

# Calibration and Uncertainty Quantification of Convective Parameters in an Idealized GCM

Oliver R. A. Dunbar<sup>1</sup>, Alfredo Garbuno-Inigo<sup>2</sup>, Tapio Schneider<sup>1</sup>,  
Andrew M. Stuart<sup>1</sup>

<sup>1</sup>California Institute of Technology, Pasadena, California, USA

<sup>2</sup>Instituto Tecnológico Autónomo de México, Ciudad de México, México.

## Key Points:

- We use time averaged climate statistics to calibrate convective parameters and quantify their uncertainties.
- We demonstrate use of the calibrate-emulate-sample algorithm to provide efficient calibration and uncertainty quantification.
- Parametric uncertainty in climate predictions is quantified by sampling from the learnt convective parameter distribution.

---

Corresponding author: Oliver Dunbar, [odunbar@caltech.edu](mailto:odunbar@caltech.edu)

**Abstract**

Parameters in climate models are usually calibrated manually, exploiting only small subsets of the available data. This precludes both optimal calibration and quantification of uncertainties. Traditional Bayesian calibration methods that allow uncertainty quantification are too expensive for climate models; they are also not robust in the presence of internal climate variability. For example, Markov chain Monte Carlo (MCMC) methods typically require  $O(10^5)$  model runs and are sensitive to internal variability noise, rendering them infeasible for climate models. Here we demonstrate an approach to model calibration and uncertainty quantification that requires only  $O(10^2)$  model runs and can accommodate internal climate variability. The approach consists of three stages: (i) a calibration stage uses variants of ensemble Kalman inversion to calibrate a model by minimizing mismatches between model and data statistics; (ii) an emulation stage emulates the parameter-to-data map with Gaussian processes (GP), using the model runs in the calibration stage for training; (iii) a sampling stage approximates the Bayesian posterior distributions by sampling the GP emulator with MCMC. We demonstrate the feasibility and computational efficiency of this calibrate-emulate-sample (CES) approach in a perfect-model setting. Using an idealized general circulation model, we estimate parameters in a simple convection scheme from synthetic data generated with the model. The CES approach generates probability distributions of the parameters that are good approximations of the Bayesian posteriors, at a fraction of the computational cost usually required to obtain them. Sampling from this approximate posterior allows the generation of climate predictions with quantified parametric uncertainties.

**Plain Language Summary**

Calibrating climate models with available data and quantifying their uncertainties are essential to make climate predictions accurate and actionable. A primary source of uncertainties in climate models comes from representation of small-scale processes such as moist convection. Parameters in convection schemes and other parameterizations are usually calibrated by hand, using only a small fraction of data that are available. As a result, the calibration process may miss information about the small-scale processes in question. This paper presents a proof-of-concept, in an idealized setting, of how parameters in climate models can be calibrated using a substantial fraction of the available data, and how uncertainties in the parameters can be quantified. We employ a new algorithm, called calibrate-emulate-sample (CES), which makes such calibration and uncertainty quantification feasible for computationally expensive climate models. CES reduces the hundreds of thousands of model runs usually required to quantify uncertainties in computer models to hundreds, thereby achieving about a factor 1000 speedup. It leads to more robust calibration and uncertainty quantification in the presence of noise arising from chaotic variability of the climate system. We show how uncertainties in climate model parameters can be translated into quantified uncertainties of climate predictions through ensemble integrations.

**1 Introduction**

The principal uncertainties in climate predictions arise from the representation of unresolvable yet important small-scale processes, such as those controlling cloud cover (Cess et al., 1989, 1990; Bony & Dufresne, 2005; Stephens, 2005; Bony et al., 2006; Vial et al., 2013; Webb et al., 2013; Brient & Schneider, 2016; Schneider, Teixeira, et al., 2017). These processes are represented by parameterization schemes, which relate unresolved quantities such as cloud statistics to variables resolved on the climate models' computational grid, such as temperature and humidity. The parameterization schemes depend on parameters that are a priori unknown, and so fixing the parameters is associated with uncertainty. The process of fixing these parameters to values that are most consistent

with data is known as calibration, which generally requires solving an optimization problem. Traditionally, however, parameters are calibrated (“tuned”) by hand, in a process that exploits only a small subset of the available observational data and relies on the knowledge and intuition of climate modelers about plausible ranges of parameters and their effect on the simulated climate of a model (Randall & Wielicki, 1997; Mauritsen et al., 2012; Golaz et al., 2013; Hourdin et al., 2013; Flato et al., 2013; Hourdin et al., 2017; Schmidt et al., 2017; Zhao et al., 2018). More recently, some broader-scale automated approaches that more systematically quantify the plausible range of parameters have begun to be explored (Couvreur et al., 2020; Hourdin et al., 2020).

Opportunities to improve climate models lie in exploiting a larger fraction of the available observational data together with high-resolution simulations, and learning from both systematically and not manually (Schneider, Lan, et al., 2017). To fully account for parametric uncertainty, we adopt a Bayesian view of the model-data relationship, which amounts to solving a Bayesian inverse problem. Given model, data, and prior information on the parameters, Bayesian inversion yields a posterior distribution of parameters. The mean or mode of the posterior distribution define the best parameter estimate, and the entire distribution provides uncertainty quantification through its spread about the mean or mode. We use the mismatch between climate statistics simulated with the model and those obtained from observations or high-resolution simulations as the data likelihood in the Bayesian inverse problem to calibrate parameterizations in a climate model and to quantify their parametric uncertainties. Our focus is on learning from time-averaged climate statistics for three reasons: (1) time-averaged statistics are what is relevant for climate predictions; (2) time-averaged statistics vary more smoothly in space than atmospheric states, leading to a smoother optimization problem than that of atmospheric state estimation in numerical weather prediction (NWP); (3) time-averaging over long time-intervals reduces the effect of the unknown initial state of the system, removing the need to determine it. Focusing on time-averaged climate statistics, rather than on instantaneous states or trajectories as in NWP, makes it possible to exploit climate observations and high-resolution simulations even when their native resolutions are very different from those of climate models.

While learning from climate statistics accumulated in time presents opportunities, it also comes with challenges. Accumulating statistics in time is computationally much more expensive than the forecasts over hours or days used in NWP. Therefore, we need algorithms for learning from data that minimize the number of climate model runs required. Traditional methods for Bayesian calibration and uncertainty quantification such as Markov chain Monte Carlo (MCMC) typically require many iterations—often more than  $10^5$ —to reach statistical convergence (e.g., Geyer, 2011). Conducting so many computationally expensive climate model runs is not feasible, rendering MCMC impractical for climate model calibration (Annan & Hargreaves, 2007). Additionally, while MCMC can be used to obtain the distribution of model parameters given data, it is not robust with respect to noise in the evaluation of the map from model parameters to data. Such noise, arising from natural variability in the chaotic climate system, can lead to trapping of the Markov chains in spurious, noise-induced local maxima of the likelihood function (Cleary et al., 2021). This presents additional challenges to using MCMC methods for climate model calibration.

Here we demonstrate a new approach to climate model uncertainty quantification that overcomes the limitations of traditional Bayesian calibration methods in a relatively simple proof-of-concept. The approach—called calibrate-emulate-sample (CES) (Cleary et al., 2021)—consists of three successive stages, which each exploit proven concepts and methods:

1. In a calibration stage, we use variants of ensemble Kalman methods. These approaches were originally developed as a derivative-free method for state estima-

tion (Evensen, 1994; Van Leeuwen & Evensen, 1996) and are now widely used in NWP (Houtekamer & Zhang, 2016). The methods were subsequently developed for simultaneous state and parameter estimation, and variants were developed to deal with strongly nonlinear systems (Bocquet & Sakov, 2012; Gu & Oliver, 2007; Li & Reynolds, 2009; Sakov et al., 2012; Bocquet & Sakov, 2014). They were eventually recognized as a general purpose tool for the solution of inverse problems whose objective is parameter estimation (Chen & Oliver, 2012; Emerick & Reynolds, 2013; Reich, 2011; Evensen, 2018). Here we use an optimization approach, referred to as ensemble Kalman inversion (Iglesias et al., 2013), which builds on the work of (Chen & Oliver, 2012; Emerick & Reynolds, 2013). However, ensemble Kalman inversion and other ensemble Kalman methods do not provide a basis for systematic uncertainty quantification, except in linear Gaussian problems (Annan & Hargreaves, 2007; Gland et al., 2009; Ernst et al., 2015).

2. In an emulation stage, we train an emulator on the climate model statistics generated during the calibration stage in order to quantify uncertainties. To emulate how the climate model statistics depend on parameters to be calibrated, we use Gaussian processes (GPs), a hierarchical machine learning method that learns smooth functions from a set of noisy training points (Kennedy & O’Hagan, 2001; Santner et al., 2018). The GP approach also learns the statistics of the uncertainty in the resulting predicted function. The training points here are provided by the climate model runs performed in the calibration stage.
3. In a sampling stage, we approximate the posterior distribution on parameters given model and data, using the GP emulator to replace the parameter-to-climate statistics map. The emulator is used to estimate error statistics and to sample from the approximate posterior distribution with MCMC. Because the GP emulator is computationally cheap to evaluate and is smooth by virtue of the smoothing properties of GPs, this avoids the issues that limit the usability of MCMC for sampling from climate models directly.

The CES approach is described in detail in Cleary et al. (2021), which provides a justification and contextualization of the approach in the literature on data assimilation and Bayesian calibration. If uncertainty quantification is not required, the calibration step provides an effective, derivative-free parameter estimation method, and the emulation and sampling stages are not required. However, when uncertainty quantification is required, the role of the calibration stage is to provide a good set of training points in the vicinity of the posterior mean or mode, to train the emulator in the next stage of the algorithm.

The purpose of this paper is to demonstrate the feasibility of the approach for estimating parameters in an idealized general circulation model (GCM). This represents a proof-of-concept in a small parameter space and limited data space; how the methods scale up to larger problems will be discussed at the end.

This paper is arranged as follows: Section 2 describes the experimental setup, including the idealized GCM and the generation of synthetic data from it. Section 3 describes the CES approach and the methods used in each stage. Section 4 describes the results of numerical experiments that use CES to calibrate parameters in the idealized GCM and quantify their uncertainties. It also demonstrates how sampling from the posterior distribution of parameters can be used to generate climate predictions with quantified uncertainties. Section 5 discusses and summarizes the results and their applicability to larger problems.

## 2 Experimental Setup

### 2.1 General Circulation Model

We use the idealized GCM described by Frierson et al. (2006) and O’Gorman and Schneider (2008b), which is based on the spectral dynamical core of the Flexible Modeling System developed at the Geophysical Fluid Dynamics Laboratory. To approximate the solution of the hydrostatic primitive equations, it uses the spectral transform method in the horizontal, with spectral resolution T21 and with 32 latitude points on the transform grid. It uses finite differences with 10 unevenly spaced sigma levels in the vertical. We chose this relatively coarse resolution to keep our numerical experiments computationally efficient, so that comparison of CES with much more expensive methods is feasible. The lower boundary of the GCM is a homogeneous slab ocean (1 m mixed-layer thickness). Radiative transfer is represented by a semi-gray, two-stream radiative transfer scheme, in which the optical depth of longwave and shortwave absorbers is a prescribed function of latitude and pressure (O’Gorman & Schneider, 2008b), irrespective of the concentration of water vapor in the atmosphere (i.e., without an explicit representation of water vapor feedback). Insolation is constant and approximates Earth’s annual mean insolation at the top of the atmosphere.

We focus our calibration and uncertainty quantification experiments on parameters in the GCM’s convection scheme, which is a quasi-equilibrium moist convection scheme that can be viewed as a simplified version of the Betts-Miller convection scheme (Betts, 1986; Betts & Miller, 1986, 1993). It relaxes temperature  $T$  and specific humidity  $q$  toward reference profiles on a timescale  $\tau$  (Frierson, 2007):

$$\frac{\partial T}{\partial t} + \dots = -f_T \frac{T - T_{\text{ref}}}{\tau} \quad (1)$$

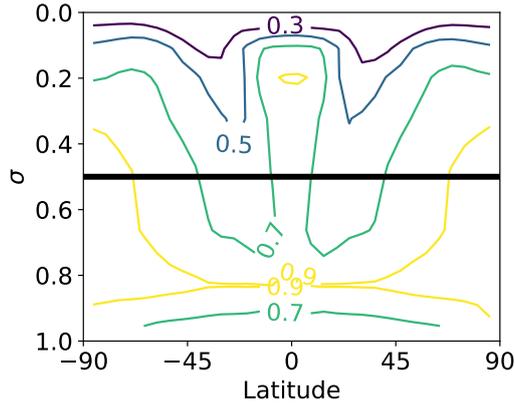
and

$$\frac{\partial q}{\partial t} + \dots = -f_T f_q \frac{q - q_{\text{ref}}}{\tau}. \quad (2)$$

Here,  $f_T(z; T, q, p)$  is a function of altitude  $z$  and of the thermodynamic state of an atmospheric column (dependent on temperature  $T$ , pressure  $p$ , and specific humidity  $q$  in the column), which determines where and when the convection scheme is active;  $f_q(T, q, p)$  is a function that modulates the relaxation of the specific humidity in non-precipitating (shallow) convection (Frierson, 2007; O’Gorman & Schneider, 2008b). The reference temperature profile is a moist adiabat,  $T_{\text{ma}}(z)$ , shifted by a state-dependent and constant-with-height offset  $\Delta T$ , which is chosen to ensure conservation of enthalpy integrated over a column:  $T_{\text{ref}}(z) = T_{\text{ma}}(z) + \Delta T$ . The reference specific humidity  $q_{\text{ref}}(z)$  is the specific humidity corresponding to a fixed relative humidity RH relative to the moist adiabat  $T_{\text{ma}}(z)$ . The two key parameters in this simple convection scheme thus are the timescale  $\tau$  and the relative humidity RH. We demonstrate how we can learn about them from synthetic data generated with the GCM.

### 2.2 Variable Selection and Generation of Synthetic Data

The idealized GCM with the simple quasi-equilibrium convection scheme has been used in numerous studies of large-scale atmosphere dynamics and mechanisms of climate changes, especially those involving the hydrologic cycle (e.g., O’Gorman & Schneider, 2008b, 2008a; Bordoni & Schneider, 2008; O’Gorman & Schneider, 2009b; Schneider et al., 2010; Merlis & Schneider, 2011; O’Gorman, 2011; Kaspi & Schneider, 2011, 2013; Levine & Schneider, 2015; Bischoff & Schneider, 2014; Wills et al., 2017; Wei & Bordoni, 2018). We know from this body of work that the convection scheme primarily affects the atmospheric thermal stratification in the tropics, with weaker effects in the extratropics (Schneider & O’Gorman, 2008). We also know that the relative humidity parameter (RH) in the moist convection scheme controls the humidity of the tropical free troposphere but likewise has a weaker effect on the humidity of the extratropical free troposphere (O’Gorman



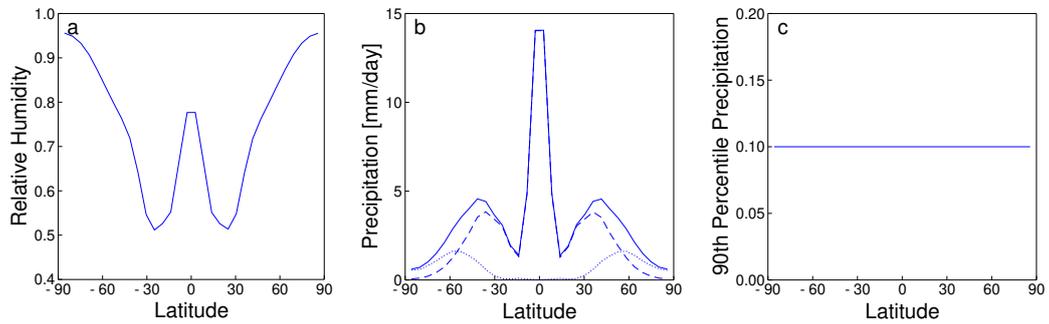
**Figure 1.** Zonal average of relative humidity averaged over one month. The black line shows the level at which data was extracted for computing data misfit functions.

205 et al., 2011). Thus, we expect tropical circulation statistics to be especially informative  
 206 about the parameters in the convection scheme. However, convection plays a central role  
 207 in extreme precipitation events at all latitudes (O’Gorman & Schneider, 2009b, 2009a),  
 208 so we expect statistics of precipitation extremes to be informative about convective pa-  
 209 rameters, and in particular to contain information about the relaxation timescale  $\tau$ .

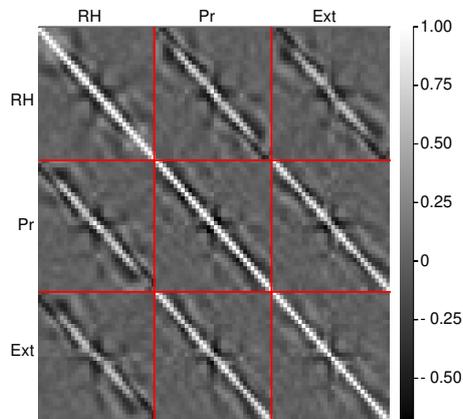
210 As climate statistics from which we want to learn about the convective parame-  
 211 ters, we choose 30-day averages of the free-tropospheric relative humidity, of the precip-  
 212 itation rate, and of a measure of the frequency of extreme precipitation. Because the GCM  
 213 is statistically zonally symmetric, we take zonal averages in addition to the time aver-  
 214 ages. The relative humidity is evaluated at  $\sigma = 0.5$  (where  $\sigma = p/p_s$  is pressure  $p$  nor-  
 215 malized by the local surface pressure  $p_s$ ), as shown in Figure 1. As a measure of the fre-  
 216 quency of precipitation extremes, we use the probability that daily precipitation rates  
 217 exceed a high, latitude-dependent threshold. The threshold is chosen as the latitude-dependent  
 218 90th percentile of daily precipitation in a long (18000 days) control simulation of the GCM  
 219 in a statistically steady state. So for the parameters in the control simulation, the precip-  
 220 itation threshold is expected to be exceeded 10% of the time at each latitude. The  
 221 convective parameters in the control simulation are fixed at their reference values  $\text{RH} =$   
 222  $0.7$  and  $\tau = 2$  h (O’Gorman & Schneider, 2008b), and we collect the parameters in the  
 223 vector  $\theta^\dagger = (\theta_{\text{RH}}^\dagger, \theta_\tau^\dagger) = (0.7, 2 \text{ h})$ . Figure 2 shows the mean relative humidity, the mean  
 224 precipitation rate (broken down into its contributions coming from the convection scheme  
 225 and from condensation at resolved scales), and the 90th percentile precipitation rate, from  
 226 the control simulation averaged over 600 batches of 30-day windows. We use the single  
 227 long control simulations of duration 18000 days only for the creation of Figure 2 and for  
 228 the estimation of noise covariances.

### 229 2.3 Definition of noise covariance

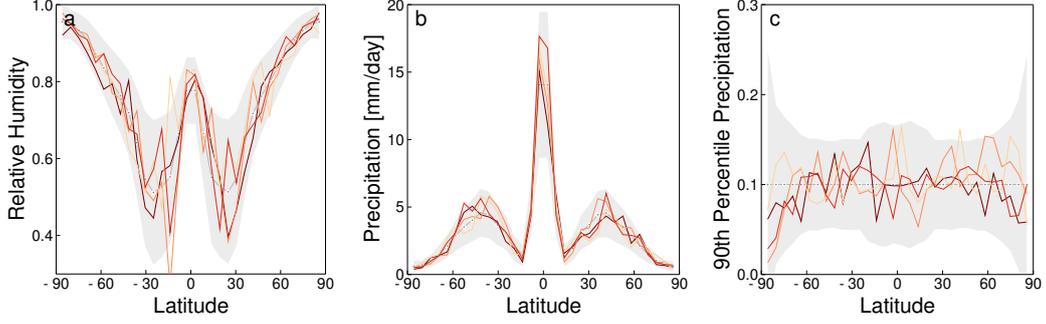
230 Estimation of model parameters requires specification of a noise covariance matrix,  
 231 reflecting errors and uncertainties in the data. The principal source of noise in our perfect-  
 232 model setting with synthetic data is sampling variability due to finite-time averaging with  
 233 unknown initial conditions. The initial condition is forgotten at sufficiently long times  
 234 because of the chaotic nature of atmospheric variability, so a central limit theorem quan-  
 235 tifies the finite-time fluctuations around infinite-time averages that are caused by uncer-  
 236 tain initial conditions. Therefore, the asymptotic distribution of the fluctuations is a mul-  
 237 tivariate normal distribution  $N(0, \Sigma(\theta))$  with zero mean and covariance matrix  $\Sigma(\theta)$ . We



**Figure 2.** Long-term mean values of synthetic data. (a) Free-tropospheric relative humidity. (b) Total daily precipitation rate (solid) and its contributions from convection (dashed) and grid-scale condensation (dotted). (c) Probability of daily precipitation exceeding a 90th percentile (which is trivially 10% in this case).



**Figure 3.** Correlation matrix associated with the internal variability estimated from the control simulation, plotted to illustrate the structure of  $\Sigma$ . The matrix blocks labeled (RH, Pr, Ext) are associated with the observed relative humidity, precipitation, and extreme precipitation.



**Figure 4.** Four noisy realizations of the synthetic data, plotted in color over the underlying mean (grey circles) and 95% confidence intervals from  $\Gamma(\theta^\dagger)$  (grey bars). (a) Relative humidity. (b) Daily precipitation rate. (c) Probability of daily precipitation exceeding the 90th percentile of the long-term mean data.

238 estimate the covariance matrix at  $\Sigma(\theta^\dagger)$ , that is, with the parameters  $\theta^\dagger$  in the control  
 239 simulation. To estimate  $\Sigma(\theta^\dagger)$ , we run the GCM for 600 windows of length 30 days (be-  
 240 cause we use 30-day averages to estimate parameters) and calculate the sample covari-  
 241 ance matrix of the 30-day means. With the 3 latitude-dependent fields evaluated at 32  
 242 latitude points,  $\Sigma(\theta^\dagger)$  is a  $96 \times 96$  symmetric matrix representing noise from internal  
 243 variability in finite-time averages. Hereafter, we make the assumption that  $\Sigma(\theta) \approx \Sigma(\theta^\dagger)$   
 244 for any  $\theta$ , and thus we treat  $\Sigma$  as a constant matrix. In practical implementations of this  
 245 method, a corresponding constant  $\Sigma$  can be estimated from climatology. We illustrate  
 246 the correlation structure of  $\Sigma$  in Figure 3.

To generate synthetic data, we also include the effect of measurement error (Kennedy  
 & O’Hagan, 2001). We add Gaussian noise to the time-averaged model statistics, with  
 a diagonal covariance structure in data space. We construct the measurement error cov-  
 covariance matrix  $\Delta$  to be diagonal with entries  $\delta_i > 0$ , where  $i$  indexes over data type  
 (the 3 observed quantities) and latitude (32 locations). Combining this measurement cov-  
 covariance matrix  $\Delta$  with the covariance matrix  $\Sigma$  arising from internal variability leads  
 to an inflated noise covariance matrix

$$\Gamma = \Sigma + \text{diag}(\delta_i^2) = \Sigma + \Delta. \quad (3)$$

There are many options to pick  $\delta_i$ . We choose it by reducing the distance of the 95% con-  
 fidence interval to its nearest physical boundary for each  $i$  by a constant factor  $C$ , so as  
 to retain physical properties (e.g., precipitation must be nonnegative). Denote the mean  
 $\mu_i$ , variance  $\Sigma_{ii}$ , and a physical boundary set  $\partial\Omega_i$  for each data  $i$ , we choose

$$\delta_i = C \min \left( \text{dist}(\mu_i + 2\sqrt{\Sigma_{ii}}, \partial\Omega_i), \text{dist}(\mu_i - 2\sqrt{\Sigma_{ii}}, \partial\Omega_i) \right).$$

247 We take  $C = 0.2$ . This value implies a significant noise inflation, with an average ratio  
 248 of the standard deviations  $\sqrt{\Gamma_{ii}}/\sqrt{\Sigma_{ii}}$  of 2.3. Figure 4 shows the resulting data mean  
 249 (grey circles), the 95% confidence interval of the inflated covariance (grey ribbon), and  
 250 four realizations of the data  $\mathbf{y}^{(1)}, \dots, \mathbf{y}^{(4)}$  (yellow to red lines), each defined by taking  
 251 a different 30-day average of the GCM, and adding a different realization of  $N(0, \Delta)$ . These  
 252 four realizations will be used throughout when presenting our results.

### 3 Methods

#### 3.1 Misfit functions for time averaged data

Both calibration and uncertainty quantification in CES rely on a misfit function (standardized error) that quantifies mismatch between model output and data. Calibration minimizes a (possibly regularized) misfit function over the parameter space; uncertainty quantification samples from the posterior distribution, using a misfit function as the negative log-likelihood. To define the desired misfit function, we introduce  $\mathcal{G}_T(\boldsymbol{\theta}; \mathbf{z}^{(0)})$  and  $\mathcal{G}_\infty(\boldsymbol{\theta})$ , which denote the mapping from the parameter vector  $\boldsymbol{\theta}$  to the 96 data points, either averaged over a finite time horizon ( $T$ ) or over an infinite time horizon ( $\infty$ ). The former average depends on the unknown initial condition  $\mathbf{z}^{(0)}$ , whereas the latter does not, because the initial condition is forgotten after a sufficiently long time. We refer to  $\mathcal{G}_T(\boldsymbol{\theta}; \mathbf{z}^{(0)})$  as the forward model and  $\mathcal{G}_\infty(\boldsymbol{\theta})$  as the infinite time-horizon forward model.

To define the misfit function, we begin from the relationship between parameters  $\boldsymbol{\theta}$  and data  $\mathbf{y}$ . Expressed in terms of finite-time averages, this relationship has the form

$$\mathbf{y} = \mathcal{G}_T(\boldsymbol{\theta}; \mathbf{z}^{(0)}) + \boldsymbol{\eta}, \quad \boldsymbol{\eta} \sim N(0, \Delta). \quad (4)$$

The data realizations  $\mathbf{y}^{(i)}$  for  $i = 1, \dots, 4$  in Section 2.3 can therefore be seen as evaluations of (4) for initial conditions  $\mathbf{z}_i^{(0)}$ , noise realization  $\boldsymbol{\eta}_i \sim N(0, \Delta)$ , and at a parameter  $\boldsymbol{\theta}^\dagger$ . This form has the undesirable feature that it involves  $\mathbf{z}^{(0)}$ , a quantity that is not of intrinsic interest. However, the central limit theorem asserts that

$$\mathcal{G}_T(\boldsymbol{\theta}; \mathbf{z}^{(0)}) \approx \mathcal{G}_\infty(\boldsymbol{\theta}) + \sigma,$$

where  $\sigma \sim N(0, \Sigma)$ . This theorem quantifies the forgetting of the initial condition after sufficiently long times, e.g., for the atmosphere,  $T \gtrsim 15$  days (Zhang et al., 2019). From the definition (3) of the total noise covariance matrix,  $\Gamma = \Sigma + \Delta$ , we can combine the observational and internal-variability noise and write

$$\mathbf{y} = \mathcal{G}_\infty(\boldsymbol{\theta}) + \boldsymbol{\gamma}, \quad \boldsymbol{\gamma} \sim N(0, \Gamma). \quad (5)$$

We can equally interpret  $\mathbf{y}^{(i)}$  for  $i = 1, \dots, 4$  from Section 2.3 as evaluations of (5) with noise  $\boldsymbol{\gamma}_i \sim N(0, \Gamma)$  and at parameter  $\boldsymbol{\theta}^\dagger$ . This removes the dependence on initial condition but is expressed in terms of infinite-time averages. Computing these averages directly is not feasible, but (in the emulation phase) we introduce a procedure that enables us to learn a surrogate model for  $\mathcal{G}_\infty$ , using carefully chosen finite-time model evaluations  $\mathcal{G}_T$ .

In the Bayesian approach to parameter learning, the aim is to determine the conditional distribution of parameters  $\boldsymbol{\theta}$  given a realization of data  $\mathbf{y}$  (written  $\boldsymbol{\theta} \mid \mathbf{y}$ ) by applying Bayes theorem, which states  $\mathbb{P}(\boldsymbol{\theta} \mid \mathbf{y})$  is proportional to the product of the likelihood  $\mathbb{P}(\mathbf{y} \mid \boldsymbol{\theta})$  and the prior  $\mathbb{P}(\boldsymbol{\theta})$ . In the case where a surrogate model for  $\mathcal{G}_\infty$  is available,  $\mathbf{y} \mid \boldsymbol{\theta}$  is defined as the pushforward of  $\boldsymbol{\theta}$  through a parameter-to-data map (5), and we define the misfit function, the negative logarithm of the likelihood, as

$$\Phi(\boldsymbol{\theta}, \mathbf{y}) = \frac{1}{2} \|\mathbf{y} - \mathcal{G}_\infty(\boldsymbol{\theta})\|_\Gamma^2, \quad (6)$$

where  $\|\cdot\|_\Gamma = \|\Gamma^{-1/2} \cdot\|_2$  is the Mahalanobis distance. Before a surrogate model for  $\mathcal{G}_\infty$  is available, this function is infeasible to evaluate, but as  $\mathbf{y} \mid \boldsymbol{\theta}$  can be defined using (4) as well, we consider the related misfit function

$$\Phi_T(\boldsymbol{\theta}, \mathbf{y}; \mathbf{z}^{(0)}) = \frac{1}{2} \|\mathbf{y} - \mathcal{G}_T(\boldsymbol{\theta}; \mathbf{z}^{(0)})\|_{\Gamma+\Sigma}^2. \quad (7)$$

We view evaluation of  $\mathcal{G}_T$  from any initial condition as a random approximation of  $\mathcal{G}_\infty$ ; hence, the additional internal-variability covariance matrix  $\Sigma$  appears in (7).

273 The CES algorithm in this context proceeds as follows: use optimization based on  
 274 (7) to calibrate parameters; on the basis of evaluations of  $\mathcal{G}_T$  made during the calibra-  
 275 tion, learn a GP surrogate for  $\mathcal{G}_\infty$ ; then use this surrogate to sample from the posterior  
 276 distribution of  $(\boldsymbol{\theta} \mid \mathbf{y})$  defined using (6). We will henceforth neglect  $\mathbf{z}^{(0)}$  in our nota-  
 277 tion, and just write  $\mathcal{G}_T(\boldsymbol{\theta})$  and  $\Phi_T(\boldsymbol{\theta}, \mathbf{y})$ . Viewing the initial condition as random makes  
 278 these objects random as well.

279 We have the following undesirable properties of the finite-time model average  $\mathcal{G}_T(\boldsymbol{\theta})$ :  
 280 (i) it is computationally expensive to evaluate for large  $T$ ; (ii) it can be nondifferentiable  
 281 or difficult to differentiate (e.g., because of non-differentiability of parameterization schemes  
 282 in climate models); and (iii) evaluations of it are not deterministic (when one drops the  
 283 explicit dependence on initial conditions). Our methodology, detailed in the upcoming  
 284 sections, is constructed to overcome these difficulties.

285 An alternative approach to emulating the map  $\mathcal{G}_\infty(\boldsymbol{\theta})$  is to emulate  $\Phi(\boldsymbol{\theta}, \mathbf{y})$  directly.  
 286 This is often more computationally efficient as one always models a scalar function. Dur-  
 287 ing investigation however, we found the efficiency comes at a cost of reduced accuracy  
 288 and interpretability; in preliminary experiments these drawbacks proved detrimental to  
 289 performance, and so we abandoned the approach.

### 290 3.2 Prior distributions

291 The priors of the physical parameters are taken to be the logit-normal and lognor-  
 292 mal distributions,  $\theta_{RH} \sim \text{Logitnormal}(0, 1)$  and  $\theta_\tau \sim \text{Lognormal}(12 \text{ h}, (12 \text{ h})^2)$ , for  
 293 the relative humidity and timescale parameter, respectively. Define the invertible trans-  
 294 formation  $\mathcal{T} : \mathbb{R}^2 \rightarrow \mathbb{R}^2$  by  $\mathcal{T}(x, y) = (\text{logit}(x), \log(y))$ . Through an abuse of nota-  
 295 tion, we relabel the transformed (or computational) parameters as  $\boldsymbol{\theta}$ , and the untrans-  
 296 formed (or physical) parameters (relative humidity and timescale) are uniquely defined  
 297 by  $\mathcal{T}^{-1}(\boldsymbol{\theta})$ . The forward map  $\mathcal{G}_T$  is defined as a composition  $\mathcal{F}_T \circ \mathcal{T}^{-1}$  where  $\mathcal{F}_T$  is the  
 298 forward map defined on the physical parameters.

299 This choice allows us to apply our methods in the transformed space, where our  
 300 priors are normally distributed and unbounded (namely  $\text{logit}(\theta_{RH}) \sim N(0, 1)$  and  $\log(\theta_\tau) \sim$   
 301  $N(10.17, 1)$ ); meanwhile the climate model is applied only to physically defined variables,  
 302  $\theta_{RH} \in [0, 1]$  and  $\theta_\tau \in [0, \infty)$ . In this way the prior distributions enforce physical con-  
 303 straints on the parameters.

### 304 3.3 Calibrate: Ensemble Kalman Inversion

305 Ensemble Kalman inversion (EKI) (Iglesias et al., 2013) is an offline variant of en-  
 306 semble Kalman filtering designed to learn parameters in a general model, rather than  
 307 states of a dynamical system. EKI can be viewed as a derivative-free optimization al-  
 308 gorithm. Given a set of data  $\mathbf{y}$ , it iteratively evolves an ensemble of parameter estimates  
 309 both so that they achieve consensus and evolve toward the optimal parameter value  $\boldsymbol{\theta}^*$   
 310 (close to  $\boldsymbol{\theta}^\dagger$  if the quality/volume of the data is good enough) that minimizes an objec-  
 311 tive function, in our case given by the misfit (7), possibly with inclusion of a regulariza-  
 312 tion term. It has great potential for use with chaotic or stochastic models due to its ensemble-  
 313 based, derivative-free approach for optimizing parameters. Theoretical work shows that  
 314 noisy continuous-time versions of EKI exhibit an averaging effect that skips over fluc-  
 315 tuations superimposed onto the ergodic averaged forward model (Duncan et al., 2021).  
 316 Empirical results suggest similar phenomena apply to EKI as implemented here, justifi-  
 317 fying the application to noisy forward model evaluations. Furthermore, the derivative-  
 318 free approach scales well to high-dimensional parameter spaces, as evidenced by the use  
 319 of Kalman filtering in numerical weather prediction, where billions of parameters char-  
 320 acterizing atmospheric states are routinely estimated (Kalnay, 2002). This makes the  
 321 algorithm appealing for complex climate models. The algorithm is mathematically proven

322 to find the optimizer, within an initial, ensemble-dependent subspace, for linear mod-  
 323 els (Schillings & Stuart, 2017a). It is known to be effective for high-dimensional nonlin-  
 324 ear models (Iglesias et al., 2013; Schneider et al., 2020c, 2020a), such as the nonlinear  
 325 map from parameters to data represented by the idealized GCM we use in our proof-of-  
 326 concept here.

The EKI algorithm we use is detailed in Iglesias et al. (2013). The algorithm it-  
 eratively updates an ensemble of parameters,  $\boldsymbol{\theta}_m^{(n)}$ , where  $m = 1, \dots, M$  denotes an en-  
 semble member, and the superscript  $n$  denotes the iteration count. The algorithm uses  
 the ensemble to update parameters according to the following equation

$$\boldsymbol{\theta}_m^{(n+1)} = \boldsymbol{\theta}_m^{(n)} + C_{\boldsymbol{\theta}\boldsymbol{G}}^{(n)} \left( \Gamma + C_{\boldsymbol{G}\boldsymbol{G}}^{(n)} \right)^{-1} \left( \boldsymbol{y} - \mathcal{G}_T(\boldsymbol{\theta}_m^{(n)}) \right),$$

327 where  $C_{\boldsymbol{G}\boldsymbol{G}}$  is the empirical covariance of the ensemble of quantities of interest from model  
 328 runs, and  $C_{\boldsymbol{\theta}\boldsymbol{G}}$  is the empirical cross-covariance of the ensemble of parameters and the  
 329 ensemble of quantities of interest. The covariance matrix of the distribution of differences  
 330 between realizations of  $\boldsymbol{y}$  and  $\mathcal{G}_T(\cdot)$  is  $\Gamma + \Sigma$ , which is approximated empirically in the  
 331 update by  $\Gamma + C_{\boldsymbol{G}\boldsymbol{G}}^{(n)}$ . When ensemble methods are used as approximate samplers (Chen  
 332 & Oliver, 2012; Emerick & Reynolds, 2013), additional independent noise is added to  
 333  $\boldsymbol{y}$  at each iteration and for every ensemble member; however, because we are solving an  
 334 optimization problem within this calibration phase, such noise is not added here.

335 We initialize the algorithm by drawing an ensemble of size  $M = 100$  by sampling  
 336 the parameter space from assumed prior distributions on the parameters.

### 337 3.4 Emulate: Gaussian Process Emulators (EKI-GP)

During the calibration stage with  $N$  iterations and an ensemble of size  $M$ , we ob-  
 tain a collection of input–output pairs

$$\{\boldsymbol{\theta}_m^{(n)}, \mathcal{G}_T(\boldsymbol{\theta}_m^{(n)})\}, \quad n = 0, \dots, N, \quad m = 1, \dots, M.$$

338 The cloud of points  $\{\boldsymbol{\theta}_m^{(n)}\}$  from the calibration stage (a) spans the prior distribution,  
 339 as initial EKI draws are from the prior, and (b) in later iterations, has a high density  
 340 around the point  $\boldsymbol{\theta}^*$  to which EKI converges. We use regression to train a GP emula-  
 341 tor mapping  $\boldsymbol{\theta}$  to  $\mathcal{G}_T(\boldsymbol{\theta})$ , using the input–output pairs  $\{\boldsymbol{\theta}_m^{(n)}, \mathcal{G}_T(\boldsymbol{\theta}_m^{(n)})\}$ , which are referred  
 342 to as training points in the context of GP regression. The emulation will be most accu-  
 343 rate in regions with more training points, that is, around  $\boldsymbol{\theta}^*$ . This is typically near the  
 344 true solution  $\boldsymbol{\theta}^\dagger$ , and it is the region where the posterior parameter distribution will have  
 345 high probability; this is precisely where uncertainty quantification requires accuracy. In  
 346 effect, EKI serves as an effective algorithm for selecting good training points for GP re-  
 347 gression.

Gaussian processes emulate the statistics of the input–output pairs, using a Gaus-  
 sian assumption. Specifically, we learn an approximation of the form

$$\mathcal{G}_T(\boldsymbol{\theta}) \approx \mathcal{N}(\mathcal{G}_{\text{GP}}(\boldsymbol{\theta}), \Sigma_{\text{GP}}(\boldsymbol{\theta})).$$

348 The approximation is learned from the input–output pairs assuming that the outputs are  
 349 produced from a mean function  $\mathcal{G}_{\text{GP}}(\boldsymbol{\theta})$ , and subject to normally distributed noise de-  
 350 fined by a covariance function  $\Sigma_{\text{GP}}(\boldsymbol{\theta})$ , both dependent on the parameters. The choice  
 351 of notation here is to imply the fact that  $\mathcal{G}_{\text{GP}}(\boldsymbol{\theta})$  serves to approximate the (unattain-  
 352 able) infinite-time average of the model  $\mathcal{G}_\infty(\boldsymbol{\theta})$ , and  $\Sigma_{\text{GP}}(\boldsymbol{\theta})$  serves to approximate the  
 353 covariance matrix  $\Sigma$ . Importantly,  $\Sigma_{\text{GP}}(\boldsymbol{\theta})$  is  $\boldsymbol{\theta}$ -dependent as it also includes the uncer-  
 354 tainty in the approximation of the emulator at  $\boldsymbol{\theta}$  (for example, the emulator uncertainty  
 355  $\Sigma_{\text{GP}}(\boldsymbol{\theta})$  will be large when  $\boldsymbol{\theta}$  is far from the inputs  $\{\boldsymbol{\theta}_m\}$  used in training).

The atmospheric quantities from which we learn about model parameters are cor-  
 related (e.g., relative humidity or daily precipitation at neighboring latitudes are cor-  
 related), resulting in a nondiagonal covariance matrix  $\Sigma$ . Any GP emulator therefore also

requires a nondiagonal covariance  $\Sigma_{\text{GP}}(\boldsymbol{\theta})$ . We can enforce this, by mapping the correlated statistics from the GCM into a decorrelated space by using a principal component analysis on  $\Sigma$ , and then training the GP with the decorrelated statistics to produce an emulator with diagonal covariance  $\tilde{\Sigma}_{\text{GP}}(\boldsymbol{\theta})$ . We use the notation  $(\tilde{\cdot})$  to denote variables in the uncorrelated space. To this end, we first decompose  $\Sigma$  as

$$\Sigma = VD^2V^T.$$

Here,  $V$  is an orthonormal matrix of eigenvectors of the covariance matrix  $\Sigma$ , and  $D$  is the diagonal matrix of the square root of the eigenvalues, or the ordered standard deviations in the basis spanned by the eigenvectors of  $\Sigma$ . We store the outputs from the pairs as columns of a matrix  $Y_{kl} = (\mathcal{G}_T(\boldsymbol{\theta}_l))_k$ , and then we change the basis of this matrix into the uncorrelated coordinates

$$\tilde{Y} = D^{-1}V^TY.$$

356 When trained on  $\tilde{Y}$ , GP returns the mean  $\tilde{\mathcal{G}}_{\text{GP}}(\boldsymbol{\theta})$  and the (diagonal) covariance matrix  
 357  $\tilde{\Sigma}_{\text{GP}}(\boldsymbol{\theta})$ . We use tools from scikit-learn (Pedregosa et al., 2011) to train the emulator.  
 358 After the diagonalization, we can train a scalar-valued GP for each of the 96 output di-  
 359 mensions, rather than having to train processes with vector-valued output. We construct  
 360 a kernel by summing an Automatic Relevance Determination (ARD) radial basis func-  
 361 tion kernel and a white-noise kernel. The ARD kernel is a standard squared exponen-  
 362 tial kernel, where each input dimension has an independent lengthscale hyperparameter.  
 363 This corresponds to regression, rather than interpolation, and the variance of the  
 364 white noise kernel reflects the noise level assumed in the regression. We train by learn-  
 365 ing 4 hyperparameters: the radial basis function variance, a lengthscale for each of the  
 366 two parameters  $\boldsymbol{\theta}$  (due to ARD), and the white-noise variance. We train using the input-  
 367 output pairs of the initial ensemble plus  $N = 5$  subsequent iterations of the EKI algo-  
 368 rithm. We use  $M = 100$  ensemble members; thus, the training requires  $(N+1) \times M =$   
 369 600 short (30-day) runs of our GCM.

370 We continue using the uncorrelated basis in the sampling stage; where required,  
 371 we transform the output of the emulator back into a correlated basis,

$$\begin{aligned} \mathcal{G}_{\text{GP}}(\boldsymbol{\theta}) &= VD\tilde{\mathcal{G}}_{\text{GP}}(\boldsymbol{\theta}), \\ \Sigma_{\text{GP}}(\boldsymbol{\theta}) &= VD\tilde{\Sigma}_{\text{GP}}(\boldsymbol{\theta})DV^T. \end{aligned}$$

### 372 **3.5 Sample: MCMC Sampling with a Gaussian Process Emulator**

373 To quantify uncertainties, we use MCMC to sample the posterior distribution of  
 374 parameters with the GP emulator. The primary reason for using the GP emulator goes  
 375 back to the seminal paper by Sacks et al. (1989) and concerns the fact that it can be eval-  
 376 uated far more quickly than the GCM at a point in parameter space; this is important  
 377 as we require more than  $10^5$  samples within the likelihood  $\mathbb{P}(\mathbf{y} | \boldsymbol{\theta})$  in a typical MCMC  
 378 run to sample the posterior distribution of parameters given data. However the emula-  
 379 tor is also important for two additional reasons: (i) it naturally includes the approxima-  
 380 tion uncertainty (within  $\tilde{\Sigma}_{\text{GP}}$ ) of using an emulator; (ii) it smooths the likelihood func-  
 381 tion because we work with an approximation of (6) based on the smooth  $\mathcal{G}_{\infty}$ , rather than  
 382 (7) based on the noisy  $\mathcal{G}_T$ ; as a result, MCMC is less likely to get stuck in local extrema.

Recall that we trained the GP in uncorrelated coordinates. Within MCMC, one can either map back into the original coordinates or continue working in the uncorrelated space. We choose to continue working in the uncorrelated space, and so we map each data realization  $\mathbf{y}$  into this space:  $\tilde{\mathbf{y}} = D^{-1}V^T\mathbf{y}$ . In the Gaussian likelihood, we can use the GP emulated mean  $\tilde{\mathcal{G}}_{\text{GP}}(\boldsymbol{\theta})$  and covariance matrix  $\tilde{\Sigma}_{\text{GP}}(\boldsymbol{\theta})$  as surrogates for the map  $\mathcal{G}_{\infty}$  and the internal variability covariance matrix  $\Sigma$  (after passing to the un-

correlated coordinates). That is, we approximate the Bayesian posterior distribution as

$$\begin{aligned} \mathbb{P}(\boldsymbol{\theta} \mid \tilde{\mathbf{y}}) &\propto \mathbb{P}(\tilde{\mathbf{y}} \mid \boldsymbol{\theta})\mathbb{P}(\boldsymbol{\theta}) \\ &\propto \frac{1}{\sqrt{\det(\tilde{\Gamma}_{GP}(\boldsymbol{\theta}))}} \exp\left(-\frac{1}{2}\|\tilde{\mathbf{y}} - \tilde{\mathcal{G}}_{GP}(\boldsymbol{\theta})\|_{\tilde{\Gamma}_{GP}(\boldsymbol{\theta})}^2\right)\mathbb{P}(\boldsymbol{\theta}) \\ &\propto \exp\left(-\frac{1}{2}\|\tilde{\mathbf{y}} - \tilde{\mathcal{G}}_{GP}(\boldsymbol{\theta})\|_{\tilde{\Gamma}_{GP}(\boldsymbol{\theta})}^2 - \frac{1}{2}\log\det\tilde{\Gamma}_{GP}(\boldsymbol{\theta})\right)\mathbb{P}(\boldsymbol{\theta}). \end{aligned} \quad (8)$$

Here,  $\tilde{\Gamma}_{GP}(\boldsymbol{\theta}) = \tilde{\Sigma}_{GP}(\boldsymbol{\theta}) + D^{-1}V^T\Delta VD^{-1}$  is the GP approximation of  $\Gamma = \Sigma + \Delta$  in the uncorrelated coordinates. We include the (often overlooked) log-determinant term, arising from the normalization constant due to dependence of  $\Gamma_{GP}$  on  $\boldsymbol{\theta}$ . In the transformed parameter space, our prior  $\mathbb{P}(\boldsymbol{\theta})$  is also Gaussian and therefore can be factored inside this exponential, adding a quadratic penalty to the negative log-likelihood. The MCMC objective function is then defined to be the negative logarithm of the posterior distribution (8); explicitly, given a Gaussian prior  $N(\mathbf{m}, C)$  on the parameters, we define it as

$$\Phi_{\text{MCMC}}(\boldsymbol{\theta}, \tilde{\mathbf{y}}) = \frac{1}{2}\|\tilde{\mathbf{y}} - \tilde{\mathcal{G}}_{GP}(\boldsymbol{\theta})\|_{\tilde{\Gamma}_{GP}(\boldsymbol{\theta})}^2 + \frac{1}{2}\log\det\tilde{\Gamma}_{GP}(\boldsymbol{\theta}) + \frac{1}{2}\|\boldsymbol{\theta} - \mathbf{m}\|_C^2.$$

383 Evaluation of  $\Phi_{\text{MCMC}}$  requires only the mean and covariance from the GP emulator. Fur-  
 384 thermore  $\Phi_{\text{MCMC}}$  is smooth and so is suitable for use within an MCMC algorithm to gen-  
 385 erate samples from the approximate posterior distribution of the parameters. Cleary et  
 386 al. (2021) contains further discussion of MCMC using GPs to emulate the forward model,  
 387 including situations where data comes from finite time-averages but the emulator is de-  
 388 signed to approximate the infinite time-horizon forward model.

389 We use the random walk metropolis algorithm for MCMC sampling. The priors  
 390 chosen were the same, physics-informed priors used to initialize EKI. We choose the pro-  
 391 posal distribution also as a Gaussian with covariance proportional to the prior covari-  
 392 ance. The MCMC run consists of a burn-in of 10,000 samples followed by 190,000 sam-  
 393 ples.

### 394 3.6 Benchmark Gaussian process (B-GP)

395 The performance of any emulator is dependent on the training points. Since we use  
 396 an adaptive procedure (EKI) to concentrate the training points, which is the novel ap-  
 397 proach introduced in Cleary et al. (2021), we also train a benchmark emulator to com-  
 398 pare our results with those resulting from more traditional, brute-force approaches to  
 399 the emulation.

400 For this purpose, we use a GP emulator trained on a uniform set of points. Even  
 401 in two dimensions, it is prohibitively costly for this set to span the support of the prior  
 402 distribution. Instead, we use knowledge of the location and size of the posterior distri-  
 403 bution to place a uniform grid of  $40 \times 40 = 1600$  training points over  $[-1.25, -0.5] \times$   
 404  $[8.0, 10.0]$  in the transformed parameter space. This corresponds to  $[0.62, 0.77] \times [0.83 \text{ h}, 6.12 \text{ h}]$   
 405 in the untransformed parameter space and captures the region of high probability of the  
 406 posterior. The use of posterior knowledge here reduces the number of training points by  
 407 a factor of 20 when compared to spanning 95% of the prior mass; no such posterior in-  
 408 formation is used in the CES algorithm, which automatically places training points where  
 409 they are needed. The benchmark emulator uses the same kernel and training setup as  
 410 in section 3.4, and we use the trained emulator in MCMC experiments in the same way  
 411 as described in Section 3.5. To distinguish the two methods, we refer to the EKI-trained  
 412 GP as EKI-GP and the benchmark (traditionally trained) GP as B-GP.

## 4 Results

To demonstrate the dependence of the parameter uncertainty on the realization of the (inflated) synthetic data, we reproduce the experiments 4 times with each of the four realizations shown in Figure 4. We denote these four sets of data  $\mathbf{y}^{(1)}, \dots, \mathbf{y}^{(4)}$ .

### 4.1 Calibrate: Ensemble Kalman Inversion

We use the first 6 iterations of EKI in the CES training process. The initial ensemble is spread over the whole parameter space but collapses within a few iterations near the true parameter values—to within 10% error in  $\theta_{RH}$  and 30 minutes error in  $\theta_\tau$  (Figure 5). That is, the algorithm evolves toward consensus and toward the true solution. Biases arise from the realization of internal variability and the realization of the observational noise in each  $\mathbf{y}^{(i)}$ .

To check for EKI convergence, we evaluate an additional 4 EKI iterations (labeled 0 to 9). At each iteration  $n$ , we compute residuals of the ensemble mean for each realization of the synthetic data  $\mathbf{y}^{(1)}, \dots, \mathbf{y}^{(4)}$  created at the true parameters  $\boldsymbol{\theta}^\dagger$ ,

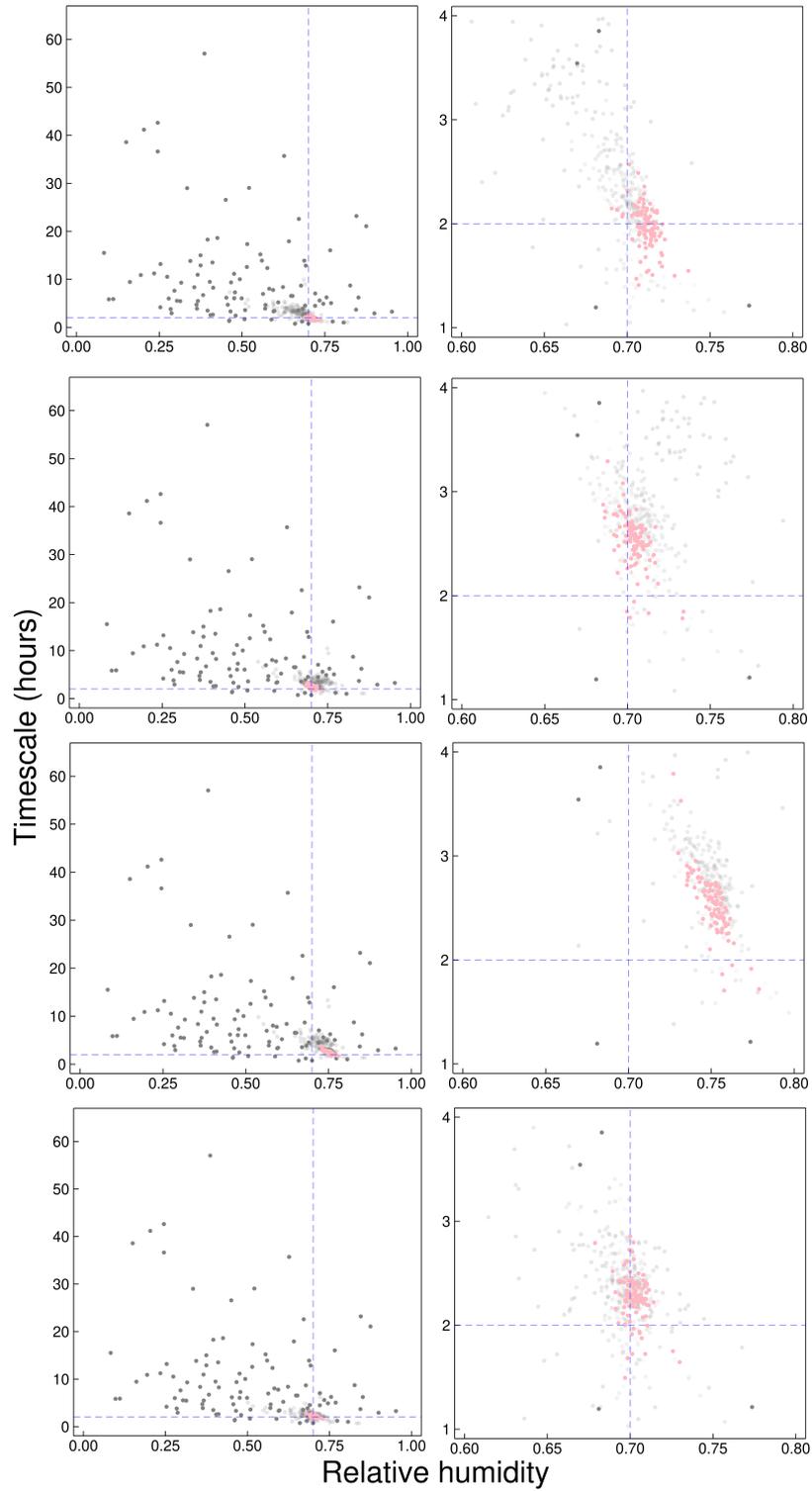
$$\text{Residual}(n; \mathbf{y}^{(i)}) = \left\| \frac{1}{M} \sum_{m=1}^M \mathcal{G}_T(\boldsymbol{\theta}_m^{(n)}) - \mathbf{y}^{(i)} \right\|_{\Gamma}^2,$$

weighting the residuals by the covariance matrix  $\Gamma$  of the synthetic data. Figure 6(a) shows the residual as a function of EKI iteration. The residual decreases quickly over the first few iterations, before plateauing for subsequent iterations. Figure 6(b) shows standard deviations of the ensemble of parameters. The standard deviations decrease monotonically from iteration to iteration, reflecting the evolution toward consensus. The behavior is qualitatively similar for all realizations; quantitative differences reflect different realizations of internal variability in the different data realizations. This behavior reflects the fact that EKI is an optimization method for calibrating parameters: it is not constructed to learn uncertainty but rather to reach consensus around a single parameter value that makes the misfit small (Schillings & Stuart, 2017b).

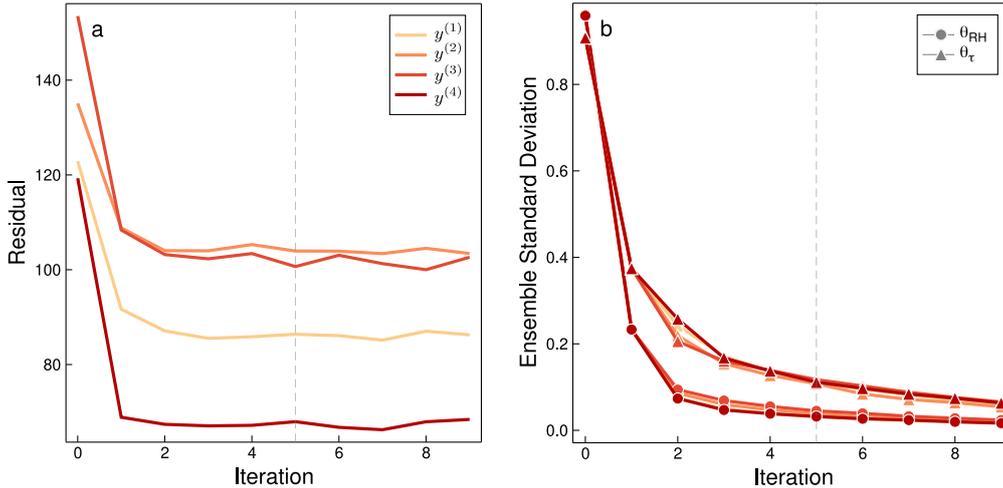
### 4.2 Emulate: Validation

Figure 7 shows the parameter values used for training points for the EKI-GP and B-GP. We use the first 6 EKI iterations (i.e., 600 training points) for training. These are plotted over the associated objective function  $\Phi_{\text{MCMC}}$ . The panels in the left column correspond to the EKI-GP using data vectors  $\mathbf{y}^{(1)}, \dots, \mathbf{y}^{(4)}$ . We see that the resulting objective functions  $\Phi_{\text{MCMC}}$  from training EKI-GP at the marked training points (black dots) lead to unimodal distributions with a minimum near to a significant number of the training points; there are also training points that fall outside of the plotting domain (see Figure 5 for their extent). The right column of Figure 7 shows the benchmark grid for B-GP, which we use as a comparison for the EKI-GP method; the contours of  $\Phi_{\text{MCMC}}$  were calculated using the same realization as their counterpart EKI-GPs. We see that for each realization the EKI-GP and B-GP produce objective functions that are qualitatively similar in terms of the magnitude of the minimum, the location of the minimum, and the approximate shape of the objective function; the quantitative differences are accounted for by differing geometry and density of training points (and hence a difference in approximation uncertainty). In both settings, the objective function  $\Phi_{\text{MCMC}}$  is smooth because the GP smoothly approximates  $\mathcal{G}_\infty$ .

EKI-GP shows similar results for the objective function as B-GP, at a fraction of the computational effort. B-GP is far less practical as a methodology than is EKI-GP because it does not scale well to high-dimensional parameter spaces; it requires many more training points than EKI-GP. Even in our two-dimensional experiments, we needed to use posterior information to reduce the number of training points for B-GP by a fac-



**Figure 5.** EKI ensemble at iterations 0 to 5 displayed as particles in parameter space. Left column: all members; right column: zoom-in near true parameter values. Each row represents optimization with a different data vector  $\mathbf{y}^{(i)}$  from Figure 4. The (initial) prior ensemble 0 is highlighted in dark grey, and the final ensemble 5 is highlighted in pink. The intersection of the dashed blue lines represents the true parameter values used to generate observational data from the GCM.



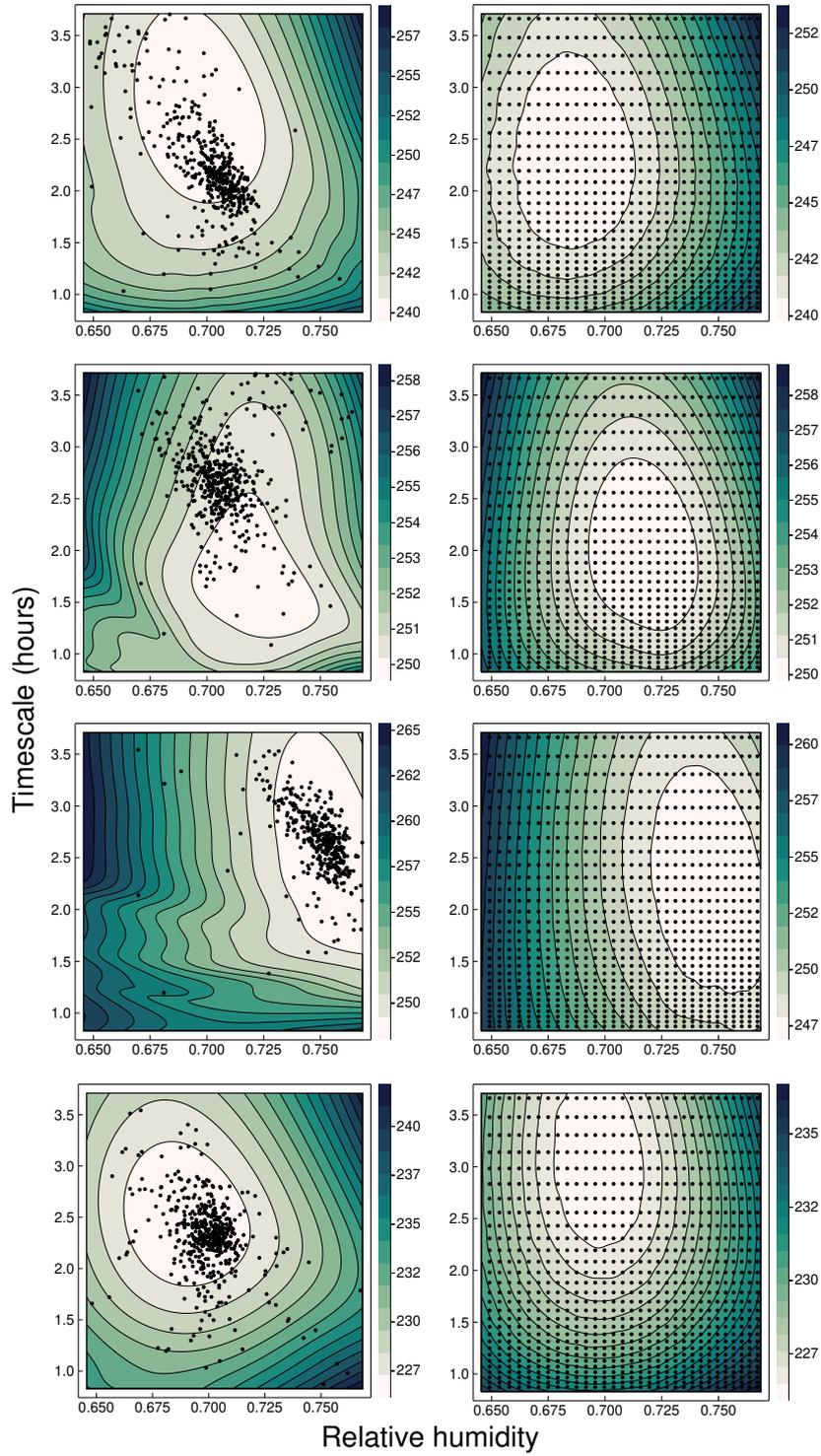
**Figure 6.** Convergence behaviour tests over 9 iterations of EKI for each realization of the data. The vertical dashed line marks the final iteration of Figure 5; we also show behaviour of 4 further iterations. (a) Ensemble-mean residuals relative to synthetic data for each EKI iteration. (b) Standard deviation of ensemble for the relative humidity parameter (circle) and timescale parameter (triangle) for each realization.

	$\sigma_{RH}$	$\sigma_{\tau}$ (hrs)
EKI (Iteration 9)	0.017	0.053
MCMC (EKI-GP)	0.099	0.265
MCMC (B-GP)	0.096	0.359

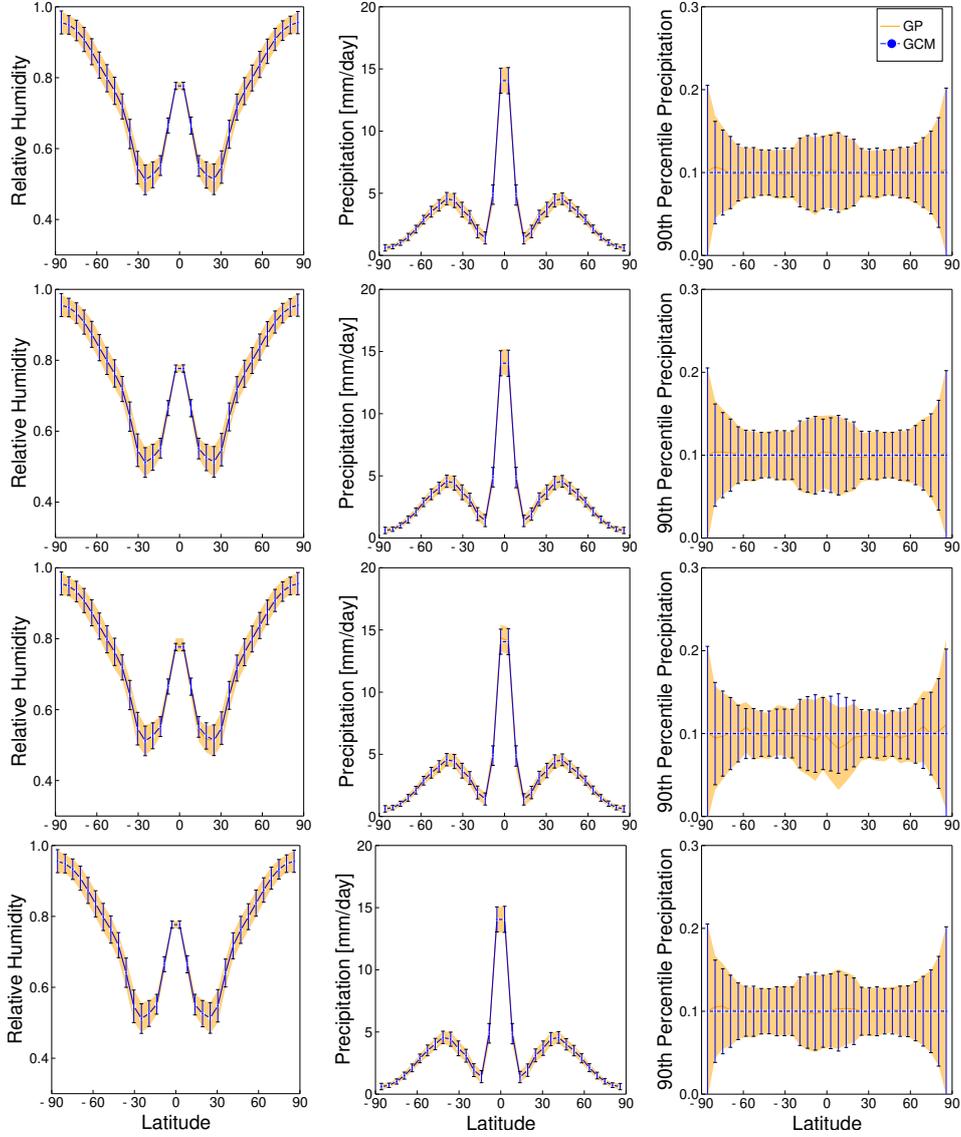
**Table 1.** Average standard deviations of parameters from EKI and MCMC experiments over  $\mathbf{y}^{(1)}, \dots, \mathbf{y}^{(4)}$ .

456 tor of 20. The B-GP comparison is included simply to demonstrate that EKI-GP achieves  
457 comparable results to those achieved by means of traditional emulation.

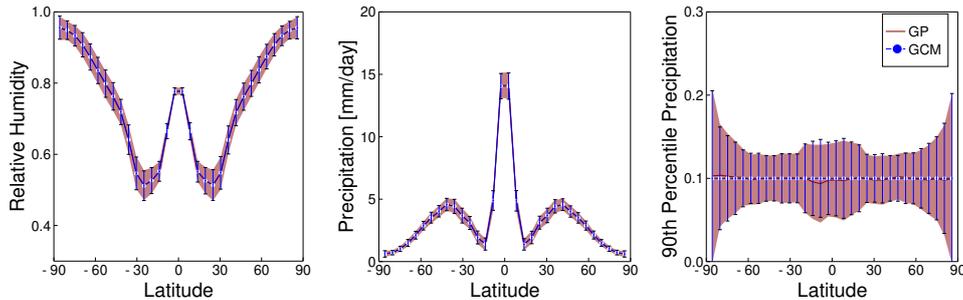
458 We validate the emulator approximation to the data by making a prediction at the  
459 true parameters  $\theta^\dagger$ . We display  $\mathcal{G}_{GP}(\theta^\dagger)$  and the 95% confidence intervals computed  
460 using the variance from  $\Sigma_{GP}(\theta^\dagger)$  in Figure 8 for EKI-GP, and in Figure 9 for B-GP. The  
461 rows of Figure 8 correspond to the EKI-GP results for  $\mathbf{y}^{(1)}, \dots, \mathbf{y}^{(4)}$ . In both figures we  
462 also show the statistics of 600 30-day samples from the control simulation at  $\theta^\dagger$ . Both  
463 the mean and 95% confidence intervals of all EKI-GP emulators (orange line and rib-  
464 bon) closely match the statistics from the GCM runs (blue dots and error bars), as does  
465 the prediction from the B-GP (dark red line and ribbon). The training data are suffi-  
466 cient to ensure that the predicted 95% confidence interval from the emulators do not pro-  
467 duce unphysical values (such as giving negative precipitation rates, or relative humid-  
468 ities outside  $[0, 1]$ ).



**Figure 7.** Training points for the GP emulators (EKI-GP and B-GP), plotted over the objective function  $\Phi_{\text{MCMC}}$  for different data realizations  $\mathbf{y}^{(1)}, \dots, \mathbf{y}^{(4)}$  (rows). Left column: particles representing members of the first 6 EKI iterations. Right column: grid (uniform in the transformed parameters) used to train the benchmark Gaussian process. In both cases, some additional training points fall outside of the plotting domain.



**Figure 8.** Comparison between the GCM statistics at the true parameters  $\theta^\dagger$  and the trained EKI-GP emulator at  $\theta^\dagger$ . The four rows correspond to using EKI with the data vectors  $\mathbf{y}^{(1)}, \dots, \mathbf{y}^{(4)}$ . Blue lines: GCM mean (dots) averaged over 600 30-day runs, with the error bars marking a 95% confidence interval from variances on the diagonal of  $\Gamma$ . Orange: predicted mean (line) and 95% confidence interval (shaded region) produced by the GP emulator.



**Figure 9.** Comparison between the GCM statistics at the true parameters  $\theta^\dagger$  and the trained B-GP emulator predictions at  $\theta^\dagger$ . Blue: GCM mean (dots) averaged over 600 30-day runs, with the error bars marking a 95% confidence interval from variances on the diagonal of  $\Gamma$ . Dark red: predicted mean (line) and 95% confidence interval (shaded region) produced by the B-GP emulator.

469

### 4.3 Sample: MCMC Sampling

470

471

472

473

474

475

476

477

478

479

We use an MCMC algorithm to generate a set of samples from the posterior distribution with the help of the GP emulator. We choose the random walk step size (which multiplies the covariance in the proposal) at the start of a run to achieve proposal acceptance rates near to 25%. (This is near optimal in a precise sense for certain high-dimensional posteriors (Roberts et al., 2004); in practice, it works well beyond this setting.) All sampling is performed in the transformed space where the prior distribution is normal. Figure 10 shows kernel density estimates of the MCMC results; the panels in the left column are for EKI-GP (for  $\mathbf{y}^{(1)}, \dots, \mathbf{y}^{(4)}$ ), and the panels in the right column are for B-GP for the same data realization. We display contours of the posterior that contain 50%, 75%, and 99% of the mass of the posterior density.

480

481

482

483

484

485

486

487

488

489

490

491

492

493

494

All sets of results converge to similar regions of the parameter space about the true parameters, and the spread of uncertainty is quantified similarly in both EKI-GP and B-GP. Table 1 shows the standard deviations of the individual parameters alongside the empirical standard deviation calculated from the ensemble spread in EKI iteration 9. The standard deviations from the MCMC posterior based on EKI-GP and B-GP are similar. We re-emphasize that the EKI is constructed as an ensemble optimization method and has the property that the ensemble evolves towards consensus among the parameters, while also matching the data: the ensemble collapses. As a result, the EKI ensemble spread is an inadequate estimate the uncertainty, as seen in Table 1. As explained in the introduction, ensemble methods are only justifiable to quantify uncertainties in the Gaussian posterior setting (Chen & Oliver, 2012; Emerick & Reynolds, 2013). Our approach is justifiable whenever the GP accurately approximates the forward model (Cleary et al., 2021). The use of EKI for the design of training points for the GP does not require accurate uncertainty quantification within EKI; it only relies on EKI approximately locating minimizers of the model-data misfit.

495

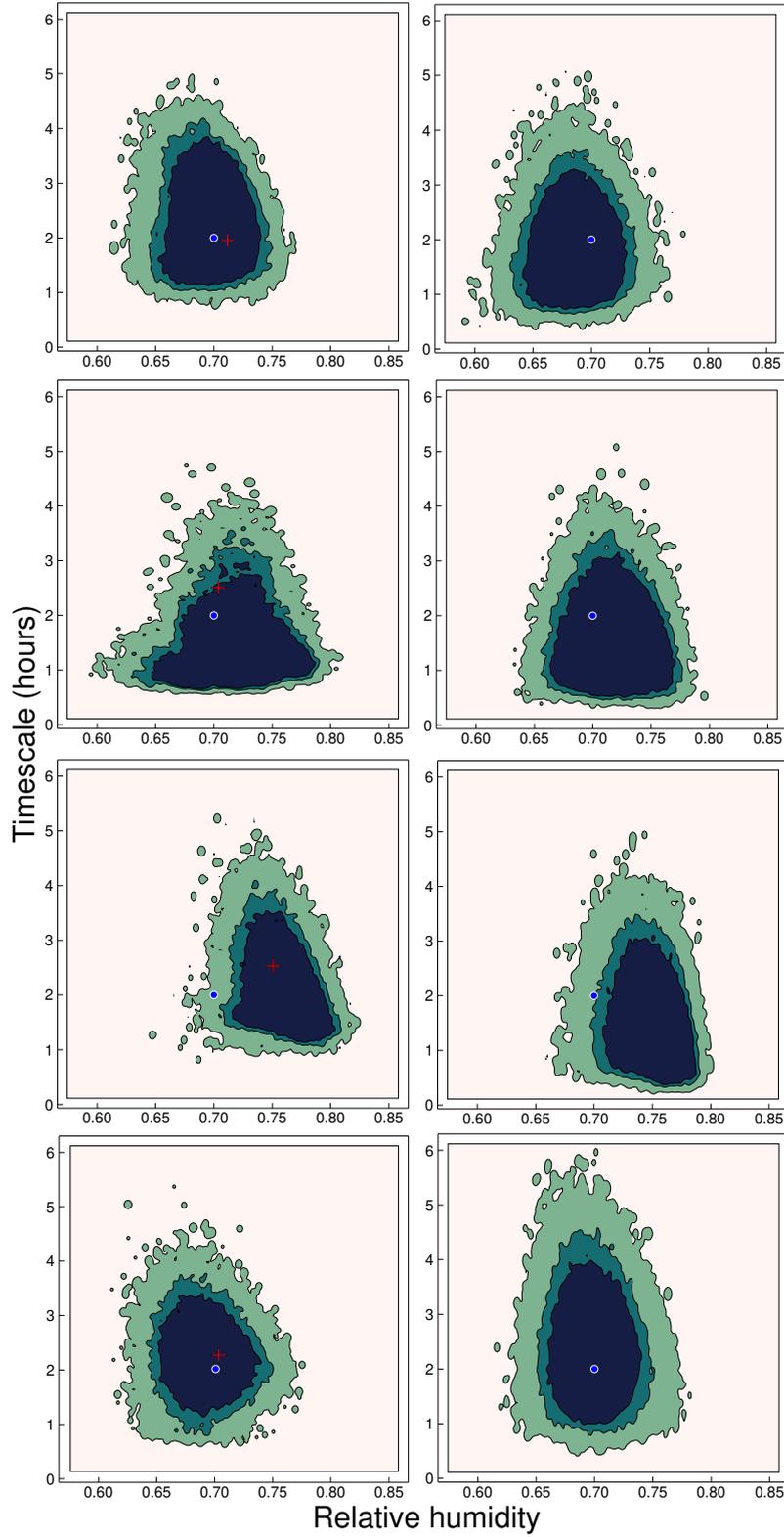
496

497

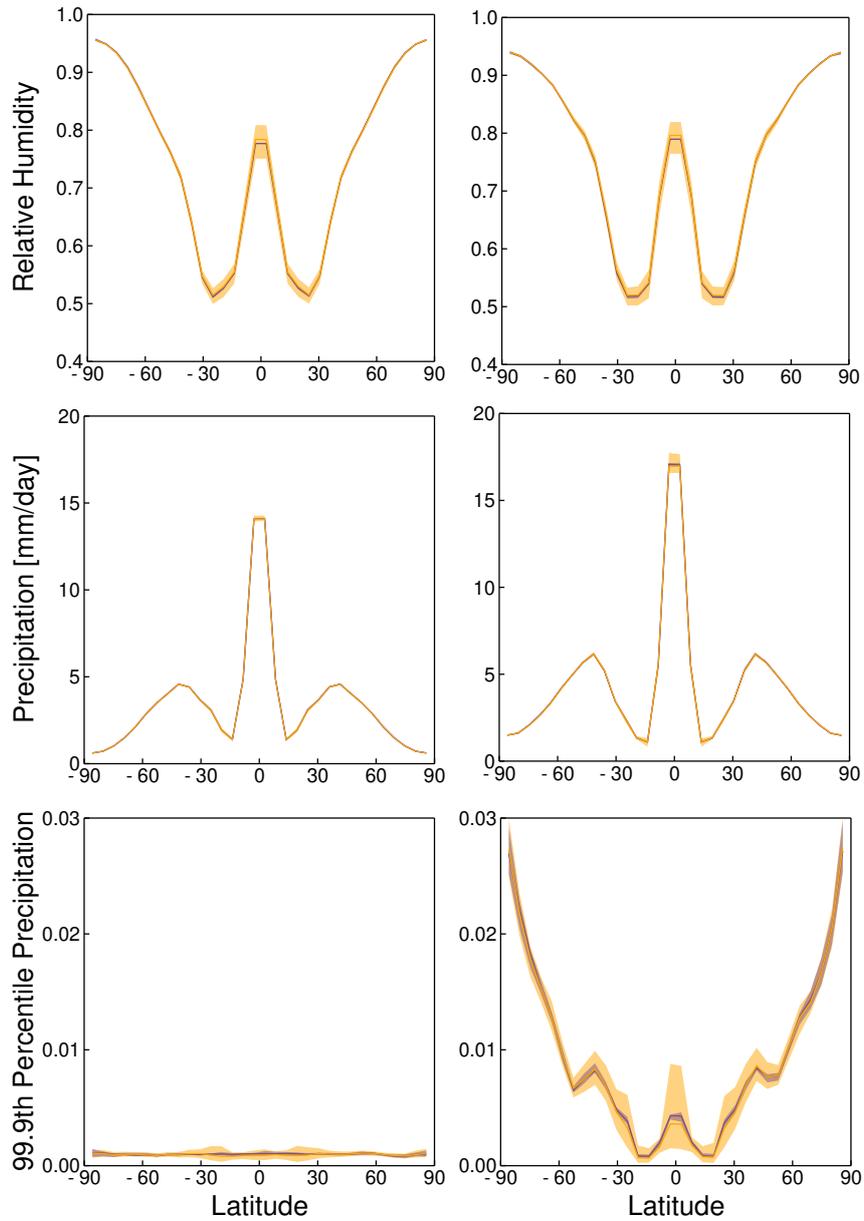
498

499

There is sampling variability because of the different data realizations. This sampling variability can be assessed by asking which probability contours contain the true parameters. For both EKI-GP and B-GP, in three of four realizations, we capture the true values within 50% of the posterior probability mass; the realization  $\mathbf{y}^{(3)}$  is captured only within the 99% contour of the posterior probability.



**Figure 10.** Density plot of MCMC samples of the posterior distribution. The contours are drawn to contain 50%, 75%, and 99% of the distribution generated from the samples. The left column show distributions learned using EKI-GP at  $\mathbf{y}^{(1)}, \dots, \mathbf{y}^{(4)}$ , and the right column using B-GP at the same data realizations. The blue dot represents the true parameters, while the red plus is the average across particles in the 6th EKI iteration.



**Figure 11.** Comparison of statistics of a 7200-day average in a climate-change simulation. Left column: control climate; right column: warmer climate. Synthetic observational data evaluated at the true fixed parameters are shown in blue, while data evaluated at 100 samples from the posterior distribution (EKI-GP) are shown in orange. (We choose the posterior from the first data realization, top-left panel of Figure 10.) The solid lines are the medians, and the shaded regions represent the 95% confidence intervals (between the 2.5th and 97.5th percentiles). Top: Relative humidity in mid-troposphere. Middle: Precipitation rate. Bottom: Frequency with which 99.9th percentile of latitude-dependent daily precipitation in the control climate is exceeded.

#### 4.4 Uncertainty Quantification in Prediction Experiments

To illustrate how the posterior distribution of parameters obtained in the sample step of the CES algorithm can be used to produce climate predictions with quantified uncertainties, we consider an idealized global-warming experiment. As in O’Gorman and Schneider (2008a, 2008b), we rescale the longwave opacity of the atmosphere everywhere by a uniform factor  $\alpha$ . In the control climate we have considered so far,  $\alpha = 1$ . We generate a warm climate by setting  $\alpha = 1.5$ , which results in a global-mean surface air temperature increase from 287 K in the control climate to 294 K in the warm climate. To see parametric uncertainty rather than internal variability noise in the resulting “climate change predictions,” we use long (7,200-day or approximately 20-year) averages in the prediction experiments.

We evaluate predictions of the latitude-dependent relative humidity and mean precipitation rate that we used in the CES algorithm. We also consider the frequency of precipitation extremes, now taken as the frequency with which the 99.9th percentile of daily precipitation in the control simulation is exceeded (rather than the 90th percentile we considered earlier). This last statistic indicates how the frequency of what are 1-in-1000 day precipitation events in the control climate change in the warmer climate.

We investigate the effect of parametric uncertainty on predictions by taking 100 samples of parameters from the posterior, creating a prediction for each sample, and comparing statistics of these runs with runs in which parameters are fixed to the true values  $\theta^\dagger$ . The climate statistics in the control climate are shown in the left column of Figure 11. The runs from posterior samples (orange) and with fixed true parameters (blue) match well. The noise due to internal variability is quantitatively represented by the blue shaded region. Unlike in the earlier figures with short (30-day) averages (e.g., Figure 9), the internal variability noise here is small relative to the parametric uncertainty because of the (long) 7200-day averaging window. The orange shaded region contains both internal variability and parametric uncertainty and is dominated by parametric uncertainty. This remains the case in the warmer climate (right column of Figure 11).

The effects of global warming on atmospheric quantities is seen by comparing the two columns of Figure 11. Relative humidity is fairly robust to the warming climate, and precipitation rates increase globally (O’Gorman & Schneider, 2008b). The most dramatic changes occur for the frequency of extreme precipitation events (O’Gorman & Schneider, 2009b). What is a 1-in-1000 day event in the control climate (e.g., occurring with frequency 0.001) occurs in the extratropics of the warmer climate an order of magnitude more frequently, with the 95% confidence interval spanning 0.01 to 0.03. That is, a 1-in-1000 day event in the control climate occurs every 30 to 100 days in the warmer climate. The parametric uncertainty is particularly large for extreme precipitation events within the tropics—behavior one would not be able to see in global warming experiments with fixed parameters. This is consistent with the known high uncertainty in predictions of tropical rainfall extremes with comprehensive climate models (O’Gorman & Schneider, 2009a).

## 5 Conclusion and Discussion

The primary goal of this article was to demonstrate that ensemble Kalman inversion (EKI), machine learning, and MCMC algorithms can be judiciously combined within the calibrate-emulate-sample framework to efficiently estimate uncertainty of model parameters in computationally expensive climate models. We provided a proof-of-concept in a relatively simple idealized GCM.

Our approach is novel because we train a machine learning (GP) emulator using input-output pairs generated from an EKI algorithm. This methodology has several advantageous features:

- 550 1. It requires a modest number of runs of the expensive forward model (typically,  $O(100)$   
551 runs).
- 552 2. It generally finds optimal or nearly optimal parameters even in the presence of in-  
553 ternal variability noise because EKI is robust with respect to such noise.
- 554 3. The resulting GP emulation is naturally most accurate around the (a priori un-  
555 known) optimal parameters because this is where EKI training points concentrate.
- 556 4. MCMC shows robust convergence to the posterior distribution, and allows iden-  
557 tification of the optimal parameters with the maximum of the posterior probabili-  
558 ty, because it uses an objective function that is smoothed by GP emulation.

559 The effectiveness of GP depends on the training points, and a user must choose how many  
560 iterations of EKI to use for training (before ensemble collapse). In practice, we find the  
561 GP performance is robust as long as we include the initial iteration of training points  
562 (drawn from the prior) in our training set. The necessity of using the initial ensemble  
563 could be side-stepped by using an ensemble method that does not collapse, such as the  
564 recently introduced ensemble Kalman sampler (EKS) (Garbuno-Inigo et al., 2020).

565 The CES algorithm is efficient, as it addresses two dominant sources of computa-  
566 tional expense. First, poor prior knowledge of model parameters requires blind explo-  
567 ration of a possibly high-dimensional parameter space to find optimal parameters and  
568 thus the region of high posterior probability. The CES framework handles this with an  
569 EKI algorithm, which we show to be successful when using time averaged data from a  
570 chaotic nonlinear model. Second, computing parametric uncertainty with a sampling tech-  
571 nique (such as MCMC) generally requires many ( $10^5$ – $10^6$ ) evaluations of an expensive  
572 forward model. We instead solve a cheap approximate problem by exploiting GP em-  
573 ulators. We train the emulators on relatively few ( $O(100)$ ) intelligently chosen evalua-  
574 tions provided by EKI, which ensures that training points are placed where they are most  
575 needed—near the minimum of the model-data misfit. The training itself introduces neg-  
576 ligible computational cost relative to the running of the forward model, and the com-  
577 putational expense of evaluating the emulator in the sampling step is also negligible. Hence,  
578 the CES framework achieves about a factor 1000 speedup over brute-force MCMC al-  
579 gorithms. Significant efforts to accelerate brute-force MCMC without approximation have  
580 been undertaken (Järvinen et al., 2010; Solonen et al., 2012), and improvements of up  
581 to a factor 5 speedup have been made with adaptive and parallelized Markov chains. How-  
582 ever, these approaches still are considerably more expensive than the CES algorithm.

583 The CES algorithm also has a smoothing property, which is beneficial even in sit-  
584 uations where a forward model is cheap enough to apply a brute-force MCMC. If the for-  
585 ward model exhibits internal variability, the objective function for the sampling algorithm  
586 will contain a data misfit of the form (7), which has a random component because it con-  
587 tains a finite-time average. Without more sophisticated sampling methods, MCMC al-  
588 gorithms get stuck in noise-induced local minima. In the CES algorithm, only EKI uses  
589 the functional (7), and EKI is well suited for this purpose. The GP emulator learns the  
590 smooth, noiseless model  $\mathcal{G}_\infty$  (in which internal variability disappears), using evaluations  
591 of  $\mathcal{G}_T$  (which are affected by internal variability). Thus, MCMC within the CES algo-  
592 rithm uses the smooth GP approximation of (6).

593 The MCMC results in this study successfully capture the true parameters and their  
594 uncertainties. The results contain natural biases arising from the use of prior distribu-  
595 tions, internal variability of the climate, and use of a single noisy sample as synthetic  
596 data. Despite the sampling variability and emulator constraints, our MCMC samples were  
597 able to capture the true parameters within an estimated 99% confidence interval in our  
598 examples, demonstrating the potential of EKI-trained GP emulators for MCMC sam-  
599 pling. Validation of the emulator (Figure 8) further supports the MCMC results, as do  
600 our comparisons with MCMC using the benchmark emulator (Table 1). The GP emu-  
601 lator both smooths the objective function and allows us to quantify uncertainty by sam-

602 pling from the posterior distribution. This contrasts with uncertainty quantification based  
 603 on the EKI ensemble, which underestimates the true uncertainties in our experiments  
 604 by an order of magnitude. As used here, EKI should be viewed as an optimization al-  
 605 gorithm and not a sampling algorithm. Adding additional spread to match the poste-  
 606 rior within EKI may be achieved for Gaussian posteriors (Chen & Oliver, 2012; Emer-  
 607 ick & Reynolds, 2013) or by means of EKS (Garbuno-Inigo et al., 2020); however, these  
 608 methods are not justifiable beyond the Gaussian setting. The MCMC algorithm within  
 609 CES, on the other hand, samples from an approximate posterior distribution and is jus-  
 610 tifiable beyond the Gaussian posterior setting (Cleary et al., 2021).

611 Good scaling of the CES algorithm with the dimension of the parameter space will  
 612 be of critical importance for moving beyond the current, low-dimensional proof-of-concept  
 613 setting. Each stage of the CES algorithm scales to higher dimensions: For the calibra-  
 614 tion stage, ensemble methods scale well to high-dimensional state and parameter spaces,  
 615 typically with  $O(10^2)$  forward model runs (Kalnay, 2003; Oliver et al., 2008), if used with  
 616 localization. In high dimensions, regularization is also typically needed in calibration al-  
 617 gorithms, and various regularization schemes can be added (Chada et al., 2020; Iglesias,  
 618 2015, 2016; Garbuno-Inigo et al., 2019; Schneider et al., 2020b). The sampling stage also  
 619 scales well to high dimensions. In particular, MCMC methods scale to non-parametric  
 620 (infinite-dimensional) problems, in which the unknown is a function (Cotter et al., 2013).  
 621 The current bottleneck for scalability lies with the GP emulator, which does not scale  
 622 easily to high-dimensional inputs. However, other supervised machine learning techniques  
 623 have potential to do so. For example, random feature maps and deep neural networks  
 624 show promise in this regard (Nelsen & Stuart, 2020; Bhattacharya et al., 2020); incor-  
 625 porating these tools in the CES algorithm is a direction of current research.

626 An alternative form of constraining parameter uncertainty is history matching, or  
 627 precalibration (Vernon et al., 2010; Edwards et al., 2011; Williamson et al., 2013). The  
 628 idea complements that of Bayesian uncertainty quantification, where instead of search-  
 629 ing for a high probability region of parameter space with respect to data, one rules out  
 630 regions of parameter space that are deemed inconsistent with the data. Couvreur et al.  
 631 (2020) and Hourdin et al. (2020) recently constrained the parameter space of a param-  
 632 eterization scheme by approximating a plausibility function over the parameter space us-  
 633 ing a Gaussian process, and then removing “implausible” regions of parameter space where  
 634 the plausibility function passes a threshold. This removal process is iterated until the  
 635 uncertainty of the emulator is small enough, or the space becomes empty. History match-  
 636 ing accomplishes a similar adaptivity task as that performed in the CES algorithm by  
 637 EKI. During early stages of history matching, however, one must sample the full param-  
 638 eter space with reasonable resolution, and emulator training is required at every itera-  
 639 tion to evaluate the plausibility function. In contrast, in the CES algorithm, EKI draws  
 640 a modest numbers of samples at every iteration and can work directly with noisy model  
 641 evaluations, lowering the computational expense. The output of history matching is a  
 642 (possibly empty) acceptable set of forward model runs; sampling this set leads to an up-  
 643 per bound on the prediction uncertainty. The benefit of the CES algorithm is that it pro-  
 644 vides samples of the posterior distribution, which lead to full estimates of prediction un-  
 645 certainty (Figure 11). For this reason, history matching has been proposed as a prepro-  
 646 cessing step for Bayesian uncertainty quantification, known as precalibration to improve  
 647 priors and assess model validity (Vernon et al., 2010; Edwards et al., 2011). The CES  
 648 algorithm targets the Bayesian posterior distributions directly.

649 In the more comprehensive climate modeling settings we target, data will be given  
 650 from earth observations and from local high-resolution simulations (Schneider, Lan, et  
 651 al., 2017). In these settings, model error leads to deficiencies when comparing model eval-  
 652 uations to data, leading to structural biases and uncertainty that must be quantified. Struc-  
 653 tural model errors can be quantified with a flexible hierarchical Gaussian process regres-  
 654 sion that learns a non-parametric form of the model deficiency, as demonstrated in pro-

655 prototype problems in Schneider et al. (2020a). This approach represents model error in an  
 656 interpretable fashion, as part of the model itself, rather than in the data space as pio-  
 657 neered in Kennedy and O’Hagan (2001).

658 The CES framework has potential for both the calibration and uncertainty quan-  
 659 tification of comprehensive climate models, and other computationally expensive mod-  
 660 els. It is computationally efficient enough to use data averaged in time (e.g., over sea-  
 661 sons), which need to be accumulated over longer model runs. Time-averaged climate statis-  
 662 tics, including mean values and higher-order statistics such as extreme value statistics,  
 663 are what typically matters in climate predictions. CES allows us to focus model calibra-  
 664 tion and uncertainty quantification on such immediately relevant statistics. Using time  
 665 averaged statistics also has the advantage that it leads to smoother, albeit still noisy, ob-  
 666 jective functions when compared with calibration of climate models by minimizing mis-  
 667 matches in instantaneous, short-term forecasts (Schneider, Lan, et al., 2017). The lat-  
 668 ter approach can improve short-term forecasts but may not translate into improved cli-  
 669 mate simulations (Schirber et al., 2013). It also suffers from the difficulty that model res-  
 670 olution and data resolution may be mismatched. Focusing on climate statistics, as we  
 671 did in our proof-of-concept here, circumvents this problem: time-aggregated climate statis-  
 672 tics are varying relatively smoothly in space and, hence, minimizing mismatches in statis-  
 673 tics between models and data does not suffer from the resolution-mismatch problem. CES  
 674 can be used to learn about arbitrary parameters in climate models from time-averaged  
 675 data. It leads to quantification of parametric uncertainties that then can be converted  
 676 into parametric uncertainties in predictions by sampling from the posterior distribution  
 677 of parameters.

## 678 Acknowledgments

679 This work was supported by the generosity of Eric and Wendy Schmidt by recommen-  
 680 dation of the Schmidt Futures program, by the Hopewell Fund, the Paul G. Allen Fam-  
 681 ily Foundation, and the National Science Foundation (NSF, award AGS1835860). A.M.S.  
 682 was also supported by the Office of Naval Research (award N00014-17-1-2079). We thank  
 683 Emmet Cleary for his preliminary work underlying some of the results shown here.

684 **Data Availability.** All computer code used in this paper is open source. The code for  
 685 the idealized GCM, the Julia code for the CES algorithm, the plot tools, and the slurm/bash  
 686 scripts to run both GCM and CES are available at <https://doi.org/10.5281/zenodo.4393029>.

## 687 References

- 688 Annan, J. D., & Hargreaves, J. C. (2007). Efficient estimation and ensemble gener-  
 689 ation in climate modelling. *Phil. Trans. R. Soc. A*, *365*, 2077–2088. doi: 10  
 690 .1098/rsta.2007.2067
- 691 Betts, A. K. (1986). A new convective adjustment scheme. Part I: Observational and  
 692 theoretical basis. *Quart. J. Roy. Meteor. Soc.*, *112*, 677–691.
- 693 Betts, A. K., & Miller, M. J. (1986). A new convective adjustment scheme. Part II:  
 694 Single column tests using GATE wave, BOMEX, ATEX and arctic air-mass  
 695 data sets. *Quart. J. Roy. Meteor. Soc.*, *112*, 693–709.
- 696 Betts, A. K., & Miller, M. J. (1993). The Betts–Miller scheme. In K. A. Emanuel &  
 697 D. J. Raymond (Eds.), *The representation of cumulus convection in numerical*  
 698 *models* (Vol. 24, pp. 107–121). Am. Meteor. Soc.
- 699 Bhattacharya, K., Hosseini, B., Kovachki, N. B., & Stuart, A. M. (2020).  
 700 Model reduction and neural networks for parametric pdes. *arXiv preprint*  
 701 *arXiv:2005.03180*.
- 702 Bischoff, T., & Schneider, T. (2014). Energetic constraints on the position of the In-  
 703 ter-tropical Convergence Zone. *J. Climate*, *27*, 4937–4951. doi: 10.1175/JCLI-D

- 704 -13-00650.1
- 705 Bocquet, M., & Sakov, P. (2012). Combining inflation-free and iterative ensemble  
706 kalman filters for strongly nonlinear systems. *Nonlinear Processes in Geo-*  
707 *physics*, *19*(3), 383–399.
- 708 Bocquet, M., & Sakov, P. (2014). An iterative ensemble Kalman smoother. *Q. J. R.*  
709 *Meteorol. Soc.*, *140*, 1521–1535. doi: 10.1002/qj.2236
- 710 Bony, S., Colman, R., Kattsov, V. M., Allan, R. P., Bretherton, C. S., Dufresne,  
711 J.-L., ... Webb, M. J. (2006). How well do we understand and evalu-  
712 ate climate change feedback processes? *J. Climate*, *19*, 3445–3482. doi:  
713 10.1175/JCLI3819.1
- 714 Bony, S., & Dufresne, J. L. (2005). Marine boundary layer clouds at the heart of  
715 tropical cloud feedback uncertainties in climate models. *Geophys. Res. Lett.*,  
716 *32*, L20806.
- 717 Bordoni, S., & Schneider, T. (2008). Monsoons as eddy-mediated regime transitions  
718 of the tropical overturning circulation. *Nature Geosci.*, *1*, 515–519. doi: 10  
719 .1038/ngeo248
- 720 Brient, F., & Schneider, T. (2016). Constraints on climate sensitivity from space-  
721 based measurements of low-cloud reflection. *J. Climate*, *29*, 5821–5835. doi: 10  
722 .1175/JCLI-D-15-0897.1
- 723 Cess, R. D., Potter, G., Blanchet, J., Boer, G., Ghan, S., Kiehl, J., ... others  
724 (1989). Interpretation of cloud-climate feedback as produced by 14 atmo-  
725 spheric general circulation models. *Science*, *245*, 513–516.
- 726 Cess, R. D., Potter, G. L., Blanchet, J. P., Boer, G. J., Del Genio, A. D., Déqué,  
727 M., ... Zhang, M.-H. (1990). Intercomparison and interpretation of climate  
728 feedback processes in 19 atmospheric general circulation models. *J. Geophys.*  
729 *Res.*, *95*, 16601–16615. doi: 10.1029/JD095iD10p16601
- 730 Chada, N. K., Stuart, A. M., & Tong, X. T. (2020). Tikhonov regularization within  
731 ensemble kalman inversion. *SIAM Journal on Numerical Analysis*, *58*(2),  
732 1263–1294.
- 733 Chen, Y., & Oliver, D. S. (2012). Ensemble randomized maximum likelihood method  
734 as an iterative ensemble smoother. *Mathematical Geosciences*, *44*(1), 1–26.
- 735 Cleary, E., Garbuno-Inigo, A., Lan, S., Schneider, T., & Stuart, A. M. (2021). Cal-  
736 ibrate, emulate, sample. *J. Comp. Phys.*, *424*, 109716. Retrieved from 10  
737 .1016/j.jcp.2020.109716
- 738 Cotter, S. L., Roberts, G. O., Stuart, A. M., & White, D. (2013). MCMC methods  
739 for functions: Modifying old algorithms to make them faster. *Statist. Science*,  
740 *28*(3), 424–446.
- 741 Couvreur, F., Hourdin, F., Williamson, D., Roehrig, R., Volodina, V., Villefranche,  
742 N., ... Xu, W. (2020). Process-based climate model development harnessing  
743 machine learning: I. A calibration tool for parameterization improvement.  
744 doi: 10.1002/essoar.10503597.1
- 745 Duncan, A. B., Stuart, A. M., & Wolfram, M.-T. (2021). Ensemble inference  
746 methods for models with noisy and expensive likelihoods. *arXiv preprint*  
747 *arXiv:2104.03384*.
- 748 Edwards, N. R., Cameron, D., & Rougier, J. (2011). Precalibrating an intermediate  
749 complexity climate model. *Climate Dynamics*, *37*(7-8), 1469–1482.
- 750 Emerick, A. A., & Reynolds, A. C. (2013). Ensemble smoother with multiple data  
751 assimilation. *Computers & Geosciences*, *55*, 3–15.
- 752 Ernst, O. G., Sprungk, B., & Starkloff, H.-J. (2015). Analysis of the ensemble  
753 and polynomial chaos Kalman filters in Bayesian inverse problems.  
754 *SIAM/ASA Journal on Uncertainty Quantification*, *3*(1), 823–851. doi:  
755 10.1137/140981319
- 756 Evensen, G. (1994). Sequential data assimilation with a nonlinear quasi-geostrophic  
757 model using monte carlo methods to forecast error statistics. *Journal of Geo-*  
758 *physical Research: Oceans*, *99*(C5), 10143–10162.

- 759 Evensen, G. (2018). Analysis of iterative ensemble smoothers for solving inverse  
760 problems. *Comp. Geosci.*. doi: 10.1007/s10596-018-9731-y
- 761 Flato, G., Marotzke, J., Abiodun, B., Braconnot, P., Chou, S. C., Collins, W., ...  
762 Rummukainen, M. (2013). Evaluation of climate models. In T. F. Stocker  
763 et al. (Eds.), *Climate change 2013: The physical science basis. contribution of*  
764 *working group i to the fifth assessment report of the intergovernmental panel*  
765 *on climate change* (pp. 741–853). Cambridge, UK, and New York, NY, USA:  
766 Cambridge University Press.
- 767 Frierson, D. M. W. (2007). The dynamics of idealized convection schemes and their  
768 effect on the zonally averaged tropical circulation. *J. Atmos. Sci.*, *64*, 1959–  
769 1976.
- 770 Frierson, D. M. W., Held, I. M., & Zurita-Gotor, P. (2006). A gray-radiation aqua-  
771 planet moist GCM. Part I: Static stability and eddy scale. *J. Atmos. Sci.*, *63*,  
772 2548–2566.
- 773 Garbuno-Inigo, A., Hoffmann, F., Li, W., & Stuart, A. M. (2019). Interacting  
774 Langevin Diffusions: Gradient Structure And Ensemble Kalman Sampler. ,  
775 1–34. Retrieved from <http://arxiv.org/abs/1903.08866>
- 776 Garbuno-Inigo, A., Hoffmann, F., Li, W., & Stuart, A. M. (2020). Interact-  
777 ing langevin diffusions: Gradient structure and ensemble Kalman sampler.  
778 *SIAM Journal on Applied Dynamical Systems*, *19*(1), 412-441. Retrieved from  
779 <https://doi.org/10.1137/19M1251655> doi: 10.1137/19M1251655
- 780 Geyer, C. J. (2011). Introduction to Markov Chain Monte Carlo. In S. Brooks,  
781 A. Gelman, G. L. Jones, & X.-L. Meng (Eds.), *Handbook of Markov Chain*  
782 *Monte Carlo*. Chapman and Hall/CRC Press. doi: 10.1201/b10905-3
- 783 Gland, F. L., Monbet, V., & Tran, V.-D. (2009). *Large sample asymptotics for the*  
784 *ensemble kalman filter* (Tech. Rep. No. RR-7014). INRIA.
- 785 Golaz, J.-C., Horowitz, L. W., & II, H. L. (2013). Cloud tuning in a coupled climate  
786 model: Impact on 20th century warming. *Geophys. Res. Lett.*, *40*, 2246–2251.  
787 doi: 10.1002/grl.50232
- 788 Gu, Y., & Oliver, D. S. (2007). An iterative ensemble kalman filter for multiphase  
789 fluid flow data assimilation. *SPE Journal*, *12*(04), 438–446.
- 790 Hourdin, F., Grandpeix, J.-Y., Rio, C., Bony, S., Jam, A., Cheruy, F., ... Roehrig,  
791 R. (2013). LMDZ5B: the atmospheric component of the IPSL climate model  
792 with revisited parameterizations for clouds and convection. *Clim. Dyn.*, *40*,  
793 2193–2222. doi: 10.1007/s00382-012-1343-y
- 794 Hourdin, F., Mauritsen, T., Gettelman, A., Golaz, J.-C., Balaji, V., Duan, Q., ...  
795 Williamson, D. (2017). The art and science of climate model tuning. *Bull.*  
796 *Amer. Meteor. Soc.*, *98*, 589–602. doi: 10.1175/BAMS-D-15-00135.1
- 797 Hourdin, F., Williamson, D., Rio, C., Couvreur, F., Roehrig, R., Villefranche,  
798 N., ... Volodina, V. (2020). Process-based climate model develop-  
799 ment harnessing machine learning: II. Model calibration from single col-  
800 umn to global. *Journal of Advances in Modeling Earth Systems*. doi:  
801 <https://doi.org/10.1029/2020MS002225>
- 802 Houtekamer, P. L., & Zhang, F. (2016). Review of the ensemble Kalman filter for  
803 atmospheric data assimilation. *Mon. Wea. Rev.*, *144*, 4489–4532. doi: 10.1175/  
804 MWR-D-15-0440.1
- 805 Iglesias, M. A. (2015). Iterative regularization for ensemble data assimilation in  
806 reservoir models. *Computational Geosciences*, *19*(1), 177–212.
- 807 Iglesias, M. A. (2016). A regularizing iterative ensemble kalman method for pde-  
808 constrained inverse problems. *Inverse Problems*, *32*(2), 025002.
- 809 Iglesias, M. A., Law, K. J. H., & Stuart, A. M. (2013, mar). Ensemble Kalman  
810 methods for inverse problems. *Inverse Problems*, *29*(4), 045001. doi: 10.1088/  
811 0266-5611/29/4/045001
- 812 Järvinen, H., Räisänen, P., Laine, M., Tamminen, J., Ilin, A., Oja, E., ... Haario,  
813 H. (2010). Estimation of ECHAM5 climate model closure parame-

- 814           ters with adaptive MCMC. *Atmos. Chem. Phys.*, *10*, 9993–10002.   doi:  
815           10.5194/acp-10-9993-2010
- 816 Kalnay, E. (2002). *Atmospheric modeling, data assimilation and predictability*. Cam-  
817           bridge University Press. doi: 10.1017/CBO9780511802270
- 818 Kalnay, E. (2003). *Atmospheric modeling, data assimilation and predictability*. Cam-  
819           bridge, UK: Cambridge Univ. Press.
- 820 Kaspi, Y., & Schneider, T. (2011). Winter cold of eastern continental boundaries in-  
821           duced by warm ocean waters. *Nature*, *471*, 621–624.
- 822 Kaspi, Y., & Schneider, T. (2013). The role of stationary eddies in shaping midlati-  
823           tude storm tracks. *J. Atmos. Sci.*, *70*, 2596–2613.
- 824 Kennedy, M. C., & O’Hagan, A. (2001). Bayesian calibration of computer models. *J.*  
825           *Roy. Statist. Soc. B*, *63*, 425–464. doi: 10.1111/1467-9868.00294
- 826 Levine, X., & Schneider, T. (2015). Baroclinic eddies and the extent of the Hadley  
827           circulation: An idealized GCM study. *J. Atmos. Sci.*, *72*, 2744–2761. doi: 10  
828           .1175/JAS-D-14-0152.1
- 829 Li, G., & Reynolds, A. C. (2009). An iterative ensemble kalman filter for data assim-  
830           ilation. In (Vol. 14, pp. 496–505). Society of Petroleum Engineers.
- 831 Mauritsen, T., Stevens, B., Roeckner, E., Crueger, T., Esch, M., Giorgetta, M., ...  
832           Tomassini, L. (2012). Tuning the climate of a global model. *J. Adv. Model.*  
833           *Earth Sys.*, *4*, M00A01. doi: 10.1029/2012MS000154
- 834 Merlis, T. M., & Schneider, T. (2011). Changes in zonal surface temperature gradi-  
835           ents and walker circulations in a wide range of climates. *J. Climate*, *24*, 4757–  
836           4768.
- 837 Nelsen, N. H., & Stuart, A. M. (2020). The random feature model for input-output  
838           maps between banach spaces. *arXiv preprint arXiv:2005.10224*.
- 839 O’Gorman, P. A. (2011). The effective static stability experienced by eddies in a  
840           moist atmosphere. *J. Atmos. Sci.*, *68*, 75–90.
- 841 O’Gorman, P. A., Lamquin, N., Schneider, T., & Singh, M. S. (2011). The relative  
842           humidity in an isentropic advection–condensation model: Limited poleward in-  
843           fluence and properties of subtropical minima. *J. Atmos. Sci.*, *68*, 3079–3093.
- 844 O’Gorman, P. A., & Schneider, T. (2008a). Energy of midlatitude transient eddies  
845           in idealized simulations of changed climates. *J. Climate*, *21*, 5797–5806.
- 846 O’Gorman, P. A., & Schneider, T. (2008b). The hydrological cycle over a wide range  
847           of climates simulated with an idealized GCM. *J. Climate*, *21*, 3815–3832.
- 848 O’Gorman, P. A., & Schneider, T. (2009a). The physical basis for increases in pre-  
849           cipitation extremes in simulations of 21st-century climate change. *Proc. Natl.*  
850           *Acad. Sci.*, *106*, 14773–14777.
- 851 O’Gorman, P. A., & Schneider, T. (2009b). Scaling of precipitation extremes over  
852           a wide range of climates simulated with an idealized GCM. *J. Climate*, *22*,  
853           5676–5685.
- 854 Oliver, D. S., Reynolds, A. C., & Liu, N. (2008). *Inverse theory for petroleum reser-*  
855           *voir characterization and history matching*. Cambridge Univ. Press.
- 856 Pedregosa, F., Varoquaux, G., Gramfort, A., Michel, V., Thirion, B., Grisel, O., ...  
857           Duchesnay, E. (2011). Scikit-learn: Machine learning in Python. *Journal of*  
858           *Machine Learning Research*, *12*, 2825–2830.
- 859 Randall, D. A., & Wielicki, B. A. (1997). Measurements, models, and hypotheses in  
860           the atmospheric sciences. *Bull. Amer. Meteor. Soc.*, *78*, 400–406.
- 861 Reich, S. (2011). A dynamical systems framework for intermittent data assimilation.  
862           *BIT Numerical Mathematics*, *51*(1), 235–249.
- 863 Roberts, G. O., Rosenthal, J. S., et al. (2004). General state space Markov chains  
864           and mcmc algorithms. *Probability surveys*, *1*, 20–71.
- 865 Sacks, J., Welch, W. J., Mitchell, T. J., & Wynn, H. P. (1989). Design and analysis  
866           of computer experiments. *Statistical science*, 409–423.
- 867 Sakov, P., Oliver, D. S., & Bertino, L. (2012). An iterative EnKF for strongly non-  
868           linear systems. *Mon. Wea. Rev.*, *140*, 1988–2004. doi: 10.1175/MWR-D-11

869 -00176.1

- 870 Santner, T. J., Williams, B. J., Notz, W., & Williams, B. J. (2018). *The design and*  
871 *analysis of computer experiments* (2nd ed.). New York, NY: Springer.
- 872 Schillings, C., & Stuart, A. M. (2017a). Analysis of the ensemble kalman fil-  
873 ter for inverse problems. *SIAM Journal on Numerical Analysis*, *55*(3),  
874 1264–1290. Retrieved from <https://doi.org/10.1137/16M105959X> doi:  
875 10.1137/16M105959X
- 876 Schillings, C., & Stuart, A. M. (2017b). Analysis of the ensemble Kalman filter for  
877 inverse problems. *SIAM J. Numer. Anal.*, *55*, 1264–1290.
- 878 Schirber, S., Klocke, D., Pincus, R., Quaas, J., & Anderson, J. L. (2013, mar).  
879 Parameter estimation using data assimilation in an atmospheric general circula-  
880 tion model: From a perfect toward the real world. *Journal of Advances in*  
881 *Modeling Earth Systems*, *5*(1), 58–70. doi: 10.1029/2012MS000167
- 882 Schmidt, G. A., Bader, D., Donner, L. J., Elsaesser, G. S., Golaz, J.-C., Hannay,  
883 C., . . . Saha, S. (2017). Practice and philosophy of climate model tuning  
884 across six u.s. modeling centers. *Geosci. Model Dev.*, *10*, 3207–3223. doi:  
885 10.5194/gmd-2017-30
- 886 Schneider, T., Lan, S., Stuart, A., & Teixeira, J. (2017). Earth system model-  
887 ing 2.0: A blueprint for models that learn from observations and targeted  
888 high-resolution simulations. *Geophys. Res. Lett.*, *44*, 12396–12417. doi:  
889 10.1002/2017GL076101
- 890 Schneider, T., & O’Gorman, P. A. (2008). Moist convection and the thermal stratifi-  
891 cation of the extratropical troposphere. *J. Atmos. Sci.*, *65*, 3571–3583.
- 892 Schneider, T., O’Gorman, P. A., & Levine, X. J. (2010). Water vapor  
893 and the dynamics of climate changes. *Rev. Geophys.*, *48*, RG3001.  
894 (doi:10.1029/2009RG000302)
- 895 Schneider, T., Stuart, A. M., & Wu, J.-L. (2020a). Ensemble Kalman inversion for  
896 sparse learning of dynamical systems from time-averaged data. *arXiv preprint*,  
897 *arXiv:2007.06175*.
- 898 Schneider, T., Stuart, A. M., & Wu, J.-L. (2020b). Imposing sparsity within ensem-  
899 ble kalman inversion. *arXiv preprint arXiv:2007.06175*.
- 900 Schneider, T., Stuart, A. M., & Wu, J.-L. (2020c). Learning stochastic closures using  
901 ensemble Kalman inversion. *arXiv preprint, arXiv:2004.08376*.
- 902 Schneider, T., Teixeira, J., Bretherton, C. S., Brient, F., Pressel, K. G., Schär, C.,  
903 & Siebesma, A. P. (2017). Climate goals and computing the future of clouds.  
904 *Nature Climate Change*, *7*, 3–5. doi: 10.1038/nclimate3190
- 905 Solonen, A., Ollinaho, P., Laine, M., Haario, H., Tamminen, J., & Järvinen, H.  
906 (2012). Efficient mcmc for climate model parameter estimation: Parallel  
907 adaptive chains and early rejection. *Bayesian Analysis*, *7*(3), 715–736. doi:  
908 10.1214/12-BA724
- 909 Stephens, G. L. (2005). Cloud feedbacks in the climate system: A critical review. *J.*  
910 *Climate*, *18*, 237–273. doi: 10.1175/JCLI-3243.1
- 911 Van Leeuwen, P. J., & Evensen, G. (1996). Data assimilation and inverse meth-  
912 ods in terms of a probabilistic formulation. *Monthly Weather Review*, *124*(12),  
913 2898–2913.
- 914 Vernon, I., Goldstein, M., & Bower, R. G. (2010, 12). Galaxy formation: A  
915 Bayesian uncertainty analysis. *Bayesian Analysis*, *5*(4), 619–669. doi:  
916 10.1214/10-BA524
- 917 Vial, J., Dufresne, J.-L., & Bony, S. (2013). On the interpretation of inter-model  
918 spread in CMIP5 climate sensitivity estimates. *Clim. Dyn.*, *41*, 3339–3362. doi:  
919 10.1007/s00382-013-1725-9
- 920 Webb, M. J., Lambert, F. H., & Gregory, J. M. (2013). Origins of differences in cli-  
921 mate sensitivity, forcing and feedback in climate models. *Clim. Dyn.*, *40*, 677–  
922 707. doi: 10.1007/s00382-012-1336-x

- 923 Wei, H.-H., & Bordoni, S. (2018). Energetic constraints on the ITCZ position in ide-  
 924 alized simulations with a seasonal cycle. *J. Adv. Model. Earth Sys.*, *10*. doi: 10  
 925 .1029/2018MS001313
- 926 Williamson, D., Goldstein, M., Allison, L., Blaker, A., Challenor, P., Jackson, L.,  
 927 & Yamazaki, K. (2013). History matching for exploring and reducing cli-  
 928 mate model parameter space using observations and a large perturbed physics  
 929 ensemble. *Climate dynamics*, *41*(7-8), 1703–1729.
- 930 Wills, R. C., Levine, X. J., & Schneider, T. (2017). Local energetic constraints on  
 931 Walker circulation strength. *J. Atmos. Sci.*, *74*, 1907–1922. doi: 10.1175/JAS  
 932 -D-16-0219.1
- 933 Zhang, F., Sun, Y. Q., Magnusson, L., Buizza, R., Lin, S.-J., Chen, J.-H., &  
 934 Emanuel, K. (2019). What is the predictability limit of midlatitude  
 935 weather? *Journal of the Atmospheric Sciences*, *76*(4), 1077 - 1091. Re-  
 936 trieved from [https://journals.ametsoc.org/view/journals/atsc/76/4/  
 937 jas-d-18-0269.1.xml](https://journals.ametsoc.org/view/journals/atsc/76/4/jas-d-18-0269.1.xml) doi: 10.1175/JAS-D-18-0269.1
- 938 Zhao, M., Golaz, J.-C., Held, I. M., Guo, H., Balaji, V., Benson, R., ... Xiang,  
 939 B. (2018). The GFDL global atmosphere and land model AM4.0/LM4.0: 2.  
 940 Model description, sensitivity studies, and tuning strategies. *J. Adv. Model.  
 941 Earth Sys.*, *10*, 735–769. doi: 10.1002/2017MS001209