such that $\mathcal{H}_c \circ \Psi(\phi; x_c, t) \approx y_c(0)$ for all times t, the covariance of (6) from t = 0 to $t = t_c \gg \tau$ reads

$$\Gamma_c = \operatorname{Cov}(y_c) \approx \operatorname{Cov}(\delta(x_c)) + \operatorname{Cov}(\eta_c), \tag{7}$$

which yields an estimate of the aggregate noise $\eta_c + \delta(x_c) \sim \mathcal{N}(0, \Gamma_c)$ from the vari-196 ability of the observation y_c over a time interval $[0, t_c]$. For non-stationary conditions or 197 finite-time averages, Γ_c depends on t_c . Estimating the magnitude of the aggregate noise 198 from the internal variability of the true dynamics ensures that the loss or negative log-199 likelihood $\mathcal{L}(\theta; y)$ penalizes models $\Psi(\phi)$ that produce unrealistic outputs, and it rep-200 resents a form of error inflation if the best-fitting model is expected to outperform the 201 aforementioned unskillful model. The heuristic (7) is most appropriate when the dynam-202 ical model $\Psi(\phi)$ is expressive enough to closely replicate the initial observations $y_c(0)$, 203 such that any mismatch in the initial condition can be lumped together with the observation error. 205

206 207

195

2.2 Design of the observational map

2.2.1 Application to problems with high-resolution data

High-resolution data are becoming increasingly common, from reanalysis products (Muñoz-Sabater et al., 2021), satellite imagery (Schmit et al., 2017), and partial differential equation (PDE) solvers such as LES (Shen et al., 2022). Although computationally generated and thus suffering from their own limitations (e.g., microphysical processes still need to be parameterized in LES), data from PDE solvers have some particularly desirable properties for the calibration of dynamical models:

214 215 216

• All variables appearing in the coarse-grained equations of motion are observable. As a consequence, the nature of the observational map \mathcal{H} used to constrain the model is largely a design choice.

• Data can be obtained systematically for all configurations x_c of interest, which may be chosen to minimize parameter uncertainty through active learning (Dunbar et al., 2022). In contrast, data drawn from physical measurements (e.g., field observations) are often sparse in the space of forcing and boundary conditions.

High-resolution data are often high-dimensional, which poses particular difficulties regarding the conditioning and tractability of linear systems of equations when solving inverse problems. The guidelines for the construction of the observational map \mathcal{H} presented here are tailored to solve these issues, with a focus on data from high-fidelity solvers.

2.2.2 Model calibration

225

238

We define *model calibration* as the minimization of the mismatch between the observed dynamics and the dynamics induced by the model. We will use this definition to construct a domain-agnostic map \mathcal{H} . As an example, consider a system ζ with coarsegrained dynamics

$$\frac{\partial \bar{\varphi}}{\partial t} + \bar{\boldsymbol{v}} \cdot \nabla \bar{\varphi} + \nabla \cdot (\overline{\boldsymbol{v}' \varphi'}) = F_{\varphi}, \qquad (8)$$

where $\overline{(\cdot)}$ denotes spatial filtering, $(\cdot)'$ subfilter-scale fluctuations, and F_{φ} is the forcing. The field $\bar{\boldsymbol{v}}$ is prescribed and $\overline{\boldsymbol{v}'\varphi'}$ is the term parameterized in $\Psi(\phi)$. Let $S(t) = [\bar{\varphi}(t), \ \overline{\boldsymbol{v}'\varphi'}(t)]^T$ be the true state augmented with subgrid-scale fluxes, and $\hat{S}(t)$ the augmented state predicted by the model. For an incompressible fluid model, S(t) would contain the fluid momentum, energy, and the subgrid advective fluxes of these fields.

Model calibration then entails finding the minimizer of the expected state mismatch $\mathbb{E}[\|\hat{S}-S\|]$ with respect to some norm and time interval, where the expectation is taken to allow for the calibration of stochastic models. Observations of the augmented state S(t), which includes subgrid-scale fluxes, are not always available. Therefore, this definition of model calibration is representative of the ideal learning scenario. In scenarios where the full state is not observable, we will consider S(t) to be an observed state formed by all relevant observable spatial fields.

2.2.3 Observations in physical space

Following our definition of model calibration, we preliminarily define the observations in the model-data relation (1) as finite-time averages of the normalized observed state s_c for a set of configurations C,

$$\tilde{y}_c = \frac{1}{T_c} \int_{t_c - T_c}^{t_c} s_c(\tau) d\tau, \qquad s_c = \begin{bmatrix} v_{c,1} \\ \cdots \\ v_{c,n_c} \end{bmatrix} = \begin{bmatrix} V_{c,1}/\sigma_{c,1} \\ \cdots \\ V_{c,n_c}/\sigma_{c,n_c} \end{bmatrix}, \qquad c = 1, \dots, |C|, \qquad (9)$$

where T_c is the averaging time, $v_{c,j} \in \mathbb{R}^{h_c}$ are the normalized spatial fields comprising s_c , $V_{c,j}$ are the components of the state S_c prior to normalization, n_c is the number of fields observed in configuration x_c , and h_c is the number of degrees of freedom of each field. As an example, the first configuration's observed state S_1 may include as fields atmospheric soundings of temperature and specific humidity $(n_1 = 2)$ measured at h_1 vertical locations above the surface, and the second configuration's state S_2 may include these fields as well as horizontal velocity profiles $(n_2 = 4)$, measured at h_2 different locations. Normalization of the observed state S_c is performed using the pooled time standard deviation $\sigma_{c,j}$ of each field $V_{c,j}$, with

$$\sigma_{c,j}^2 = h_c^{-1} \operatorname{tr} \Big[\operatorname{Cov}(V_{c,j}) \Big].$$
(10)

Covariances are computed over a time $t_c \ge T_c$ following the heuristic of Section 2.1 to capture the expected magnitude of the data mismatch,

$$\operatorname{Cov}(V_{c,j}) = \frac{1}{t_c} \int_0^{t_c} V_{c,j} V_{c,j}^T d\tau - \frac{1}{t_c^2} \Big(\int_0^{t_c} V_{c,j} d\tau \Big) \Big(\int_0^{t_c} V_{c,j} d\tau \Big)^T.$$
(11)

We resort to pooled normalization, instead of normalizing each of the dimensions of the observed state S_c by their standard deviation, because some of the dimensions of the spatial fields $V_{c,j}$ may not vary with a given forcing, resulting in zero-variance components. For example, in the atmospheric boundary layer, observations of liquid water specific hu-

₂₄₃ midity will always be zero below the lifting condensation level.

Stacking the observations from all configurations together, the full observation vector \tilde{y} is

$$\tilde{y} = \begin{bmatrix} \tilde{y}_1 \\ \dots \\ \tilde{y}_{|C|} \end{bmatrix} \in \mathbb{R}^{\tilde{d}}, \qquad \tilde{d} = \sum_{c=1}^{|C|} \tilde{d}_c = \sum_{c=1}^{|C|} n_c h_c.$$
(12)

Following again the heuristic in Section 2.1, the noise covariance associated with each observation vector $\tilde{y}_c \in \mathbb{R}^{\tilde{d}_c}$ is $\tilde{\Gamma}_c = \text{Cov}(s_c)$, computed as in equation (11). Given that the noise is estimated independently for each configuration, the full noise covariance is the block diagonal matrix

$$\tilde{\Gamma} = \begin{bmatrix} \tilde{\Gamma}_1 & 0 \\ & \ddots \\ 0 & & \tilde{\Gamma}_{|C|} \end{bmatrix} \in \mathbb{R}^{\tilde{d} \times \tilde{d}}, \qquad \tilde{\Gamma}_c = \operatorname{Cov}(s_c) \in \mathbb{R}^{\tilde{d}_c \times \tilde{d}_c}, \tag{13}$$

where $\tilde{\Gamma}_c$ is the noise covariance matrix of configuration c.

2.2.4 Observations in a reduced space

245

Each covariance matrix $\tilde{\Gamma}_c$, possibly associated with high-dimensional observations and a finite sampling interval, is likely to be rank-deficient and have a large condition number $\kappa = \mu_{c,1}/\mu_{c,r_c}$, where $\mu_{c,i}$ is the *i*-th largest eigenvalue of $\tilde{\Gamma}_c$ and r_c is the approximate rank of the matrix (Hansen, 1998). Numerically rank-deficient problems arise when \tilde{d}_c is greater than or equal to the number of samples used to construct $\tilde{\Gamma}_c$, or when there exist eigenvalues $\mu_{c,i}$ such that $\mu_{c,i}/\mu_{c,1} \leq \epsilon_m$, where ϵ_m is a measure of data or machine precision. An efficient regularization method for rank-deficient problems is to project the data from each configuration onto a lower-dimensional encoding, adding Tikhonov regularization to limit the condition number of the resulting global covariance matrix. If the lower-dimensional encoding is obtained through principal component analysis (PCA),

$$y_c = P_c^T \tilde{y}_c, \qquad \Gamma_c = d_c P_c^T \tilde{\Gamma}_c P_c + \kappa_*^{-1} \mu_1 I_{d_c}, \qquad (14)$$

where $y_c \in \mathbb{R}^{d_c}$, P_c is the projection matrix formed by the d_c leading eigenvectors of 246 Γ_c , I_{d_c} is the identity matrix, μ_1 is the leading eigenvalue of the unregularized global co-247 variance and κ_* is the limiting condition number of the global covariance, which should be chosen to be $\kappa_* < \epsilon_m^{-1/2}$. The encoding dimension d_c should be chosen such that $d_c \leq$ 248 249 $r_c \leq \tilde{d}_c$, where r_c is the approximate rank of $\tilde{\Gamma}_c$. The actual value of d_c may be cho-250 sen through the discrepancy principle, generalized cross validation, or based on the preser-251 vation of a given fraction of the total variance, among other criteria (Reichel & Rodriguez, 252 2013; Hansen, 1998). The Tikhonov inflation term regularizes problems where PCA is 253 performed between eigenvalues that are close in value, or where the range of configura-254 tion variances tr(Γ_c) is large (Hansen, 1990). In projection (14), since the number of re-255 tained principal components may differ among configurations for a given truncation cri-256 terion, each block covariance matrix is scaled by d_c . 257

Projection (14) enables the use of arbitrarily correlated observations by regularizing the linear system $\Gamma^{-1}(\mathcal{G}(\theta) - y)$ that appears in the gradient of the loss

$$\nabla \mathcal{L}(\theta; y) \propto (D\mathcal{G}(\theta))^T \Gamma^{-1}(\mathcal{G}(\theta) - y), \tag{15}$$

and lowering its computational cost. Here, $D\mathcal{G}(\theta) \in \mathbb{R}^{d \times p}$ is the Jacobian matrix of \mathcal{G} evaluated at θ . Although the ensemble Kalman algorithms presented in Section 3 do not compute the gradient (15) explicitly, they do rely on approximations of it, so this regularization effect still applies.

Since $\tilde{\Gamma}$ in equation (13) is block diagonal, PCA can be performed in parallel for different configurations. The projection (14) maximizes the projected variance for each configuration; it is different than performing PCA on $\tilde{\Gamma}$ in that it does not discriminate based on the total variance of each configuration. Disparities between the two approaches are discussed in Appendix A. Finally, the regularized observation vector and noise co-variance matrix read

$$y = \begin{bmatrix} y_1 \\ \cdots \\ y_{|C|} \end{bmatrix} \in \mathbb{R}^d, \qquad \Gamma = \begin{bmatrix} \Gamma_1 & 0 \\ & \ddots & \\ 0 & & \Gamma_{|C|} \end{bmatrix} \in \mathbb{R}^{d \times d}, \tag{16}$$

which define a regularized inverse problem of the form (3)–(5). A schematic of the inverse problem construction process is given in Figure 1. The construction of y_c from each dynamical system configuration $\zeta(x_c)$ defines the observational map \mathcal{H}_c , used to obtain the forward model evaluation $\mathcal{G}_c : \mathbb{R}^p \to \mathbb{R}^{d_c}$ for the same configuration from the dynamical model. The construction of each (y_c, Γ_c) pair, and the evaluation of $\mathcal{G}_c(\cdot)$, can be done in parallel.



Figure 1: Schematic of the strategy used to construct a regularized inverse problem from observations of a dynamical system ζ . The two branches represent different configurations of the dynamical system. From left to right: (a) the observed state is obtained following Section 2.2.2 or from any observable fields for each configuration c; (b) the observed state is normalized; (c) mean and covariance of the normalized state are computed; (d) \tilde{y}_c and $\tilde{\Gamma}_c$ are projected onto a lower dimension and regularized; (e) the statistical summaries of each configuration are aggregated, defining the global inverse problem (3)–(5).

268

2.3 Bayesian interpretation of the loss and batching

Once the data and noise estimate encodings (16) have been defined, iterative methods to solve inverse problem (3)–(5) require evaluating the loss $\mathcal{L}(\theta; y)$ at each iteration, which entails running the dynamical model in all configurations C and can be very computationally demanding. A less onerous alternative is to use a mini-batch of configurations $B \subset C$ to evaluate the average configuration loss,

$$L(\theta; y_B) = \frac{1}{2|B|} \sum_{c=1}^{|B|} ||y_c - \mathcal{G}_c(\theta)||_{\Gamma_c}^2 = \frac{1}{2} \sum_{c=1}^{|B|} ||y_c - \mathcal{G}_c(\theta)||_{|B|\Gamma_c}^2,$$
(17)

which acts as a surrogate of the configuration-averaged loss $L(\theta; y) = \mathcal{L}(\theta; y)/|C|$. The 269 use of $L(\theta; y_B)$ in lieu of $L(\theta; y)$ may be regarded as using noisy evaluations of the loss 270 for each parameter update. From a Bayesian perspective, using $L(\theta; y)$ in expression (5) 271 leads to the same MAP estimator as $\mathcal{L}(\theta; y)$ but a wider uncertainty about it, since we 272 no longer consider configurations independent. This is important when interpreting the 273 posterior uncertainty. To employ the loss (17), we only need to use the scaling $\Gamma_c \rightarrow$ 274 $|B|\Gamma_c$; to approximate the aggregate loss $\mathcal{L}(\theta, y)$ when batching, we can use $\Gamma_c \to (|B|/|C|)\Gamma_c$ 275 instead. 276

Batching is widely employed in data assimilation (Houtekamer & Mitchell, 2001) 277 and deep learning, where it has been shown to help avoid convergence to local minima 278 that generalize poorly (M. Li et al., 2014; Keskar et al., 2016). Understanding the be-279 havior of algorithms when using mini-batches is crucial for online learning, where obser-280 vations become available sequentially and the full loss cannot be sampled. Moreover, it 281 provides insight into the appropriateness of training sequentially on seasonal or geograph-282 ically sparse data in Earth system modeling applications. We explore the effect of batch-283 ing on the solution of the inverse problem in Section 4.2, training sequentially on ran-284 domly sampled configurations with markedly different dynamics. 285

²⁸⁶ **3** Ensemble Kalman methods

We consider two highly parallelizable gradient-free algorithms to solve the inverse 287 problem defined in Section 2: ensemble Kalman inversion (EKI, Iglesias et al., 2013) and 288 unscented Kalman inversion (UKI, Huang, Schneider, & Stuart, 2022). Both algorithms 289 are based on the extended Kalman filter and draw heavily on Gaussian conditioning for 290 their derivation: underlying their update rules is the approximation of the parameter dis-291 tribution as Gaussian. They afford a Bayesian interpretation when augmented with prior 292 information at every iteration (Huang, Huang, et al., 2022); how to do this is discussed 293 in Section 3.2. If prior information is not used, which may be desirable when training 294 for instance neural networks, they can be regarded as derivative-free methods to obtain 295 the MLE. 296

EKI and UKI have been used succesfully in a wide variety of inverse problems (Iglesias et al., 2013; Iglesias, 2016; Xiao et al., 2016; Kovachki & Stuart, 2019; Huang, Schneider, & Stuart, 2022). We demonstrate them here in the context of training models that may experience numerical instabilities for a priori unknown parameter combinations, starting with a brief review of the algorithms.

3.1 Ensemble Kalman inversion (EKI)

Ensemble Kalman inversion searches for the optimal θ^* given an inverse problem (3)–(5) through iterative updates of an initial parameter ensemble $\Theta_0 = [\theta_0^{(1)}, \ldots, \theta_0^{(J)}]$, used to obtain empirical estimates of covariances between parameters and the model output at each step of the algorithm. We form the initial ensemble by randomly sampling J parameter vectors $\theta_0^{(j)} \in \mathbb{R}^p$ from a Gaussian $\mathcal{N}(m_0, \Sigma_0)$. The EKI update equation for the ensemble at iteration n is (Schillings & Stuart, 2017)

$$\Theta_{n+1} = \Theta_n + \operatorname{Cov}(\theta_n, \mathcal{G}_n) \left[\operatorname{Cov}(\mathcal{G}_n, \mathcal{G}_n) + \Delta t^{-1} \Gamma \right]^{-1} \varepsilon(\Theta_n),$$
(18)

where $\Theta_n \in \mathbb{R}^{p \times J}$, Δt is the nominal learning rate of the algorithm, and $\varepsilon(\Theta_n) \in \mathbb{R}^{d \times J}$ encodes the mismatch between the forward model evaluations and the data,

$$\varepsilon(\Theta_n) = [y_{n+1}^{(1)} - \mathcal{G}(\theta_n^{(1)}), \dots, y_{n+1}^{(J)} - \mathcal{G}(\theta_n^{(J)})],$$
(19)

where

302

$$y_{n+1}^{(j)} = y + \xi_{n+1}^{(j)}, \qquad \xi_{n+1}^{(j)} \sim \mathcal{N}(0, \Delta t^{-1} \Gamma).$$
 (20)

All covariances in update (18) are estimated as sample covariances of the J ensemble members,

$$\operatorname{Cov}(\theta_n, \mathcal{G}_n) = \frac{1}{J} \left(\Theta_n - \frac{1}{J} \sum_j \theta_n^{(j)} \mathbf{1}^T \right) \left(\mathcal{G}_{\Theta_n} - \frac{1}{J} \sum_j \mathcal{G}(\theta_n^{(j)}) \mathbf{1}^T \right)^T,$$
(21)

$$\operatorname{Cov}(\mathcal{G}_n, \mathcal{G}_n) = \frac{1}{J} \Big(\mathcal{G}_{\Theta_n} - \frac{1}{J} \sum_j \mathcal{G}(\theta_n^{(j)}) \mathbf{1}^T \Big) \Big(\mathcal{G}_{\Theta_n} - \frac{1}{J} \sum_j \mathcal{G}(\theta_n^{(j)}) \mathbf{1}^T \Big)^T,$$
(22)

where $\mathcal{G}_{\Theta_n} = [\mathcal{G}(\theta_n^{(1)}), \dots, \mathcal{G}(\theta_n^{(J)})]$, and $\mathbf{1} \in \mathbb{R}^J$ is the all-ones vector. Note that the

sample covariances (21) and (22) have at most ranks $\min(\min(d, p), J-1)$ and $\min(d, J-1)$

1), respectively. From definitions (14) and (16), rank(Γ) = d by construction, so the linear system in (18) is well-defined even for J < d.

Through iterative application of the update (18), the ensemble Θ minimizes the 307 projection of the model-data mismatch on the linear span of its J members. In this study, 308 we limit the use of EKI and UKI to the calibration of dynamical models for which us-309 ing an ensemble size $J \sim p$ is feasible. For models with a large number of parameters, 310 localization or sampling error correction techniques can be used to maintain performance 311 with $J \ll p$ members (Lee, 2021; Tong & Morzfeld, 2022), like in EnKF for data as-312 313 similation (Anderson, 2012). The update (18) also drives the ensemble toward consensus, in the sense that $|\operatorname{Cov}(\theta_n, \mathcal{G}_n)| \to 0$ as $n \to \infty$; a popular method to control col-314 lapse speed is additive inflation (Anderson & Anderson, 1999; Tong & Morzfeld, 2022). 315 This collapse property precludes obtaining information about parameter uncertainties 316 directly from EKI. However, the sequence of parameter-output pairs $\{\Theta_n, \mathcal{G}_{\Theta_n}\}$ can be 317 used to train emulators for uncertainty quantification (Cleary et al., 2021). 318

319

3.1.1 Addressing model failures within the ensemble

For some parameters θ_f , simulations may be physically or numerically unstable. For instance, the Courant–Friedrichs–Lewy condition in fluid solvers may change nonlinearly with model parameters, or the initialized weights from a neural network parameterization may lead to unstable trajectories. In such situations, we need to modify update (18) to account for model failures within the ensemble.

Here we propose a novel failsafe EKI update based on the successful parameter ensemble. Let $\Theta_{s,n} = [\theta_{s,n}^{(1)}, \ldots, \theta_{s,n}^{(J_s)}]$ be the successful ensemble, for which each evaluation $\mathcal{G}(\theta_{s,n}^{(j)})$ is stable or physically consistent, and let $\theta_{f,n}^{(k)}$ be the ensemble members for which the evaluation of the forward model $\mathcal{G}(\theta_{f,n}^{(k)})$ fails. We update the successful ensemble $\Theta_{s,n}$ to $\Theta_{s,n+1}$ using expression (18), and redraw each failed ensemble member from a Gaussian defined by the successful ensemble

$$\theta_{f,n+1}^{(k)} \sim \mathcal{N}\left(m_{s,n+1}, \Sigma_{s,n+1}\right),\tag{23}$$

where

$$m_{s,n+1} = \frac{1}{J_s} \sum_{j=1}^{J_s} \theta_{s,n+1}^{(j)}, \qquad \Sigma_{s,n+1} = \operatorname{Cov}(\theta_{s,n+1}, \theta_{s,n+1}) + \kappa_*^{-1} \mu_{s,1} I_p$$
(24)

are the sample mean and regularized sample covariance matrix of the updated successful ensemble. In expression (24), κ_* is a limiting condition number and $\mu_{s,1}$ is the largest eigenvalue of the sample covariance $\text{Cov}(\theta_{s,n+1}, \theta_{s,n+1})$. This update has proved very effective for us in practice, even in situations where $J_s < J/2$; we use it throughout Section 4. The failsafe update may be combined with other conditioning techniques at initialization. For instance, the initial ensemble Θ_0 may be drawn recursively until the number of failed members is reduced below an acceptable threshold.

332

3.2 Bayesian regularization in ensemble Kalman methods

EKI implicitly regularizes the inverse problem by searching for the optimal solution θ^* over the finite-dimensional space spanned by the initial ensemble. Although UKI does not share this property, both algorithms can be equipped with Bayesian regularization by considering the augmented data-model relation (Chada et al., 2020)

$$y_a = \mathcal{G}_a(\theta) + \xi := \begin{bmatrix} y \\ m_p \end{bmatrix} = \begin{bmatrix} \mathcal{G}(\theta) \\ \theta \end{bmatrix} + \begin{bmatrix} \hat{\delta} + \hat{\eta} \\ \lambda \end{bmatrix}, \qquad (25)$$

instead of expression (3). Here, $m_p \in \mathbb{R}^p$ is the parameter prior mean, $\lambda \sim \mathcal{N}(0, 2\Lambda)$ defines the degree of regularization, $\hat{\delta} + \hat{\eta} \sim \mathcal{N}(0, 2\Gamma)$, and $\xi \sim \mathcal{N}(0, \Gamma_a)$ is the augmented error defined by relation (25). In practice, using expression (25) amounts to substituting $\{\mathcal{G}, y, \Gamma\}$ by $\{\mathcal{G}_a, y_a, \Gamma_a\}$ in both algorithms. The Kalman inversion solution to the inverse problem defined by relation (25) then satisfies

$$\theta^* = \arg\min_{\theta} \left[\mathcal{L}(\theta; y) + \frac{1}{2} ||\theta - m_p||_{\Lambda}^2 \right].$$
(26)

From a Bayesian perspective, the solution (26) approximately maximizes the posterior 333 density (5) for the Gaussian prior $\rho_{\text{prior}} \sim \mathcal{N}(0, \Lambda)$. This is particularly interesting for 334 UKI, which provides parametric uncertainty estimates (Huang, Huang, et al., 2022). When 335 using a nominal learning rate $\Delta t \neq 1$, the scaling $\Lambda \to \Delta t \cdot \Lambda$ must be used to retain 336 the Bayesian interpretation of Λ as the prior variance, due to the fact that Δt effectively 337 modifies the noise in update (18) to be $\Delta^{-1}\Gamma$. As noted before, if the original data-model 338 relation (3) is used instead of the augmented relation (25), UKI and EKI provide max-339 imum likelihood estimators. 340

341

3.3 Unscented Kalman inversion (UKI)

Unscented Kalman inversion seeks a Gaussian approximation of the posterior $\rho(\theta|y, \Gamma)$ around the MAP (given relation (25)), or an approximation of the likelihood around the MLE (given (3)), by deterministically evolving an initial Gaussian estimate $\mathcal{N}(m_0, \Sigma_0)$ through updates

$$m_{n+1} = m_n + \operatorname{Cov}_q(\theta_n, \mathcal{G}_n) \left[\operatorname{Cov}_q(\mathcal{G}_n, \mathcal{G}_n) + \Delta t^{-1} \Gamma \right]^{-1} \varepsilon(m_n),$$
(27)

$$\Sigma_{n+1} = (1 + \Delta t)\Sigma_n - \operatorname{Cov}_q(\theta_n, \mathcal{G}_n) \left[\operatorname{Cov}_q(\mathcal{G}_n, \mathcal{G}_n) + \Delta t^{-1}\Gamma \right]^{-1} \operatorname{Cov}_q(\theta_n, \mathcal{G}_n)^T, \quad (28)$$

where m_n and Σ_n are the mean and covariance estimates of the Gaussian after n iterations of the algorithm, and $\varepsilon(m_n) = y - \mathcal{G}(m_n)$ is the data-model mismatch of the mean estimate. The covariances $\operatorname{Cov}_q(\theta_n, \mathcal{G}_n)$ and $\operatorname{Cov}_q(\mathcal{G}_n, \mathcal{G}_n)$ in expressions (27) and (28) are computed through quadratures over 2p + 1 sigma points defined as

$$\hat{\theta}_n^{(j)} = m_n + a\sqrt{p} [\sqrt{\Sigma_n (1 + \Delta t)}]_j, \qquad 1 \le j \le p,$$

$$\hat{\theta}_n^{(j+p)} = m_n - a\sqrt{p} [\sqrt{\Sigma_n (1 + \Delta t)}]_j, \qquad 1 \le j \le p,$$
(29)

where $[\sqrt{\Gamma}]_j$ is the *j*-th column of the Cholesky factor of Γ , $a = \min(\sqrt{4/p}, 1)$ is a hyperparameter defined in Huang, Schneider, and Stuart (2022), and $\hat{\theta}_n^{(0)} = m_n$ is the central sigma point. The quadratures are then defined as

$$\operatorname{Cov}_{q}(\theta_{n}, \mathcal{G}_{n}) = \sum_{j=1}^{2p} w_{j}(\hat{\theta}_{n}^{(j)} - m_{n})(\mathcal{G}(\hat{\theta}_{n}^{(j)}) - \mathcal{G}(m_{n}))^{T},$$
(30)

$$\operatorname{Cov}_{q}(\mathcal{G}_{n},\mathcal{G}_{n}) = \sum_{j=1}^{2p} w_{j}(\mathcal{G}(\hat{\theta}_{n}^{(j)}) - \mathcal{G}(m_{n}))(\mathcal{G}(\hat{\theta}_{n}^{(j)}) - \mathcal{G}(m_{n}))^{T},$$
(31)

where $w_j = (2a^2p)^{-1}$ are the quadrature weights.

A limitation of this algorithm is that the number of sigma points scales linearly with p, which precludes its use when training models with a large number of parameters. However, for situations where using an ensemble of 2p + 1 members is tractable, UKI improves upon EKI by providing uncertainty quantification, instead of collapsing to a point estimate. In particular, when updates (27) and (28) are applied to the augmented datamodel relation (25), UKI ensures that Σ_n in the limit $n \to \infty$ converges towards a Gaussian estimate of parametric uncertainty (Huang, Schneider, & Stuart, 2022),

$$\Sigma_{\infty} \approx \operatorname{Cov}_{q}(\theta_{\infty}, \mathcal{G}_{a,\infty}) \left[\Delta t \cdot \operatorname{Cov}_{q}(\mathcal{G}_{a,\infty}, \mathcal{G}_{a,\infty}) + \Gamma_{a} \right]^{-1} \operatorname{Cov}_{q}(\theta_{\infty}, \mathcal{G}_{a,\infty})^{T}, \qquad (32)$$

which involves the augmented forward model $\mathcal{G}_{a}(\cdot)$ and covariance Γ_{a} defined in Section 344 3.2. Σ_{∞} approximates the covariance of the posterior (5) around m_{∞} if the full loss is 445 evaluated at every UKI iteration and $\Delta t = 1$ (Huang, Huang, et al., 2022). When batch-446 ing, an equivalent approximation can be recovered by using $\Delta t = |C|/|B|$ to compen-447 sate for sampling errors in the construction of the empirical covariances (30) and (31); 448 this is demonstrated in Section 4.2.

Finally, note that the limit (32) does not depend on Σ_0 , only on the Bayesian prior covariance Λ . This enables using a tight initial guess (i.e., $\operatorname{tr}(\Sigma_0) \ll \operatorname{tr}(\Lambda)$), which can reduce the fraction of model failures within the ensemble. To ensure robustness to the model failures that may still arise, we propose a modification of the UKI dynamics robust to model failures, similar to the one proposed for EKI, in Appendix B.

³⁵⁴ 4 Application to an atmospheric subgrid-scale model

In this section, the framework and algorithms discussed in Sections 2 and 3 are used 355 to learn closure parameters within an EDMF scheme of atmospheric turbulence and con-356 vection. The EDMF scheme is derived by spatially filtering the Navier-Stokes equations 357 for an anelastic fluid, and then decomposing the subgrid flow into n > 1 distinct sub-358 domains with moving boundaries (Cohen et al., 2020). In practice, the subdomain de-359 composition requires the use of n-1 additional equations per grid-mean prognostic field, 360 and n-1 additional equations tracking the volume fraction of each subdomain within 361 the grid (Tan et al., 2018). We retain second-order moments for one of the subdomains, 362 the environment. Covariances within the other subdomains (updrafts) are neglected, which 363 circumvents the need for turbulence closures therein. In the end, the EDMF equations require closure for the turbulent diffusivity and dissipation in the environment, and the 365 mass, momentum, and tracer fluxes between environment and updrafts. In what follows, 366 we consider an EDMF scheme with a single updraft (n = 2). 367

We consider the EDMF scheme discussed in Cohen et al. (2020): Lopez-Gomez et 368 al. (2020), which is implemented in a single-column model (SCM). Within this SCM, we 369 first seek to learn 16 closure parameters: 5 describing turbulent mixing, dissipation, and 370 mixing inhibition by stratification (Lopez-Gomez et al., 2020), 3 describing the momen-371 tum exchange between subdomains (He et al., 2021), 7 describing entrainment fluxes be-372 tween updrafts and the environment (Cohen et al., 2020), and another one defining the 373 surface area fraction occupied by updrafts. In Section 4.4, we substitute the empirical 374 dynamical entrainment closure proposed in Cohen et al. (2020) by a neural network, and 375 train the resulting physics-based machine-learning model. 376

To showcase the versatility of the algorithms, UKI is used for approximate Bayesian 377 inference of empirical parameters (using relation (25)), and EKI is used for both MAP 378 estimation of empirical parameters (relation (25), Sections 4.2, 4.3) and MLE estima-379 tion of neural network parameters (relation (3), Section 4.4). In all cases, we employ our 380 failsafe modifications of the algorithms (Section 3.1.1 and Appendix B). The name, prior 381 range U, and reference to the definition of each empirical parameter in the literature are 382 given in Table 1. The prior mean is taken to be equal to the parameter values used in 383 Lopez-Gomez et al. (2020) and Cohen et al. (2020). The prior in unconstrained space 384 $\mathcal{N}(m_p,\Lambda)$ is obtained from the physical prior mean and range through transformations 385 defined in Appendix C. Finally, we initialize EKI ensembles from the prior, $\mathcal{N}(m_0, \Sigma_0) \equiv$ 386 $\mathcal{N}(m_p, \Lambda)$, and all UKI sigma points from a tighter initial guess $\mathcal{N}(m_p, \Lambda/16)$ to demon-387 strate the ability of UKI to decouple from the initial guess. 388

389

4.1 Description of LES data and model configurations

The data used for training and testing the EDMF scheme are taken from the LES library described in Shen et al. (2022). This library contains high-resolution simulations

Symbol	Description	Prior range	Prior mean
$\overline{c_m}$	Eddy viscosity coefficient	(0.01, 1.0)	0.14, LG2020
c_d	Turbulent dissipation coefficient	(0.01, 1.0)	0.22, LG2020
c_b	Static stability coefficient	(0.01, 1.0)	0.63, LG2020
$Pr_{t,0}$	Neutral turbulent Prandtl number	(0.5, 1.5)	0.74, LG2020
κ_*	Ratio of rms turbulent velocity to friction velocity	(1.0, 4.0)	1.94, LG2020
c_{ε}	Entrainment rate coefficient	(0, 1)	0.13, C2020
c_{δ}	Detrainment rate coefficient	(0, 1)	0.51, C2020
c_{γ}	Turbulent entrainment rate coefficient	(0, 10)	0.075, C2020
β	Detrainment relative humidity power law	(0, 4)	2.0, C2020
μ_0	Entrainment sigmoidal activation parameter	$(10^{-6}, 10^{-2})$	$4 \cdot 10^{-4}, C2020$
χ_i	Updraft-environment buoyancy mixing ratio	(0, 1)	0.25, C2020
c_{λ}	Turbulence-induced entrainment coefficient	(0, 10)	0.3, C2020
a_s	Updraft surface area fraction	(0.01, 0.5)	0.1, C2020
α_b	Updraft virtual mass loading coefficient	(0, 10)	0.12, H2021
α_a	Updraft advection damping coefficient	(0, 100)	0.001, H2021
α_d	Updraft drag coefficient	(0, 50)	10.0, H2021

Table 1: Parameters ϕ considered for calibration in this study. The prior mean values are taken from LG2020 (Lopez-Gomez et al., 2020), C2020 (Cohen et al., 2020) and H2021 (He et al., 2021), where a physical description of the parameters may be found.

of low-level clouds spanning the stratocumulus-to-cumulus transition in the East Pacific
Ocean. The large-scale forcing used for these simulations is derived from the cfSites output of the HadGEM2-A model, retrieved from the Coupled Model Intercomparison Project
Phase 5 (CMIP5) archive. In particular, the monthly climatology of the cfSites output
is computed over the 5-year period 2004-2008, and used to initialize and force large-eddy
simulations for a period of 6 days. Radiative forcing is computed interactively using the
Rapid Radiative Transfer Model (RRTM, Mlawer et al., 1997).

The SCM runs are initialized from the coarse-grained LES fields after 5.75 days of simulation and are run for 6 hours. This runtime was chosen to be much longer than the equilibration time of the SCM to the steady forcing; experiments using a runtime of 12 hours resulted in no statistical changes of the results. Large-scale forcing is identical to that of the LES, and the radiative heating rates are given by the horizontal mean of the rates experienced by the high-resolution simulations. The observational map used to define the inverse problem follows the guidelines of Section 2.2, using time and horizontally averaged vertical profiles from the last $T_c = 3$ hours of simulation, at a vertical resolution of $\Delta z = 50$ m; this is also the resolution of the SCM simulations, which employ 80 vertical levels. Following the strategy in Figure 1, we extract the observations from each configuration as

$$S_c = [\bar{u}, \bar{v}, \bar{s}, \bar{q}_l, \bar{q}_t, \overline{w'q'_t}, \overline{w's'}]^T,$$
(33)

where (.) denotes time and horizontal averaging, \bar{u} and \bar{v} are the horizontal velocity com-399 ponents, \bar{s} is the entropy, \bar{q}_t is the total specific humidity, $w'q'_t$ and w's' are vertical fluxes 400 of moisture and entropy, and \bar{q}_l is the liquid water specific humidity. The pooled vari-401 ances for normalization and covariance matrix $\tilde{\Gamma}_c$ associated with the observed state S_c 402 are obtained from the full 6 day statistics of the LES to capture the internal variabil-403 ity of the system. Finally, a low-dimensional encoding is obtained from the normalized 404 time-averaged observations through truncated PCA as in equation (14), truncating the 405 dimension of the noise covariance matrix so as to preserve 99% of the total noise vari-406

ance. Calibration results using fewer observed fields at a coarser resolution are discussed
 in Section 4.3.

As training data we include a total of 60 LES configurations from the Atmospheric 409 Model Intercomparison Project (AMIP) experiment, spanning the months of January, 410 April, July and October, and locations from the coasts of Peru and California to the trop-411 ical Pacific. Results are also shown for a validation set, which includes January and July 412 simulations from an AMIP4K experiment, where sea surface temperature is increased 413 by 4 K with respect to AMIP. This temperature increase leads to 10-20% weaker large-414 415 scale subsidence, higher cloud tops, and reduced cloud cover; see Shen et al. (2022) for a detailed comparison. Validation results are representative of the generalizability of the 416 trained model for the simulation of a warming climate; the model was not trained on these 417 warmer conditions. 418

419

4.2 Calibration using mini-batch loss evaluations

To demonstrate the effectiveness of Kalman inversion in settings where evaluating 420 all configurations of interest per iteration may be too expensive or impossible (e.g., due 421 to sequential data availability), we present calibration results using mini-batches. Batch-422 ing introduces noise in the loss evaluations due to sampling error. For this reason, the 423 behavior of Kalman inversion algorithms using mini-batches is representative of their ro-424 bustness to other sources of noise, such as noise in the data or stochasticity of the dy-425 namical model. To correct for sampling noise due to batching, we use $\Delta t = |C|/|B|$ as 426 discussed in Section 3.3. 427

For training, data are fed to the algorithm by drawing |B| configurations randomly 428 and without replacement from the training set at every iteration. Configurations are reshuf-429 fled at the end of every epoch (i.e., every full pass through the training set). Figure 2 430 shows the evolution of the training and validation errors for UKI and EKI, using train-431 ing batches of 5 and 20 configurations. Since the total number of configurations in the 432 training set is 60, an epoch requires 12 iterations when using |B| = 5 and 3 when us-433 ing |B| = 20. For many geophysical applications, the cost of evaluating an ensemble 434 of long-term statistics $\mathcal{G}(\cdot)$ from a forward model is significantly higher than perform-435 ing the inversion updates (18) or (27). In these situations, a training epoch has similar 436 computational cost for any value of |B|. 437

The training error is evaluated in normalized physical space with respect to the current batch,

$$MSE(\theta; \tilde{y}_B) = \frac{1}{\tilde{d}_B} ||\tilde{y}_B - \tilde{\mathcal{G}}_B(\theta)||^2 = \frac{1}{\sum_{c=1}^{|B|} \tilde{d}_c} \sum_{c=1}^{|B|} ||\tilde{y}_c - \tilde{\mathcal{G}}_c(\theta)||^2,$$
(34)

where $\tilde{y}_B \in \mathbb{R}^{\tilde{d}_B}$. The validation error is defined similarly, but it is computed over the 438 entire validation set at every iteration. Thus, variations in the validation error are only 439 due to changes in the model parameters; there is no random data sampling. The train-440 ing and validation errors decrease sharply during the first epoch (Fig. 2). Subsequent 441 epochs fine-tune the model parameters, further reducing the data-model mismatch. It 442 is remarkable and important that the validation error decreases by about the same mag-443 nitude as the training error, demonstrating that the parameterization approach that lever-444 ages a physical model generalizes successfully out of the present-climate training sam-445 ple to a warmer climate. 446

Both EKI and UKI display efficient training in the low batch-size regime: the validation error tends to be lower for smaller batches after a fixed number of epochs. Hence, decreasing batch size in EKI and UKI can help reduce the computational cost of training models. The optimal batch size will depend on the CPU and wall-clock time con-



Figure 2: Batch (a) training and (b) validation MSE as defined in equation (34). Lines represent the error of the ensemble mean $\bar{\theta}$, MSE($\bar{\theta}$; \tilde{y}_B), and the shading represents the ensemble standard deviation of MSE(θ ; \tilde{y}_B) around the optimal point estimate $\bar{\theta}$. All errors are normalized with respect to the largest initial MSE_v($\bar{\theta}$; \tilde{y}_B), so they can be compared. Results are shown for EKI and UKI, using J = 2p + 1 and training batch sizes |B| = 5, 20. Errors for |B| = 5 are averaged using a rolling mean of 20 configurations to enable comparison with |B| = 20. In (b), the inset focuses on the validation error evolution for a longer training period.

straints of the user. Although using smaller batches reduces CPU time, it requires more
 serial operations, so using larger batches can reduce wall-clock time.

The sampling noise due to the use of different configurations (e.g., stratocumulus 453 versus cumulus regimes) increases for smaller batches. Although both algorithms achieve 454 convergence for a wide range of batch sizes, we find that EKI is more robust than UKI 455 to high levels of noise. This is shown in the inset of Figure 2b for |B| = 5, and in Ap-456 pendix D for |B| = 1. Other differences between UKI and EKI are observed in Figure 457 2. The consensus property of EKI leads to a collapse of the model error spread after a 458 few iterations, converging to a point estimate. On the other hand, the UKI ensemble con-459 verges to an MSE spread characteristic of the parameter uncertainty as approximated 460 by the distribution $\mathcal{N}(m_n, \Sigma_n)$. 461

The evolution of the parameter estimate (m_n, Σ_n) is depicted in Figure 3 for the 462 turbulent dissipation c_d , updraft advection damping α_a and surface area fraction a_s . The 463 initial parameter estimate depends on the stochastic initialization for EKI. The UKI es-464 timate provides information about parameter uncertainty, whereas EKI only provides 465 a point estimate (i.e., m_n). From the UKI estimate, we observe that the training set con-466 strains the likely values of the turbulent dissipation (c_d) and surface area fraction (a_s) 467 to a significantly smaller region than the prior. However, the magnitude of updraft ad-468 vection damping (α_a) is not identifiable using this dataset: the corresponding diagonal 469 element of Σ_n converges to the prior variance used in the regularized problem (25) (Fig-470 ure 3b). 471

The covariance estimate Σ_n also provides information about correlations between model parameters and total reduction of uncertainty (Figure 4). For the current stratocumulusto-cumulus transition dataset, our EDMF model shows moderate correlations between parameters regulating the turbulence kinetic energy budget in the boundary layer (c_b, c_m, c_d , see Lopez-Gomez et al., 2020). We also find entrainment to be negatively correlated with



Figure 3: Parameter evolution of the turbulent dissipation (a), updraft advection damping (b), and updraft surface area fraction (c). All values are given in physical space. The solid lines describe the trajectories of the mean estimate, $\mathcal{T}^{-1}(m_n)$. For UKI, the marginal $\pm \sigma$ uncertainty band is included in shading. This uncertainty is equal to $\pm \mathcal{T}^{-1}(\sqrt{(\Sigma_n)_{i,i}})$ for the *i*-th parameter. The black dashed lines are the $\pm \sigma$ uncertainty bands of the prior used for regularization. Legend as in Figure 2.

surface updraft area fraction, detrainment and drag. These correlations can be used to improve parameterizations at the process level by identifying or developing a set of uncorrelated parameters. Figure 4b shows how Σ_n converges to a quasi-steady state estimate of the posterior covariance after ~ 30 iterations.

Vertical profiles of $\bar{q}_l, \overline{w'q'_t}$ and \bar{u} from the validation set are compared to the ref-481 erence LES profiles in Figure 5. The calibrated model yields smoother and more accu-482 rate profiles than the model before training. In particular, calibration significantly re-483 duces biases in liquid water specific humidity and moisture transport for both stratocu-484 mulus and cumulus cloud regimes in the 4 K-warmer AMIP4K experiment. These re-485 sults confirm that the dynamical model can be trained using a low-dimensional encod-486 ing of the time statistics, as proposed in Section 2. They also highlight the generalizabil-487 ity of sparse physics-based models. 488

489

4.3 Calibration using partial observations

Another application of synthetic high-resolution data is the study of calibration sensitivity to data resolution and partial loss of information. Such sensitivity studies can inform the technical requirements of future observing systems or field campaigns (Suselj et al., 2020), and are easily implemented with ensemble and unscented Kalman inversion through modifications of the observational map \mathcal{H} .

Here, we employ the EKI and UKI algorithms for this task by using coarser train-495 ing data at a vertical resolution of $\Delta z = 200$ m. In addition, we consider only a sub-496 set of fields for which observational data may be obtained in practice: the liquid water 497 potential temperature θ_l , the total water specific humidity \bar{q}_t and the liquid water spe-498 cific humidity \bar{q}_l (National Academies of Sciences, Engineering, and Medicine, 2018; Suselj 499 et al., 2020). Figure 6 compares calibration results using this reduced setup with the re-500 sults obtained using the full high-resolution observations of Section 4.2. The loss of in-501 formation is evident in the inability of the algorithms to find the same minimum reached 502



Figure 4: Parameter correlations estimated from UKI using |B| = 20 (a), and evolution of the total parameter variance from UKI using |B| = 20, 10 and 5, normalized by the prior variance tr(Λ) = 16 (b). Note that the initial covariance estimate used in UKI (with tr(Σ_0) = 1) is decoupled from the prior. Symbols follow Table 1.

with richer observations. Nevertheless, Kalman inversion significantly reduces the validation error from the prior even with sparser data and a limited number of fields.

The identifiability of individual parameters as a function of the observational map 505 \mathcal{H} can be inferred from the UKI Σ_n diagnostic. Figure 6 shows that the partial obser-506 vations of temperature and humidity are enough to constrain the entrainment coefficient 507 in the EDMF scheme. However, the loss of information with respect to the original ob-508 servations leads to much poorer constraints on the turbulent dissipation coefficient. The 509 same comparison can be performed for any parameter of interest to inform observational 510 requirements to constrain models at the process level. This diagnostic is an important 511 advantage of UKI over EKI; identifiability is not directly inferable from ensemble Kalman 512 inversion due to the ensemble collapse. However, this information can be recovered through 513 the emulation of the forward map (Cleary et al., 2021). 514

The use of partial observations also highlights the benefits of learning from time 515 statistics instead of tendencies. Learning from statistics not only ensures that the cal-516 ibrated dynamical model is stable, which requires a leap of faith when training on in-517 stantaneous tendencies (Bretherton et al., 2022). It also couples the evolution of ther-518 modynamic and dynamical fields, which can improve the forecast of fields unseen dur-519 ing training. An example is shown in Figure 7. The model calibrated using thermody-520 namic profiles improves upon the prior model in the forecast of horizontal velocities within 521 the boundary and cloud layers. A common reason to use tendencies for calibration is that 522 they enable the use of supervised learning techniques, which are easy to implement for 523 neural network architectures (e.g., Bretherton et al., 2022). In the next subsection, we 524 demonstrate the power of UKI and EKI to calibrate hybrid models with embedded neu-525 ral network parameterizations. 526

527 528

4.4 Calibration of a hybrid model with embedded neural network closures

We consider now a hybrid EDMF scheme that substitutes the dynamical entrainment and detrainment closures proposed by Cohen et al. (2020) with a three-layer dense



Figure 5: Prior, posterior and LES profiles of liquid water specific humidity (\bar{q}_l), subgridscale moisture flux ($\overline{w'q'_l}$) and zonal velocity (\bar{u}) for cfSites 5 (top) and 14 (bottom) using July forcing from the AMIP4K experiment as in Shen et al. (2022). The gray shading represents the internal variability of the LES simulations over 6 days of steady forcing, and the full lines represent 3-hour time-averaged profiles. EKI prior and posterior results are point estimates evaluated at the parameter vector closest to the ensemble mean. The UKI posterior shading spans the central 68% of the profile posterior distribution. All Kalman methods used |B| = 5 and J = 2p + 1.

neural network. We define the fractional entrainment (ϵ) and detrainment (δ) rates as

$$\begin{bmatrix} \epsilon \\ \delta \end{bmatrix} = \frac{1}{z} \text{NN}_3(\Pi_1, \dots, \Pi_6), \tag{35}$$

where z is the height, and the hidden layers of NN₃ have 5 and 4 nodes, from input to outputs. Our closure (35) seeks to learn local expressions for the z-normalized entrainment/detrainment rates, which have been shown to vary weakly in empirical studies of shallow cumulus convection (Siebesma, 1996; de Roode et al., 2000). The neural network inputs Π_1, \ldots, Π_6 are 6 nondimensional groups on which entrainment and detrainment may depend, defined as

$$\Pi_1 = \frac{z(b_{up} - b_{en})}{(w_{up} - w_{en})^2 + w_d^2},$$
(36a)

$$\Pi_2 = \frac{a_{up}w_{up}^2 + (1 - a_{up})w_{en}^2}{2(1 - a_{up})e_{en} + a_{up}w_{up}^2 + (1 - a_{up})w_{en}^2},$$
(36b)

$$\Pi_3 = \sqrt{a_{up}},\tag{36c}$$



Figure 6: Evolution of the validation error (a) and estimates of the turbulent dissipation (b) and entrainment coefficient (c) for calibration processes using observations of the state (33) at 50 m resolution (UKI_f, EKI_f), or from $\bar{\theta}_l, \bar{q}_t$ and \bar{q}_l at 200 m resolution (UKI_o, EKI_o). All inversion processes use |B| = 20. Shading is defined as in Figures 2 and 3.

$$\Pi_4 = \mathrm{RH}_{up} - \mathrm{RH}_{en},\tag{36d}$$

$$\Pi_5 = z/H_{up},\tag{36e}$$

$$\Pi_6 = gz/R_d T_{\rm ref}.\tag{36f}$$

In expressions (36), $w_d = (H_{inv}\overline{w'b'}|_s)^{1/3}$ is the Deardorff convective velocity, H_{inv} is the inversion height, $\overline{w'b'}|_s$ is the surface buoyancy flux, g is the gravitational acceleration, R_d is the ideal gas constant for dry air and T_{ref} is a reference temperature. The subscripts up and en denote updraft and environment: a_{up} is the updraft area fraction, H_{up} the updraft top height, and e_{en} the environmental turbulence kinetic energy. The relative humidity RH, vertical velocity with respect to the grid mean w, and buoyancy b are defined for both updraft and environment.

The neural network closure (35) introduces 63 additional coefficients with respect 536 to the entrainment and detrainment closure calibrated in Sections 4.2 and 4.3, for a to-537 tal of 79 parameters. As the closure complexity increases, it is most practical to use EKI 538 for calibration, since it enables the use of ensembles with J < 2p+1 members. In Fig-539 ure 8, we present training and validation errors for the hybrid model using ensemble sizes 540 J = 50, 100, and 159, and for the empirical EDMF scheme with J = 2p + 1 = 33 en-541 semble members. We initialize the neural network weights as $\theta_{\rm NN} \sim \mathcal{N}(\theta_{\rm NN}^0, I)$ with 542 $\theta_{\rm NN}^0 \sim U(-0.05, 0.05)$. In all cases, we use Bayesian regularization as discussed in Sec-543 tion 4.2 for all model parameters except for the neural network weights. We calibrate 544 all parameters of the empirical and hybrid models, to compare the optimal performance 545 of both closures. 546

Both the empirical and hybrid EDMF schemes generalize well to the validation set, 547 with training and validation errors reaching levels of about 5% of the largest a priori val-548 idation error. The strong generalization to 4 K-warmer cloud regimes contrasts with re-549 sults from approaches that try to learn unresolved tendencies directly, without encod-550 ing the physics (Rasp et al., 2018). Using a physics-based hybrid approach, all learned 551 closures are consistently placed within the coarse-grained dynamics of the system (Cohen 552 et al., 2020), which also vastly reduces data requirements. Further, targeting closure terms 553 that isolate a single physical process lends itself to interpretability in a manner difficult 554



Figure 7: Prior, posterior and LES profiles of liquid water specific humidity (\bar{q}_l) , vertical moisture flux $(w'q'_l)$ and zonal velocity (\bar{u}) for cfSite 3 using July forcing (top) and cfSite 14 using January forcing (bottom) from the AMIP4K experiment (Shen et al., 2022). Posterior results are shown for a model calibrated using the high-resolution state (33) (Full Obs.), and coarse-resolution observations of $\bar{\theta}_l$, \bar{q}_l and \bar{q}_l (Partial Obs.). Shadings and legend as in Figure 5. Results obtained using UKI with |B| = 20.

for purely machine-learning based parameterizations that simultaneously model many
 physical processes. After training, relationships between EDMF variables and targeted
 physical quantities like entrainment can be teased out using partial dependence plots or
 ablation studies. In addition, the learned relationships are point-wise and causal.

The inset in Figure 8b shows how the higher-complexity hybrid model moderately 559 overfits to the training set after ~ 10 epochs, a behavior that is not observed with the 560 empirical model. Hence, in the low-data regime $(d \leq p)$, adoption of techniques such 561 as early stopping (Prechelt, 1998) or sparsity-inducing regularization (Schneider et al., 562 2020) becomes necessary. The compact support property of EKI, which mandates that 563 the solution be in the linear span of the initial ensemble, also regularizes the learned hy-564 brid model with decreasing J; for J = 50 < p overfitting is significantly reduced. Thus, 565 reducing the ensemble size is an efficient regularization technique when training large 566 machine-learning models that tend to overfit, at the expense of reduced expressivity. Ad-567 ditional EKI-specific regularization techniques for deeper networks are discussed in Kovachki 568 and Stuart (2019). 569

Another difference between the empirical and the hybrid models is that for the latter, we do not know a priori the parameter ranges for which the model trajectories remain physical. During the training sessions shown in Figure 8, the hybrid models experienced a maximum of 25 (J = 50), 30 (J = 100) and 72 (J = 159) failures in a single iteration, all occurring during the first epoch. The use of the failsafe update proposed in Section 3.1.1 proved crucial to enable training in the presence of model failures, and it reduced the number of failures to a small fraction of the J ensemble members after a few EKI iterations.



Figure 8: Batch (a) training and (b) validation normalized MSE for the hybrid (EDMF+NN) and empirical (EDMF) models. Lines, shading and inset as in Figure 2. Results are shown for calibration with EKI, using J = 50,100 and 2p + 1 = 159 ensemble members for the hybrid model. The empirical model training uses J = 2p + 1 = 33. All inversion processes use batch size |B| = 10.

Profiles of \bar{q}_l , \bar{q}_t and $w'q'_t$ are shown in Figure 9 for the trained empirical and hybrid EDMF models. To produce the profiles with the hybrid model, we retain the parameters learned at the iteration with lowest validation error from a training session spanning 25 epochs, effectively similar to early stopping. As expected from the validation error, the hybrid model slightly improves upon the skill of the empirical model, predicting more accurate profiles of \bar{q}_l within the cloud layer. This is, of course, at the cost of a significantly higher parameter complexity of the closure.

As shown here, ensemble Kalman inversion allows for rapid prototyping and comparison of closures within an overarching black-box model. Importantly, this comparison can be done during training in terms of the online performance of the fully calibrated dynamical model.

589 5 Discussion and conclusions

Ensemble Kalman methods such as ensemble and unscented Kalman inversion are 590 powerful tools for training possibly expensive models. By leveraging covariances between 591 the model output and its parameters, they do not impose any constraint on the data used 592 for learning, or the architecture of the closures to be calibrated. This means that ensem-593 ble Kalman methods can be used to learn all parameters within complex overarching mod-594 els, regardless of where those parameters appear in the formulation of the model. Fur-595 thermore, the Gaussian approximation of the parameter distribution makes them far more 596 efficient than standard Bayesian inference techniques, at the cost of neglecting uncer-597



Figure 9: Prior, posterior and LES profiles of liquid water specific humidity (\bar{q}_l) , total water specific humidity (\bar{q}_t) and vertical moisture flux $(w'q'_t)$ for cfSite 14 using July forcing (top) and cfSite 8 using January forcing (bottom) from the AMIP4K experiment (Shen et al., 2022). Definitions of prior, posterior and shading as in Figure 5. Posterior results are shown for the EDMF model with empirical closures (EDMF), and with the neural network entrainment closure (35) (EDMF+NN), using early stopping and 25 epochs of training. Results obtained using EKI with |B| = 10.

tainty beyond the second moment of the posterior, and the possible convergence to local minima (as for stochastic gradient descent and other optimization methods).

This enables training physics-based machine-learning parameterizations, as demonstrated here by substituting an internal component of the EDMF model by a neural network, which required no change in the data or framework used for training. The benefits of combining physics and data are demonstrated by the performance of our trained hybrid closure in simulations of clouds typical of conditions 4 K warmer than the clouds in the training set.

To use these algorithms, parameter learning must be framed as an inverse prob-606 lem. This allows great flexibility, but raises the problem of choosing a reasonable obser-607 vational map \mathcal{H} and noise covariance Γ to define an inverse problem. Through a domain-608 agnostic strategy and a reasonable heuristic about the expected model error, we have 609 demonstrated a systematic way of constructing a well-defined inverse problem from high-610 dimensional data. This strategy is designed to maximize the information content through 611 a lossy principal component encoding \mathcal{H} and to allow the use of time averages as obser-612 vations, making it amenable to harnessing, e.g., satellite observations in addition to com-613

putationally generated data. The success of this strategy is demonstrated in a variety
 of settings, using empirical and hybrid models.

The flexibility of the inverse problem allows to define the observational map \mathcal{H} through any observable diagnostic of the model, be it differentiable or not. For instance, Barthélémy et al. (2021) use a neural network as the mapping \mathcal{H} , to train a low-resolution dynamical model directly from features at high resolution. One could also envision the construction of \mathcal{H} through other statistics of the model dynamics, such as the variance or skewness. These choices may be preferable for particular tasks, such as the prediction of extreme events or the correct representation of emergent phenomena.

Given an inverse problem, we have shown that EKI and UKI are robust to noise and amenable to batching strategies. This establishes the ability of the Kalman algorithms to train models using sequentially sampled data. The same robustness can be expected for other sources of noise, such as stochasticity in the model (Schneider, Stuart, & Wu, 2021). In addition, we have proposed modifications of the EKI and UKI updates that enable calibrating models that may fail during training, which is often the case for Earth system models.

Although similar, each ensemble Kalman algorithm presents its own relative strengths 630 in our analysis. Calibration through EKI appears to be more robust to noise, and the 631 number of ensemble members may be chosen to be lower than for UKI when the param-632 eter space is high-dimensional. Indeed, Kovachki and Stuart (2019) show successful re-633 sults for EKI when the number of parameters (e.g., $p \sim 10^6$) is two orders of magni-634 tude higher than the ensemble size. Using fewer ensemble members than parameters also 635 introduces a regularization effect. On the other hand, UKI provides information about 636 parametric uncertainty and correlations, which can be used to improve models at the pro-637 cess level, and to rapidly compare the added value of increasingly precise observing sys-638 tems. Other ensemble Kalman methods, such as the sparsity-inducing EKI (Schneider 639 et al., 2020) or the ensemble Kalman sampler (Garbuno-Inigo et al., 2020), can provide 640 solutions to the inverse problem with other useful properties. In addition, all these en-641 semble methods generate parameter-output pairs that can be used to train emulators 642 for uncertainty quantification that can capture non-Gaussian posteriors (Cleary et al., 643 2021). 644

⁶⁴⁵ Finally, ensemble Kalman methods may be used for the rapid comparison of pa⁶⁴⁶ rameterizations in terms of the online skill of an overarching Earth system model. The
⁶⁴⁷ same framework could be used to train Gaussian processes, random feature models (Nelsen
⁶⁴⁸ & Stuart, 2021), Fourier neural operators (Z. Li et al., 2020), or stochastic closures (Guillaumin
⁶⁴⁹ & Zanna, 2021), for example. These are some of the exciting research avenues that we
⁶⁵⁰ will be exploring in the future.

651

Appendix A Configuration-based principal component analysis

Performing PCA on each configuration allows retaining principal modes from low-652 variance configurations while filtering out trailing modes from high-variance configura-653 tions. The importance of this is demonstrated in Figure A1 for three configurations of 654 our LES solver (Pressel et al., 2015) based on observational campaigns of a stable bound-655 ary layer, a stratocumulus-topped boundary layer, and shallow cumulus convection (Beare 656 et al., 2006; Stevens et al., 2005; Siebesma et al., 2003). Performing global PCA is equiv-657 alent to using a cutoff $\mu_{c,i} > \mu_c^*$ in Figure A1a, where we need to choose between ne-658 glecting most modes from certain configurations (e.g., GABLS in Figure A1a) or retain-659 ing highly oscillatory modes from others (e.g., Bomex), as measured by the number of 660 zero-crossings of the eigenmode (Hansen, 1998). Highly oscillatory modes may have a 661 disproportionate contribution to the loss when calibrating imperfect models. On the other 662 hand, performing PCA on each $\tilde{\Gamma}_c$ alleviates this problem by aligning the eigenspectra 663

⁶⁶⁴ before applying the cutoff, as shown in Figure A1b. Appropriate conditioning of the global

⁶⁶⁵ covariance matrix is still enforced when applying configuration-based PCA through the Tilthonour negularizer in equation (14)

 $_{666}$ Tikhonov regularizer in equation (14).



Figure A1: (a) Scatter plot of covariance eigenvalues $\mu_{c,i}$ and the number of zerocrossings of their corresponding eigenmode for three different configurations of an LES solver. (b) The same plot, with eigenvalues normalized by the leading eigenvalue of each configuration ($\mu_{c,1}$). Trailing eigenvalues are associated with high-wavenumber oscillatory modes with frequent sign changes.

Appendix B Addressing model failures with unscented Kalman inversion

In the presence of model failures, we perform the UKI quadratures over the successful sigma points. Consider the set of off-center sigma points $\{\hat{\theta}\} = \{\hat{\theta}_s\} \cup \{\hat{\theta}_f\}$ where $\hat{\theta}_s^{(j)}, j = 1, \ldots, J_s$ are successful members and $\hat{\theta}_f^{(k)}$ are not. For ease of notation, consider an ordering of $\{\hat{\theta}\}$ such that $\{\hat{\theta}_s\}$ are its first J_s elements, and note that we deal with the central point $\hat{\theta}^{(0)}$ separately. We estimate the covariances $\text{Cov}_q(\mathcal{G}_n, \mathcal{G}_n)$ and $\text{Cov}_q(\theta_n, \mathcal{G}_n)$ from the successful ensemble,

$$\operatorname{Cov}_{q}(\theta_{n}, \mathcal{G}_{n}) \approx \sum_{j=1}^{J_{s}} w_{s,j} (\hat{\theta}_{s,n}^{(j)} - \bar{\theta}_{s,n}) (\mathcal{G}(\hat{\theta}_{s,n}^{(j)}) - \bar{\mathcal{G}}_{s,n})^{T},$$
(B1)

$$\operatorname{Cov}_{q}(\mathcal{G}_{n},\mathcal{G}_{n}) \approx \sum_{j=1}^{J_{s}} w_{s,j}(\mathcal{G}(\hat{\theta}_{s,n}^{(j)}) - \bar{\mathcal{G}}_{s,n})(\mathcal{G}(\hat{\theta}_{s,n}^{(j)}) - \bar{\mathcal{G}}_{s,n})^{T},$$
(B2)

where the weights at each successful sigma point are scaled up, to preserve the sum of weights,

$$w_{s,j} = \left(\frac{\sum_{i=1}^{2p} w_i}{\sum_{k=1}^{J_s} w_k}\right) w_j.$$
 (B3)

In equations (B1) and (B2), $\bar{\theta}_{s,n}$ and $\bar{\mathcal{G}}_{s,n}$ must be modified from the original formulation if the central point $\hat{\theta}^{(0)} = m_n$ results in model failure,

$$\bar{\theta}_{s,n} = \begin{cases} m_n & \text{if } \hat{\theta}^{(0)} \text{ successful,} \\ \frac{1}{J_s} \sum_{j=1}^{J_s} \hat{\theta}_{s,n}^{(j)} & \text{otherwise,} \end{cases}$$
(B4)

$$\bar{\mathcal{G}}_{s,n} = \begin{cases} \mathcal{G}(m_n) & \text{if } \hat{\theta}^{(0)} \text{ successful,} \\ \frac{1}{J_s} \sum_{j=1}^{J_s} \mathcal{G}(\hat{\theta}_{s,n}^{(j)}) & \text{otherwise.} \end{cases}$$
(B5)

These modified UKI quadrature rules are used throughout Section 4 to deal with model failures. Since UKI can be initialized from a tighter prior than EKI, due to the absence

of ensemble collapse, failures are much easier to avoid than with EKI.

⁶⁷² Appendix C Parameter transformation and prior

Given a prior range $[\phi_i, \phi_f]$ for a parameter $\phi \in \mathbb{R}$, we define the transformation

$$\theta = \mathcal{T}(\phi) = \ln \frac{\phi - \phi_i}{\phi_f - \phi},\tag{C1}$$

such that the interval midpoint is mapped to $\theta = 0$, and the bounds to $\pm \infty$. An unconstrained Gaussian prior may then be defined for θ given the prior mean in physical (constrained) parameter space ϕ_p as

$$\theta_0 \sim \mathcal{N}(\mathcal{T}(\phi_p), \sigma_0^2),$$
 (C2)

where σ_0^2 is a free parameter controlling the size of the region within the interval $[\phi_i, \phi_f]$ 673 containing most of the probability. This means that the magnitude of σ_0 is already nor-674 malized with respect to the prior range, so we will generally choose $\sigma_0 \sim \mathcal{O}(1)$. The 675 p-dimensional prior $\mathcal{N}(m_0, \Sigma_0)$ is then constructed as an uncorrelated multivariate nor-676 mal with marginal distributions given by expression (C2). The normalization induced 677 by (C1) also enables the use of isotropic regularization in equations (25)-(26), even though 678 the physical parameters ϕ may differ in order of magnitude. For more examples of pa-679 rameter transformations in the context of EKI and UKI, see Huang, Schneider, and Stu-680 art (2022), Schneider, Dunbar, et al. (2021), and Dunbar et al. (2022). 681

Appendix D Calibration using very noisy loss evaluations

The Kalman inversion results are expected to deteriorate above some noise thresh-683 old, as the signal-to-noise ratio in the training process decreases. We explored the sen-684 sitivity of UKI and EKI to noise by sampling a single configuration per iteration from 685 the training set described in Section 4.1. As shown in Figure D1, UKI fails to converge 686 to the minimum found with larger batches in this limit. The validation error is characterized by large oscillations due to strong changes in the value of model parameters like 688 the entrainment coefficient c_{ϵ} or the eddy diffusivity coefficient c_m . On the other hand, 689 EKI proves robust to noise even in this limit, converging to the minimum found by UKI 690 employing larger batches. 691

In the context of Kalman inversion, decreasing the step size Δt is equivalent to in-692 creasing the noise variance, as shown in updates (18) and (27). We investigate the time 693 step role in the small batch limit by performing the ensemble Kalman inversion with $\Delta t =$ 694 $|C|^{-1} = 1/60$. The smaller time step increases the parameter uncertainty, which leads 695 to a reduction in parameter oscillations and estimates closer to the prior. This is accom-696 panied by a moderate reduction in validation error oscillations. We performed additional 697 inversions using even smaller time steps, which resulted in a convergence of the param-698 eter estimates towards the prior and a minor reduction in validation error with respect 699 to the initialization. We conclude that decreasing Δt in UKI can reduce oscillations due 700 to high levels of noise, but it does not result in the same robustness as EKI. 701

702 Notation

 $\phi \in \mathbb{R}^p$ Learnable parameters, in physical space.



Figure D1: Evolution of the validation error (a) and estimates of the entrainment (b), and eddy diffusivity (c) coefficients. Results shown for UKI using batch sizes of 10 and 1, and EKI using a batch size of 1. Parameter uncertainty only shown for UKI₁₀ and UKI₁, $\Delta t = 1/60$ for clarity. All results shown use $\Delta t = |C|/|B|$ unless otherwise specified. Shading as in Figures 2 and 3.

- $\theta \in \mathbb{R}^p$ Transformed learnable parameters, in unconstrained space.
- $\theta^* \in \mathbb{R}^p$ Optimal unconstrained parameter estimate (MAP or MLE).
- 706 φ_0 Initial dynamical state.
- F_{φ} Dynamical forcing.
- 708 $x_c = \{\varphi_0, F_{\varphi}\}_c$ Configuration of the dynamical system.
- 709 $\zeta(x_c): \varphi_0 \to \varphi(t)$ True dynamical system evolution.
- 710 $\Psi(\phi; x_c): \varphi_0 \to \hat{\varphi}(t)$ Dynamical model evolution.
- 711 \mathcal{H}_c Observational map for configuration c.
- $y_c \in \mathbb{R}^{d_c}$ Observation vector for configuration c.
- $\eta_c \in \mathbb{R}^{d_c}$ Observation error for map \mathcal{H}_c .
- $\delta(x_c) \in \mathbb{R}^{d_c}$ Model or representation error for configuration c.
- ₇₁₅ $\Gamma_c \in \mathbb{R}^{d_c \times d_c}$ Covariance of the Gaussian noise $\eta_c + \delta(x_c)$.
- $\mathcal{G}_c: \mathbb{R}^p \to \mathbb{R}^{d_c}$ Forward model for configuration c.
- $C = \{x_c, c = 1, \dots, |C|\}$ Set of configurations.
- 718 $y = [y_1, \dots, y_{|C|}]^T \in \mathbb{R}^d$ Global observation vector.
- 719 $\delta = [\delta(x_1), \dots, \delta(x_{|C|})]^T$ Global representation error.
- $\eta = [\eta_1, \dots, \eta_{|C|}]^T$ Global observation error.
- $_{721}$ $\Gamma \in \mathbb{R}^{d \times d}$ Global noise covariance matrix.
- 722 $\mathcal{T}: U \to \mathbb{R}^p$ Parameter transformation to unconstrained space.
- $\mathcal{G}: \mathbb{R}^p \to \mathbb{R}^d$ Forward model.
- $\rho(\theta|y, \Gamma)$ Parameter posterior probability density, given Γ and y.
- $\rho_{\text{prior}}(\theta)$ Parameter prior probability density, independent of Γ .
- $\mathcal{L}: \mathbb{R}^p \times \mathbb{R}^d \to \mathbb{R}$ Loss or negative log-likelihood given Γ .
- $S_c(t) \in \mathbb{R}^{\tilde{d}_c}$ Observed state.
- ⁷²⁸ $V_{c,j}(t) \in \mathbb{R}^{h_c}$ Spatial field j within the observed state S_c .
- $s_c(t) \in \mathbb{R}^{\tilde{d}_c}$ Normalized observed state.
- $v_{c,j}(t) \in \mathbb{R}^{h_c}$ Spatial field j within the normalized state s_c .
- $\sigma_{c,j} \in \mathbb{R}$ Pooled time standard deviation of $V_{c,j}$.

- $T_c \in \mathbb{R}$ Time-averaging window used in map \mathcal{H}_c . 732
- $\tilde{y}_c \in \mathbb{R}^{d_c}$ Counterpart of y_c prior to encoding. 733
- $\tilde{y} \in \mathbb{R}^{\tilde{d}}$ Global observation vector prior to encoding. 734
- $\tilde{\Gamma}_{c} \in \mathbb{R}^{\tilde{d}_{c} \times \tilde{d}_{c}}$ Counterpart of Γ_{c} prior to encoding. 735
- $\tilde{\Gamma} \in \mathbb{R}^{\tilde{d} \times \tilde{d}}$ Counterpart of Γ prior to encoding. 736
- $I_d \in \mathbb{R}^{d \times d}$ Identity matrix of size $d \times d$. 737
- $\mu_{c,i} \in \mathbb{R}$ *i*-th largest eigenvalue of $\tilde{\Gamma}_c$. 738
- $\kappa \in \mathbb{R}$ Approximate condition number of a matrix. 739
- $r_{c} \in \mathbb{R}$ Approximate rank of matrix $\tilde{\Gamma}_{c}$. 740
- $\epsilon_m \in \mathbb{R}$ Machine or data precision. 741
- 742
- $\kappa_* < \epsilon_m^{-1/2}$ Limiting matrix condition number. $P_c \in \mathbb{R}^{\tilde{d_c} \times d_c}$ Truncated PCA projection matrix. 743
- $D\mathcal{G}(\theta) \in \mathbb{R}^{d \times p}$ Jacobian of forward model at θ . 744
- $B = \{x_c, c = 1, \dots, |B|\}$ Mini-batch of configurations. 745
- $L: \mathbb{R}^p \times \mathbb{R}^d \to \mathbb{R}$ Configuration-averaged loss. 746
- $y_B \in \mathbb{R}^{d_B}$ Observation vector for batch B. 747
- $\tilde{y}_B \in \mathbb{R}^{\tilde{d}_B}$ Counterpart of y_B prior to encoding. 748
- $\tilde{\tilde{\mathcal{G}}}_B : \mathbb{R}^p \to \mathbb{R}^{\tilde{d}_B}$ Forward model corresponding to observations \tilde{y}_B . 749
- $\Theta_n \in \mathbb{R}^{p \times J}$ Parameter ensemble at iteration n. 750
- $m_n \in \mathbb{R}^p$ Mean parameter estimate at iteration n. 751
- $\Sigma_n \in \mathbb{R}^{p \times p}$ Parameter covariance estimate at iteration n. 752
- $\mathcal{G}_{\Theta_n} \in \mathbb{R}^{d \times J}$ Forward model evaluation ensemble at iteration *n*. 753
- $\varepsilon(\Theta_n) \in \mathbb{R}^{d \times J}$ Data-model mismatch ensemble at iteration n. 754
- $\Delta t \in \mathbb{R}$ Nominal learning rate. 755
- $\Theta_{s,n} \in \mathbb{R}^{p \times J_s}$ Successful parameter ensemble at iteration n. 756
- $\theta_{f,n}^{(k)} \in \mathbb{R}^p$ k-th failed parameter vector at iteration n. 757
- $m_p \in \mathbb{R}^p$ Parameter prior mean. 758
- $\Lambda \in \mathbb{R}^{p \times p}$ Gaussian prior covariance. 759
- $y_a \in \mathbb{R}^{d+p}$ Observation vector augmented with m_p . 760
- $\mathcal{G}_a(\theta) \in \mathbb{R}^{d+p}$ Forward model augmented with θ . 761
- $\xi \in \mathbb{R}^{d+p}$ Aggregate noise in the augmented data-model relation. 762
- $\Gamma_a \in \mathbb{R}^{(d+p) \times (d+p)}$ Covariance of the aggregate noise ξ . 763
- $\hat{\theta}_n^{(j)} \in \mathbb{R}^p$ *j*-th sigma point for UKI quadrature. 764
- Π_i *j*-th nondimensional input to neural network. 765

Acknowledgments 766

We thank Daniel Z. Huang and Zhaovi Shen for insightful discussions, and Julien 767 Brajard and an anonymous reviewer for prompting a clearer and more precise formu-768 lation of the problem and methods discussed in this study. I.L. was supported by a fel-769 lowship from the Resnick Sustainability Institute at Caltech, and an Amazon AI4Science 770 fellowship. H.L.L.E was supported by an Aker scholarship and a Fulbright fellowship. 771 This research was additionally supported by the generosity of Eric and Wendy Schmidt 772 by recommendation of the Schmidt Futures program, by the National Science Founda-773 tion (grant AGS-1835860), by the Defense Advanced Research Projects Agency (Agree-774 ment No. HR00112290030), and by the Heising-Simons Foundation. Part of this research 775 was carried out at the Jet Propulsion Laboratory, California Institute of Technology, un-776 der a contract with the National Aeronautics and Space Administration. The software 777 package implementing ensemble Kalman methods can be accessed at https://doi.org/ 778 10.5281/zenodo.6382968, the one implementing the EDMF scheme at https://doi 779

.org/10.5281/zenodo.6392397, and the software used to calibrate the EDMF scheme

may be accessed at https://doi.org/10.5281/zenodo.6382865. The data from Shen

rez et al. (2022) used for model training is available at https://doi.org/10.22002/D1.20052.

Localization and sampling error correction in ensemble

783 **References**

784

Anderson, J. L.

(2012).

Monthly Weather Review, 140, 2359–2371. Kalman filter data assimilation. 785 doi: 10.1175/MWR-D-11-00013.1 786 Anderson, J. L., & Anderson, S. L. (1999).A Monte Carlo implementa-787 tion of the nonlinear filtering problem to produce ensemble assimilations 788 Monthly Weather Review, 127, 2741-2758. and forecasts. doi: 10.1175/ 789 1520-0493(1999)127(2741:AMCIOT)2.0.CO;2 790 Barthélémy, S., Brajard, J., Bertino, L., & Counillon, F. (2021).Super-resolution 791 data assimilation. doi: 10.48550/arxiv.2109.08017 792 Beare, R. J., Macvean, M. K., Holtslag, A. A. M., Cuxart, J., Esau, I., Golaz, J.-793 C., ... Sullivan, P. (2006).An intercomparison of large-eddy simulations of 794 the stable boundary layer. Boundary-Layer Meteorology, 118, 247–272. doi: 10.1007/s10546-004-2820-6 796 Beucler, T., Pritchard, M., Rasp, S., Ott, J., Baldi, P., & Gentine, P. (2021). Enforc-797 ing analytic constraints in neural networks emulating physical systems. Physi-798 cal Review Letters, 126, 98302. doi: 10.1103/PhysRevLett.126.098302 799 Bocquet, M., & Sakov, P. (2013). Joint state and parameter estimation with an iter-800 ative ensemble Kalman smoother. Nonlinear Processes in Geophysics, 20, 803-801 818. doi: 10.5194/npg-20-803-2013 802 Brajard, J., Carrassi, A., Bocquet, M., & Bertino, L. (2021).Combining data 803 assimilation and machine learning to infer unresolved scale parametrization. 804 Philosophical Transactions of the Royal Society A: Mathematical, Physical and 805 Engineering Sciences, 379, 20200086. doi: 10.1098/rsta.2020.0086 806 Brenowitz, N. D., Beucler, T., Pritchard, M., & Bretherton, C. S. (2020). Interpret-807 ing and stabilizing machine-learning parametrizations of convection. Journal of 808 the Atmospheric Sciences, 77, 4357-4375. doi: 10.1175/JAS-D-20-0082.1 809 Brenowitz, N. D., & Bretherton, C. S. (2018). Prognostic validation of a neural net-810 work unified physics parameterization. Geophysical Research Letters, 45, 6289-811 6298. doi: 10.1029/2018GL078510 812 Bretherton, C. S., Henn, B., Kwa, A., Brenowitz, N. D., Watt-Meyer, O., McGib-813 bon, J., ... Harris, L. (2022).Correcting coarse-grid weather and cli-814 mate models by machine learning from global storm-resolving simulations. 815 Journal of Advances in Modeling Earth Systems, 14, e2021MS002794. doi: 816 10.1029/2021MS002794 817 Brynjarsdóttir, J., & O'Hagan, A. (2014). Learning about physical parameters: the 818 importance of model discrepancy. Inverse Problems, 30, 114007. doi: 10.1088/ 819 0266-5611/30/11/114007 820 Burgers, G., Jan van Leeuwen, P., & Evensen, G. (1998). Analysis scheme in the en-821 semble Kalman filter. Monthly Weather Review, 126, 1719–1724. doi: 10.1175/ 822 1520-0493(1998)126(1719:ASITEK)2.0.CO;2 823 Chada, N. K., Stuart, A. M., & Tong, X. T. (2020). Tikhonov regularization within 824 ensemble Kalman inversion. SIAM Journal on Numerical Analysis, 58, 1263-825 1294. doi: 10.1137/19M1242331 826 Chen, Y., & Oliver, D. S. (2012). Ensemble randomized maximum likelihood method 827 as an iterative ensemble smoother. Mathematical Geosciences, 44, 1–26. doi: 828 10.1007/s11004-011-9376-z 829 Cleary, E., Garbuno-Inigo, A., Lan, S., Schneider, T., & Stuart, A. M. (2021). Cal-830 ibrate, emulate, sample. Journal of Computational Physics, 424, 109716. doi: 831 10.1016/j.jcp.2020.109716 832

833	Cohen, Y., Lopez-Gomez, I., Jaruga, A., He, J., Kaul, C. M., & Schneider, T.
834	(2020). Unified entrainment and detrainment closures for extended eddy-
835	diffusivity mass-flux schemes. Journal of Advances in Modeling Earth Systems,
836	12, e2020MS002162. doi: 10.1029/2020MS002162
837	Cohn, S. E. (1997). An introduction to estimation theory. Journal of the Meteorolog-
838	ical Society of Japan. Ser. II, 75, 257-288. doi: 10.2151/jmsj1965.75.1B_257
839	Cotter, S. L., Roberts, G. O., Stuart, A. M., & White, D. (2013). MCMC methods
840	for functions: Modifying old algorithms to make them faster. Statistical Sci-
841	ence. 28, 424-446, doi: 10.1214/13-STS421
842	de Roode, S. R., Duvnkerke, P. G., & Siebesma, A. P. (2000). Analogies be-
8/3	tween mass-flux and Beynolds-averaged equations
844	<i>snheric Sciences</i> , 57, 1585-1598, doi: 10.1175/1520-0469(2000)057(1585:
845	ABMFAR\2.0.CO:2
846	Dunbar, O., Howland, M. F., Schneider, T., & Stuart, A. (2022). Ensemble-
847	based experimental design for targeted high-resolution simulations to in-
848	form climate models. Earth and Space Science Open Archive. 24. doi:
849	10.1002/essoar.10510142.1
850	Emerick, A. A., & Reynolds, A. C. (2013). Ensemble smoother with multiple data
851	assimilation. Computers ℓ Geosciences, 55, 3–15, doi: 10.1016/i.cageo.2012.03
852	.011
853	Evensen G (1994) Sequential data assimilation with a nonlinear quasi-geostrophic
854	model using Monte Carlo methods to forecast error statistics. <i>Journal of Geo</i>
855	physical Research: Oceans. 99, 10143–10162, doi: 10.1029/94JC00572
856	Garbuno-Inigo A Hoffmann F Li W & Stuart A M (2020) Interact-
857	ing Langevin diffusions: Gradient structure and ensemble Kalman sam-
858	pler. SIAM Journal on Applied Dynamical Systems, 19, 412-441. doi:
859	10.1137/19M1251655
860	Guillaumin, A. P., & Zanna, L. (2021). Stochastic-deep learning parameterization of
861	ocean momentum forcing Journal of Advances in Modeling Earth Systems 13
862	e2021MS002534. doi: 10.1029/2021MS002534
863	Hansen, P. C. (1990). Truncated singular value decomposition solutions to discrete
864	ill-posed problems with ill-determined numerical rank. SIAM Journal on Sci-
865	entific and Statistical Computing, 11, 503-518. doi: 10.1137/0911028
866	Hansen, P. C. (1998). Rank-Deficient and Discrete Ill-Posed Problems. Society for
867	Industrial and Applied Mathematics. doi: 10.1137/1.9780898719697
868	He, J., Cohen, Y., Lopez-Gomez, L., Jaruga, A., & Schneider, T. (2021). An im-
869	proved perturbation pressure closure for eddy-diffusivity mass-flux schemes.
870	Earth and Space Science Open Archive, 37. doi: 10.1002/essoar.10505084.2
871	Houtekamer, P. L. & Mitchell, H. L. (1998). Data assimilation using an ensem-
872	ble Kalman filter technique. Monthly Weather Review. 126, 796–811. doi: 10
873	.1175/1520-0493(1998)126(0796:DAUAEK)2.0.CO:2
874	Houtekamer, P. L., & Mitchell, H. L. (2001). A sequential ensemble Kalman fil-
875	ter for atmospheric data assimilation. Monthly Weather Review, 129, 123–137.
876	doi: 10.1175/1520-0493(2001)129(0123:ASEKFF)2.0.CO:2
877	Huang, D. Z., Huang, J., Reich, S., & Stuart, A. M. (2022). Efficient derivative-free
878	Bayesian inference for large-scale inverse problems. doi: 10.48550/arxiv.2204
879	.04386
880	Huang, D. Z., Schneider, T., & Stuart, A. M. (2022). Iterated Kalman methodology
881	for inverse problems. Journal of Computational Physics, 463, 111262. doi: 10
882	.1016/i.jcp.2022.111262
883	Iglesias, M. A. (2016). A regularizing iterative ensemble Kalman method for PDF-
884	constrained inverse problems. Inverse Problems 32, 025002, doi: 10.1088/0266
885	-5611/32/2/025002
886	Iglesias, M. A., Law, K. J. H., & Stuart, A. M. (2013). Ensemble Kalman methods
887	for inverse problems. Inverse Problems, 29, 045001. doi: 10.1088/0266-5611/

888	29/4/045001
889	Kaipio, J., & Somersalo, E. (2006). Statistical and computational inverse problems
890	(Vol. 160). Springer Science & Business Media.
891	Kennedy, M., & O'Hagan, A. (2001). Bayesian calibration of computer models.
892	Journal of the Royal Statistical Society Series B, 63, 425-464. doi: 10.1111/
893	1467-9868.00294
894	Keskar, N. S., Mudigere, D., Nocedal, J., Smelyanskiy, M., & Tang, P. T. P. (2016).
895	On large-batch training for deep learning: Generalization gap and sharp min-
896	<i>ima.</i> doi: 10.48550/arXiv.1609.04836
897	Kovachki, N. B., & Stuart, A. M. (2019). Ensemble Kalman inversion: a derivative-
898	free technique for machine learning tasks. Inverse Problems, 35, 095005. doi:
899	10.1088/1361-6420/ab1c3a
900	Lee, Y. (2021). Sampling error correction in ensemble Kalman inversion. doi: 10
901	.48550/arxiv.2105.11341
902	Li, M., Zhang, T., Chen, Y., & Smola, A. J. (2014). Efficient mini-batch training for
903	stochastic optimization. ACM. doi: 10.1145/2623330.2623612
904	Li, Z., Kovachki, N., Azizzadenesheli, K., Liu, B., Bhattacharya, K., Stuart, A.,
905	& Anandkumar, A. (2020). Fourier neural operator for parametric partial
906	differential equations. doi: 10.48550/arxiv.2010.08895
907	Ling, J., Kurzawski, A., & Templeton, J. (2016). Reynolds averaged turbulence
908	modelling using deep neural networks with embedded invariance. Journal of
909	Fluid Mechanics, 807, 155-166. doi: 10.1017/jfm.2016.615
910	Lopez-Gomez, I., Cohen, Y., He, J., Jaruga, A., & Schneider, T. (2020). A gener-
911	alized mixing length closure for eddy-diffusivity mass-flux schemes of turbu-
912	lence and convection. Journal of Advances in Modeling Earth Systems, 12,
913	$e^{2020MS002161}$. doi: $10.1029/2020MS002161$
914	Lopez-Gomez, I., McGovern, A., Agrawal, S., & Hickey, J. (2022). Global extreme
915	heat forecasting using neural weather models. doi: $10.48550/arxiv.2205.10972$
916	Lorenz, E. N. (1975). Climatic predictability. In The physical basis of climate and
917	climate modelling (Vol. 16, p. 132–136). World Meteorological Organization.
918	Metropolis, N., Rosenbluth, A. W., Rosenbluth, M. N., Teller, A. H., & Teller, E.
919	(1953). Equation of state calculations by fast computing machines. The
920	Journal of Chemical Physics, 21, 1087-1092. doi: 10.1063/1.1699114
921	Meyer, D., Hogan, R. J., Dueben, P. D., & Mason, S. L. (2022). Machine learning
922	emulation of 3D cloud radiative effects. Journal of Advances in Modeling Earth
923	Systems, 14, e2021MS002550. doi: 10.1029/2021MS002550
924	Mlawer, E. J., Taubman, S. J., Brown, P. D., Iacono, M. J., & Clough, S. A. (1997).
925	Radiative transfer for inhomogeneous atmospheres: RRTM, a validated
926	correlated-k model for the longwave. Journal of Geophysical Research: At-
927	<i>mospheres</i> , <i>102</i> , 16663–16682. doi: 10.1029/97JD00237
928	Moral, P. D., Doucet, A., & Jasra, A. (2006). Sequential Monte Carlo samplers.
929	Journal of the Royal Statistical Society. Series B (Statistical Methodology), 68,
930	411-430.
931	Morzfeld, M., Adams, J., Lunderman, S., & Orozco, R. (2018). Feature-based data
932	assimilation in geophysics. Nonlinear Processes in Geophysics, 25, 355-374.
933	doi: $10.5194/npg-25-355-2018$
934	Munoz-Sabater, J., Dutra, E., Agusti-Panareda, A., Albergei, C., Arduini, G., Bal-
935	samo, G., Inepaut, JN. (2021). ERAD-Land: a state-of-the-art global
936	12 12 12 12 12 12 12 12 12 12 12 12 12 1
937	4045-4000. doi: 10.0154/c550-10-4045-2021 National Academics of Sciences Engineering and Medicine (2012) Therising an
938	Our Changing Planet: A Decadal Strategy for Farth Observation from Space
939	Washington DC: The National Academics Pross. doi: 10.17226/24029
940	National Network Δ M (2021) The random feature model for input output
941	maps botwoon Banach spaces SIAM Lowmal on Scientific Commuting 19
942	maps between Danach spaces. DIAM Journal on Detentific Computing, 43,

943	A3212–A3243. doi: 10.1137/20M133957X
944	Pathak, J., Subramanian, S., Harrington, P., Raja, S., Chattopadhyay, A., Mar-
945	dani, M., Anandkumar, A. (2022). FourCastNet: A global data-driven
946	high-resolution weather model using adaptive fourier neural operators. doi:
947	10.48550/arxiv.2202.11214
948	Prechelt, L. (1998). Early stopping - but when? Springer Berlin Heidelberg. doi: 10
949	.1007/3-540-49430-8_3
950	Pressel, K. G., Kaul, C. M., Schneider, T., Tan, Z., & Mishra, S. (2015). Large-
951	eddy simulation in an anelastic framework with closed water and entropy
952	balances. Journal of Advances in Modeling Earth Systems, 7, 1425-1456. doi:
953	10.1002/2015 MS000496
954	Rasp, S., Pritchard, M. S., & Gentine, P. (2018). Deep learning to represent subgrid
955	processes in climate models. Proceedings of the National Academy of Sciences,
956	115, 9684-9689. doi: 10.1073/pnas.1810286115
957	Rasp, S., & Thuerey, N. (2021). Data-driven medium-range weather prediction with
958	a Resnet pretrained on climate simulations: A new model for WeatherBench.
959	Journal of Advances in Modeling Earth Systems, 13, e2020MS002405. doi:
960	10.1029/2020 MS002405
961	Ravuri, S., Lenc, K., Willson, M., Kangin, D., Lam, R., Mirowski, P., Mohamed,
962	S. (2021). Skilful precipitation nowcasting using deep generative models of
963	radar. Nature, 597, 672-677. doi: 10.1038/s41586-021-03854-z
964	Reichel, L., & Rodriguez, G. (2013). Old and new parameter choice rules for discrete
965	ill-posed problems. Numerical Algorithms, 63, 65-87. doi: 10.1007/s11075-012
966	-9612-8
967	Reichstein, M., Camps-Valls, G., Stevens, B., Jung, M., Denzler, J., Carvalhais, N.,
968	& Prabhat. (2019). Deep learning and process understanding for data-driven
969	Earth system science. <i>Nature</i> , 566, 195–204. doi: 10.1038/s41586-019-0912-1
970	Schillings, C., & Stuart, A. M. (2017). Analysis of the ensemble Kalman filter for in-
971	verse problems. SIAM Journal on Numerical Analysis, 55, 1264-1290. doi: 10
972	.1137/16M105959X
973	Schmit, T. J., Griffith, P., Gunshor, M. M., Daniels, J. M., Goodman, S. J.,
974	& Lebair, W. J. (2017). A closer look at the ABI on the GOES-R se-
975	ries. Bulletin of the American Meteorological Society, 98, 681-698. doi:
976	10.1175/BAMS-D-15-00230.1
977	Schneider, T., Dunbar, O. R. A., Wu, J., Böttcher, L., Burov, D., Garbuno-Iñigo,
978	A., Shaman, J. (2021). Epidemic management and control through risk-
979	dependent individual contact interventions.
980	doi: 10.48550/arxiv.2109.10970
981	Schneider, T., & Griffies, S. M. (1999). A conceptual framework for predictabil-
982	ity studies. Journal of Climate, 12, 3133-3155. doi: 10.1175/1520-0442(1999)
983	012(3133:ACFFPS)2.0.CO;2
984	Schneider, T., Lan, S., Stuart, A., & Teixeira, J. (2017). Earth system model-
985	ing 2.0: A blueprint for models that learn from observations and targeted
986	high-resolution simulations. Geophysical Research Letters, 44, 312-396. doi: 10.1000/0015CL05C101
987	10.1002/2017GL076101
988	Schneider, T., Stuart, A. M., & Wu, JL. (2020). Ensemble Kalman inversion for
989	sparse learning of dynamical systems from time-averaged data. doi: 10.48550/
990	arxiv.2007.06175
991	Schneider, T., Stuart, A. M., & Wu, JL. (2021). Learning stochastic closures using
992	ensemble Kalman inversion. Transactions of Mathematics and Its Applications,
993	5, tnab003. doi: 10.1093/imatrm/tnab003
994	Seitert, A., & Rasp, S. (2020). Potential and limitations of machine learning for
995	modeling warm-rain cloud microphysical processes. Journal of Advances in
996	Modeling Earth Systems, 12, e2020MS002301. doi: 10.1029/2020MS002301

997	Shen, Z., Sridhar, A., Tan, Z., Jaruga, A., & Schneider, T. (2022). A library of
998	large-eddy simulations forced by global climate models. Journal of Advances in
999	Modeling Earth Systems, e2021MS002631. doi: 10.1029/2021MS002631
1000	Siebesma, A. P. (1996). On the mass flux approach for atmospheric convection. In
1001	Workshop on new insights and approaches to convective parametrization, 4-7
1002	november 1996 (p. 25-57). Shinfield Park, Reading: ECMWF.
1003	Siebesma, A. P., Bretherton, C. S., Brown, A., Chlond, A., Cuxart, J., Duynkerke,
1004	P. G., others (2003). A large eddy simulation intercomparison study of
1005	shallow cumulus convection. Journal of the Atmospheric Sciences, 60, 1201–
1006	1219. doi: $10.1175/1520-0469(2003)60(1201:ALESIS)2.0.CO;2$
1007	Souza, A. N., Wagner, G. L., Ramadhan, A., Allen, B., Churavy, V., Schloss, J.,
1008	Ferrari, R. (2020). Uncertainty quantification of ocean parameteriza-
1009	tions: Application to the K-profile-parameterization for penetrative convection.
1010	Journal of Advances in Modeling Earth Systems, 12, e2020MS002108. doi:
1011	10.1029/2020 MS002108
1012	Stevens, B., Moeng, CH., Ackerman, A. S., Bretherton, C. S., Chlond, A., de
1013	Roode, S., Zhu, P. (2005). Evaluation of large-eddy simulations via ob-
1014	servations of nocturnal marine stratocumulus. Monthly Weather Review, 133,
1015	1443-1462. doi: 10.1175/MWR2930.1
1016	Suselj, K., Posselt, D., Smalley, M., Lebsock, M. D., & Teixeira, J. (2020). A new
1017	methodology for observation-based parameterization development. Monthly
1018	Weather Review, 148, 4159–4184. doi: 10.1175/MWR-D-20-0114.1
1019	Sønderby, C. K., Espeholt, L., Heek, J., Dehghani, M., Oliver, A., Salimans, T.,
1020	Kalchbrenner, N. (2020). Metnet: A neural weather model for precipitation
1021	forecasting. doi: 10.48550/arXiv.2003.12140
1022	Tan, Z., Kaul, C. M., Pressel, K. G., Cohen, Y., Schneider, T., & Teixeira, J. (2018).
1023	An extended eddy-diffusivity mass-flux scheme for unified representation of
1024	subgrid-scale turbulence and convection. Journal of Advances in Modeling
1025	Earth Systems, 10, 770-800. doi: 10.1002/2017MS001162
1026	Tarantola, A. (2005). Inverse Problem Theory and Methods for Model Parameter
1027	<i>Estimation.</i> Society for Industrial and Applied Mathematics. doi: 10.1137/1
1028	.9780898717921
1029	Tong, X. T., & Morzfeld, M. (2022). Localization in ensemble Kalman inversion.
1030	doi: 10.48550/arXiv.2201.10821
1031	van Leeuwen, P. J. (2015). Representation errors and retrievals in linear and non-
1032	linear data assimilation. Quarterly Journal of the Royal Meteorological Society,
1033	141, 1612-1623. doi: 10.1002/qj.2464
1034	Villefranque, N., Blanco, S., Couvreux, F., Fournier, R., Gautrais, J., Hogan, R. J.,
1035	Williamson, D. (2021). Process-based climate model development har-
1036	nessing machine learning: III. The representation of cumulus geometry and
1037	their 3D radiative effects. Journal of Advances in Modeling Earth Systems, 13,
1038	e2020MS002423. doi: 10.1029/2020MS002423
1039	Weyn, J. A., Durran, D. R., Caruana, R., & Cresswell-Clay, N. (2021). Sub-seasonal
1040	forecasting with a large ensemble of deep-learning weather prediction models.
1041	Journal of Advances in Modeling Earth Systems, 13, e2021MS002502. doi:
1042	10.1029/2021MS002502
1043	Xiao, H., Wu, JL., Wang, JX., Sun, R., & Roy, C. J. (2016). Quantifying and
1044	reducing model-form uncertainties in Reynolds-averaged Navier–Stokes sim-
1045	ulations: A data-driven, physics-informed Bayesian approach. Journal of
1046	Computational Physics, 324, 115-136. doi: 10.1016/j.jcp.2016.07.038
1047	Lama, L., & Bolton, I. (2020). Data-driven equation discovery of ocean mesoscale
1048	ciosures. Geophysical Research Letters, 47, 62020GL088376. doi: 10.1029/
1049	ZUZUGLU00010 Zhao W. I. Contino D. Doichstein M. Zhang V. Zhao G. Way, V. O'
1050	Linao, w. L., Gentine, P., Reichstein, M., Zhang, Y., Zhou, S., Wen, Y., Qiu,
1051	G. I. (2019). Physics-constrained machine learning of evapotranspiration.

Geophysical Research Letters, 46, 14496-14507. doi: 10.1029/2019GL085291

1052