# Ensemble-Based Experimental Design for Targeted High-Resolution Simulations to Inform Climate Models

Oliver Dunbar<sup>1</sup>, Michael F. Howland<sup>2</sup>, Tapio Schneider<sup>1</sup>, and Andrew Stuart<sup>1</sup>

<sup>1</sup>California Institute of Technology <sup>2</sup>Massachusetts Institute of Technology

November 22, 2022

#### Abstract

Targeted high-resolution simulations driven by a general circulation model (GCM) can be used to calibrate GCM parameterizations of processes that are globally unresolvable but can be resolved in limited-area simulations. This raises the question of where to place high-resolution simulations to be maximally informative about the uncertain parameterizations in the global model. Here we construct an ensemble-based parallel algorithm to locate regions that maximize the uncertainty reduction, or information gain, in the uncertainty quantification of GCM parameters with regional data. The algorithm is based on a Bayesian framework that exploits a quantified posterior distribution on GCM parameters as a measure of uncertainty. The algorithm is embedded in the recently developed calibrate-emulate-sample (CES) framework, which performs efficient model calibration and uncertainty quantification with only  $O(10^2)$  forward model evaluations, compared with  $O(10^5)$  forward model evaluations typically needed for traditional approaches to Bayesian calibration. We demonstrate the algorithm with an idealized GCM, with which we generate surrogates of high-resolution data. In this setting, we calibrate parameters and quantify uncertainties in a quasi-equilibrium convection scheme. We consider (i) localization in space for a statistically stationary problem, and (ii) localization in space and time for a seasonally varying problem. In these proof-of-concept applications, the calculated information gain reflects the reduction in parametric uncertainty obtained from Bayesian inference when harnessing a targeted sample of data. The largest information gain results from regions near the intertropical convergence zone (ITCZ) and indeed the algorithm automatically targets these regions for data collection.

# Ensemble-Based Experimental Design for Targeted High-Resolution Simulations to Inform Climate Models

# Oliver R. A. Dunbar<sup>1</sup>, Michael F. Howland<sup>1,2</sup>, Tapio Schneider<sup>1</sup>, Andrew M. Stuart<sup>1</sup>

 $^1{\rm California}$ Institute of Technology, Pasadena, USA.<br/>  $^2{\rm Civil}$  and Environmental Engineering, Massachusetts Institute of Technology, Cambridge, MA, USA

# Key Points:

1

2

3

4

5 6

7

12

8	• Climate models can be calibrated with limited-area high-resolution simulations	5;
9	we address their optimal placement in space and time.	
10	• We propose an algorithm that places high-resolution simulations so that they a	are
11	maximally informative about climate model parameters.	

• The algorithm is benchmarked in a idealized aquaplanet general circulation model.

 $Corresponding \ author: \ Oliver \ Dunbar, \ \texttt{odunbar@caltech.edu}$ 

#### 13 Abstract

Targeted high-resolution simulations driven by a general circulation model (GCM) can 14 be used to calibrate GCM parameterizations of processes that are globally unresolvable 15 but can be resolved in limited-area simulations. This raises the question of where to place 16 high-resolution simulations to be maximally informative about the uncertain parame-17 terizations in the global model. Here we construct an ensemble-based parallel algorithm 18 to locate regions that maximize the uncertainty reduction, or information gain, in the 19 uncertainty quantification of GCM parameters with regional data. The algorithm is based 20 on a Bayesian framework that exploits a quantified posterior distribution on GCM pa-21 rameters as a measure of uncertainty. The algorithm is embedded in the recently devel-22 oped calibrate-emulate-sample (CES) framework, which performs efficient model cali-23 bration and uncertainty quantification with only  $\mathcal{O}(10^2)$  forward model evaluations, com-24 pared with  $\mathcal{O}(10^5)$  forward model evaluations typically needed for traditional approaches 25 to Bayesian calibration. We demonstrate the algorithm with an idealized GCM, with which 26 we generate surrogates of high-resolution data. In this setting, we calibrate parameters 27 and quantify uncertainties in a quasi-equilibrium convection scheme. We consider (i) lo-28 calization in space for a statistically stationary problem, and (ii) localization in space 29 and time for a seasonally varying problem. In these proof-of-concept applications, the 30 calculated information gain reflects the reduction in parametric uncertainty obtained from 31 Bayesian inference when harnessing a targeted sample of data. The largest information 32 gain results from regions near the intertropical convergence zone (ITCZ) and indeed the 33 algorithm automatically targets these regions for data collection. 34

#### <sup>35</sup> Plain Language Summary

Climate models depend on dynamics across many spatial and temporal scales. It 36 is infeasible to resolve all of these scales. Instead, the physics at the smallest scales is 37 represented by parameterization schemes that link what is unresolvable to variables re-38 solved on the grid scale. A dominant source of uncertainty in climate predictions comes 39 from uncertainty in calibrating empirical parameters in such parameterization schemes, 40 and this uncertainty is generally not quantified. Targeted high-resolution simulations of 41 small-scale processes in limited areas are one means by which uncertainties in param-42 eterizations can be reduced and quantified. Here we demonstrate an algorithm that op-43 timizes placement of high-resolution simulations to maximize the information they pro-44 vide about uncertain parameters in parameterization schemes. Because the sensitivity 45 of simulated climate statistics, such as precipitation rates, to parameterizations varies 46 in space and time, how informative high-resolution simulations are about the parame-47 terizations also varies in space and time. In proof-of-concept simulations with an ide-48 alized global atmosphere model, we show that our novel algorithm successfully identi-49 fies the informative regions and times. 50

#### 51 **1** Introduction

Parameterizations of subgrid-scale processes, such as the turbulence and convec-52 tion controlling clouds, are the principal cause of physical uncertainties in climate pre-53 dictions (Cess et al., 1989, 1990; Bony & Dufresne, 2005; Stephens, 2005; Bony et al., 54 2006; Vial et al., 2013; Webb et al., 2013; Brient & Schneider, 2016). While these pro-55 cesses are too small in scale to become globally resolvable in climate models for the fore-56 seeable future, many of them can be resolved in limited-area simulations (Schneider, Teix-57 eira, et al., 2017). For example, the turbulence and convection (though currently not the 58 microphysics) controlling clouds can be resolved in large-eddy simulations (LES) over 59 limited areas (Siebesma et al., 2003; Stevens et al., 2005; Khairoutdinov et al., 2009; Math-60 eou & Chung, 2014; Schalkwijk et al., 2015; Pressel et al., 2015, 2017). High-resolution 61 simulations have been used to calibrate climate model parameterizations at selected sites, 62

primarily in low latitudes (Liu et al., 2001; Siebesma et al., 2003; Stevens et al., 2005; 63 Siebesma et al., 2007; Hohenegger & Bretherton, 2011; M. Zhang et al., 2013; de Rooy 64 et al., 2013; Romps, 2016; Tan et al., 2018; Smalley et al., 2019; Couvreux et al., 2021; 65 Hourdin et al., 2021). Similarly, limited-area high-resolution simulations (e.g., Wang et 66 al., 1996; Fox-Kemper & Menemenlis, 2013) have been used to calibrate subgrid-scale 67 parameterizations of upper-ocean turbulence (Souza et al., 2020; Li & Fox-Kemper, 2017; 68 Van Roekel et al., 2018; Reichl et al., 2016; Li et al., 2019; Campin et al., 2011; Reichl 69 & Hallberg, 2018). 70

71 However, high-resolution simulations can be used more systematically, by driving them with a large-scale GCM and using the mismatch between statistics of the high-resolution 72 simulations and GCM output to calibrate parameterizations (Schneider, Lan, et al., 2017). 73 For example, Shen et al. (2020, 2021) drive atmospheric LES in domains  $\mathcal{O}(10 \text{ km})$  in 74 the horizontal with output from climate models, to produce simulated local data with 75 which climate model parameterizations can be calibrated. In principle, LES can be em-76 bedded into coarse-resolution GCMs to sample thousands of sites across the globe. Pa-77 rameterizations in the GCM can then learn automatically from the high-resolution sim-78 ulations, and the high-resolution simulations can even be spun off on-the-fly during the 79 integration of a coarse-resolution model (Schneider, Lan, et al., 2017). This allows sys-80 tematic calibration and uncertainty quantification of parameterizations. 81

Pursuing such automatic calibration and uncertainty quantification of parameter-82 izations raises the question of where to embed the high-resolution simulations so that 83 they are maximally informative about the parameterizations. This is a question of ex-84 perimental design (Chaloner & Verdinelli, 1995), akin to the question of how to choose 85 sites for supplementary weather observations to optimally improve weather forecasts (Lorenz 86 & Emanuel, 1998; Bishop & Toth, 1999; Emanuel et al., 1995). When climate statistics 87 are non-stationary, there is a related question of what time periods should the expen-88 sive high-resolution simulation be integrated over. 89

Here we present a mathematical framework for addressing both the experimental 90 design question of where, and when, to embed high-resolution simulations in a coarser 91 model. We adopt a Bayesian inverse problem setting (see, e.g., Kaipio and Somersalo 92 (2006), Tarantola (2005), Stuart (2010), and Dashti and Stuart (2013) for reviews). In 93 this setting, parameters (or parametric or nonparametric functions) in parameterizations 94 are treated as having probability distributions. Data (e.g., from high-resolution simu-95 lations) are used to reduce the uncertainty reflected by these distributions, balancing con-96 tributions of the data with that of prior knowledge about parameters (e.g., physical con-97 straints). This setting is well suited for our needs, as it provides the complete joint pos-98 terior distribution for parameters, including the correlation structure of uncertainties among 99 parameters. Distribution information is beneficial, for example, because it enables model-100 based predictions of rare events with quantified uncertainties (Dunbar et al., 2021). Anal-101 ysis of the posterior distribution also may focus scientific development (e.g., improve-102 ment of parameterization schemes, Souza et al. (2020)) on areas where uncertainties can 103 most effectively be minimized. In this paper, we use the posterior distribution to deter-104 mine regions and times where local data (e.g., from high-resolution simulations) are max-105 imally effective at reducing parameter uncertainties. 106

Construction of the full posterior distribution of the parameters is well known to 107 be a computationally intensive task, requiring  $\mathcal{O}(10^5)$  evaluations of the model in which 108 the parameters appear with commonly used Markov chain Monte Carlo (MCMC) meth-109 ods (see Geyer (2011) for an overview). The recent development of the calibrate-emulate-110 sample (CES) framework accelerates Bayesian learning by a factor of  $10^3$  (Cleary et al., 111 2021a; Dunbar et al., 2021). The calibration stage uses a variant of ensemble Kalman 112 inversion (Iglesias et al., 2013) building on (Chen & Oliver, 2012; Emerick & Reynolds, 113 2013; Reich, 2011) to obtain a collection of samples of the model about an optimal set 114 of parameters. The emulation stage features the training of a Gaussian process (Williams 115

<sup>116</sup> & Rasmussen, 2006; Kennedy & O'Hagan, 2000, 2001a) to emulate the model using the <sup>117</sup> samples from the calibration stage. The sample stage then samples a posterior distri-<sup>118</sup> bution with MCMC methods, replacing the computationally expensive model with the <sup>119</sup> cheap emulator. We build on the CES framework and show how Bayesian experimen-<sup>120</sup> tal design approaches can be incorporated within CES at negligible additional compu-<sup>121</sup> tational expense. In particular, we do not require additional forward model evaluations <sup>122</sup> over what is already required in CES to perform uncertainty quantification.

To target high-resolution simulations, we use tools from experimental design, which 123 provides methods for assessing the efficacy of learning about parameters from different 124 designs (e.g., data from different locations or time periods) (Ryan et al., 2016). We de-125 termine the optimal designs where the model is most sensitive to the parameters by us-126 ing the posterior distribution for the parameters as a utility to be optimized. We choose 127 a utility function that assigns a score of the information entropy loss between posterior 128 and prior for each design (Chaloner & Verdinelli, 1995; Alexanderian et al., 2014; Alexan-129 derian & Saibaba, 2018); it is scale invariant and scalable to problems with high-dimensional 130 input parameter spaces. The region with maximal utility determines where to acquire 131 high-resolution data, and hence where to divert scarce computational resources for max-132 imal effect. 133

We demonstrate the effectiveness of this targeted learning approach through simulations with an idealized moist GCM, with which we generate surrogates for high-resolution data and in which we are calibrating parameters in a quasi-equilibrium convection scheme (D. M. Frierson, 2007; O'Gorman & Schneider, 2008). We showcase our algorithms by showing that the recovered posterior distributions are reflective of the utility of information at different points and different times along the annual cycle.

In Section 2, we define the inverse problems for parameter calibration and the optimal design algorithm; details of efficient uncertainty quantification (CES) are left to Appendix A. In Section 3, we briefly describe the GCM used for demonstrating the algorithm. Results of the optimal design algorithm are described in section 4, first in a setting in which the GCM statistics are statistically stationary, then with seasonally varying GCM statistics. We end with discussion and future directions in Section 5.

# $_{146}$ 2 Methodology

Our goal is to target data acquisition to regions and times at which uncertainty reduction (information gain) is maximized. We do so by first learning the temporally and spatially varying sensitivities of the model statistics with respect to model parameters. We then use this knowledge to target data acquisition to regions and times at which the model is maximally sensitive to new data. We work in a framework similar to Dunbar et al. (2021), which focuses on accelerated uncertainty quantification within a GCM.

Our point of departure in Section 2.1 is to specify the inverse problem for uncer-153 tainty quantification of parameters from data at a specific design. Related to applica-154 tion, this can be seen as the stage of learning parameter uncertainties from high-resolution 155 simulation data at a certain region or time. Treating such sdata as computationally ex-156 pensive to obtain, in Section 2.2 we investigate how to efficiently choose which region 157 or time to learn from. To do this we formulate a set of related inverse problems, whose 158 solutions allow us to assess the quality of different choices. In Section 2.3 we connect these 159 two stages to form the targeted uncertainty quantification algorithm . 160

# <sup>161</sup> 2.1 Inverse problem

We study calibration of parameters in a GCM by formulating parameter learning as a Bayesian inverse problem. Define  $\mathcal{G}_T(\boldsymbol{\theta}; \boldsymbol{v}^{(0)})$  to be the forward map sending the pa-

rameters  $\boldsymbol{\theta}$  to time-aggregated simulated climate statistics (averaged over a window of 164 length T > 0 from an initial state  $v^{(0)}$ . We assume that the aggregation  $\mathcal{G}_T(\boldsymbol{\theta}, \cdot)$  is sta-165 tistically stationary, and samples of such aggregated climate statistics are referred to as 166 data throughout. We consider a situation in which data are only locally available, at a 167 particular spatial or spatio-temporal location, indexed by k, which we refer to as the de-168 sign point. This is relevant to our application of targeted high-resolution simulations with 169 limited spatial and temporal extent. We make use of a restriction operation  $W_k$  to a point 170 k, and define the local forward map,  $S_T(\boldsymbol{\theta}; k, \boldsymbol{v}^{(0)}) = W_k \mathcal{G}_T(\boldsymbol{\theta}; \boldsymbol{v}^{(0)}).$ 171

For any given k, assume we have local data  $z_k$  available. In the application of interest,  $z_k$  is produced by a high-resolution simulation run. We can construct the forward map  $S_T(\cdot; k, \cdot)$ , to form an inverse problem for the GCM learning from the local data as

$$\boldsymbol{z}_k = \mathcal{S}_T(\boldsymbol{\theta}; k, \boldsymbol{v}^{(0)}) + \delta_k, \tag{1}$$

where  $\delta_k$  is a stochastic term to capture discrepancies between model  $S_T(\cdot; k, \cdot)$  and data  $\boldsymbol{z}_k$ , (e.g., Kennedy & O'Hagan, 2001a). The initial condition  $\boldsymbol{v}^{(0)}$  appears in this formulation but is treated a nuisance variable. This view is justified in the context of learning about atmospheric parameterizations for climate models, where lower frequency information is informative (Schneider, Lan, et al., 2017). We use time-averaged data to filter out the high frequency information, and take T is larger than the dynamical system's Lyapunov timescale (for the atmosphere, this equates to  $T \gtrsim 15$  days (F. Zhang et al., 2019)). To deal with the initial condition, one can view finite-time averaging as a perturbation of an infinite-time average by means of a central limit theorem. Following (Dunbar et al., 2021), we write  $S_T(\boldsymbol{\theta}; k, \boldsymbol{v}^{(0)}) \approx S_{\infty}(\boldsymbol{\theta}; k) + \sigma_k$ , where  $\sigma_k \sim N(0, \Sigma(\boldsymbol{\theta}))$  is normal noise, independent from  $\delta_k$ , with mean zero and with a covariance matrix  $\Sigma(\boldsymbol{\theta})$  reflecting chaotic internal variability. The inverse problem then becomes

$$\boldsymbol{z}_{k} = \mathcal{S}_{\infty}(\boldsymbol{\theta}; k) + \delta_{k} + \sigma_{k}, \qquad \sigma \sim N(0, W_{k} \Sigma(\boldsymbol{\theta}) W_{k}^{T}).$$
<sup>(2)</sup>

This is now a desirable form of the inverse problem since the dependence on the initial condition has been removed.

<sup>174</sup> Solving (2) involves finding the posterior distribution of  $\boldsymbol{\theta}$  given the data  $\boldsymbol{z}_k$ , de-<sup>175</sup> noted ( $\boldsymbol{\theta} \mid \boldsymbol{z}_k$ ). Although we cannot evaluate  $\mathcal{S}_{\infty}$  directly, we use the emulate phase of <sup>176</sup> the calibrate-emulate-sample (CES) algorithm (Cleary et al., 2021b) to construct a sur-<sup>177</sup> rogate of  $\mathcal{S}_{\infty}$  from carefully chosen evaluations of  $\mathcal{S}_T$ . This has been shown to be effi-<sup>178</sup> cient with respect to the required number of evaluations of  $\mathcal{S}_T$  (Cleary et al., 2021b; Dun-<sup>179</sup> bar et al., 2021). Details of this algorithm are provided in Appendix A.

180

## 2.2 Experimental design

We imagine a situation where evaluating  $\boldsymbol{z}_k$  has a large computational cost. In the 181 relevant application of targeted high-resolution simulations,  $\boldsymbol{z}_k$  is data obtained by run-182 ning a high-resolution simulation at design point k. Our starting point is to assume that 183 a limited computational budget restricts us to evaluate  $\boldsymbol{z}_k$  at a single design point k at 184 a time. We want to choose the k that leads to the most informative inverse problem (2). 185 We take a Bayesian point of view, namely, the optimal k is the one for which the pos-186 terior distribution of  $(\boldsymbol{\theta} \mid \boldsymbol{z}_k)$  learned from the inverse problem (2) has the smallest un-187 certainty. 188

To answer this conclusively, one would need to evaluate  $z_k$  at all design points k, which here is too computationally expensive. Instead, we investigate only the sensitivity of the forward model statistics  $\mathcal{G}_T$  to its parameters  $\boldsymbol{\theta}$  to assess the marginal information provided at each design point k. This marginal information at k is used as a proxy for the information content that would exist when learning from data  $z_k$ . The benefits of this approach are that (i) we do not require any evaluations of  $z_k$  to select the optimal location; (ii) the measure of information content is naturally constructed from the

uncertainty reflected by the Bayesian posterior distribution; and (iii) we can perform this 196 efficiently, and in a embarrassingly parallel fashion, requiring only O(100) GCM runs, 197 determined by the product of the ensemble size and the number of iterations typically 198 needed in the calibration stage of the CES algorithm (see Appendix A). The approach 199 necessarily will contain a bias from the prior distribution of the parameters, and it im-200 plicitly assumes unbiased model statistics  $\mathcal{G}_T$ . The latter in practice requires the inclu-201 sion of models for structural model error within  $\mathcal{G}_T$ , for example, learned error models 202 that enforce conservation laws and sparsity (M. E. Levine & Stuart, 2021; Schneider et 203 al., 2021). 204

Each evaluation of the forward map involves a simulation with the GCM and thus depends on an initial condition  $v^{(0)}$  and parameters  $\theta$ . Together this gives rise to the definition of time-aggregated model statistics y,

$$\boldsymbol{y} = \mathcal{G}_T(\boldsymbol{\theta}; \boldsymbol{v}^{(0)}). \tag{3}$$

Using the central limit theorem as before, we may write this relationship as

$$\boldsymbol{y} = \mathcal{G}_{\infty}(\boldsymbol{\theta}) + \sigma, \qquad \sigma \sim N(0, \Sigma(\boldsymbol{\theta})),$$
(4)

where  $\Sigma(\boldsymbol{\theta})$  is the internal variability covariance matrix for parameters  $\boldsymbol{\theta}$ . To proceed, we choose a control value  $\boldsymbol{\theta}^*$ , for example we take the mean of the prior distribution, and, fixing  $\boldsymbol{\theta} = \boldsymbol{\theta}^*$ , we generate a realization of  $\boldsymbol{y}$ . Given this realization of  $\boldsymbol{y}$ , we temporarily forget  $\boldsymbol{\theta}^*$ , and for any design point k, we consider a restriction of an inverse problem to k

$$W_k \boldsymbol{y} = W_k \mathcal{G}_{\infty}(\boldsymbol{\theta}) + \sigma_k, \qquad \sigma_k \sim N(0, W_k \Sigma(\boldsymbol{\theta}) W_k^T). \tag{5}$$

The posterior distributions of  $\boldsymbol{\theta} \mid W_k \boldsymbol{y}$  for all k obtained by solving (5) informs us about the sensitivities of  $\mathcal{G}_{\infty}$  with respect to parameters, when only data at different k is available. To simplify solution of the inverse problem, we approximate the internal variability covariance matrix  $\Sigma(\boldsymbol{\theta})$  by a fixed covariance matrix  $\Sigma(\boldsymbol{\theta}^*)$ . This covariance matrix can be obtained by running a collection of control simulations with parameters fixed to (the known)  $\boldsymbol{\theta}^*$  but with different initial conditions.

The utility U of a design  $W_k$  is a scalar function reflecting the quality of a given design. The design that maximizes the utility function is known as the optimal design. We choose a utility function by measuring information gain (or uncertainty reduction) in  $(\boldsymbol{\theta} \mid W_k \boldsymbol{y})$  relative to the prior, in a form of Bayesian optimal design. We use the utility function arising from the linear Bayesian design (Chaloner & Verdinelli, 1995), which is the determinant of the information matrix (inverse posterior covariance matrix),

$$U(W_k) = \left(\det\left(\operatorname{Cov}(\boldsymbol{\theta} \mid W_k \boldsymbol{y})\right)\right)^{-1}.$$
(6)

The posterior covariance matrix  $\operatorname{Cov}(\boldsymbol{\theta} \mid W_k \boldsymbol{y})$  can be estimated as the empirical co-211 variance matrix of samples drawn from the posterior distribution  $(\boldsymbol{\theta} \mid W_k \boldsymbol{y})$  for a de-212 sign  $W_k$ . We refer to (6) as the D-utility because it fulfills the so-called D-optimality 213 criterion. It is invariant under arbitrary linear transformations of the parameters, for ex-214 ample, when parameters are on different dimensional scales, unlike trace-based measures 215 (e.g., A-optimal utility functions). For linear forward maps and Gaussian priors, max-216 imization of this D-utility is equivalent to maximization of the expected Kullback-Leibler 217 divergence (KLD), a relative entropy measure (Ryan et al., 2014; Huan & Marzouk, 2013; 218 Cook et al., 2008; Kim et al., 2014). While KLD has beneficial mathematical properties, 219 especially for highly non-Gaussian posteriors (Paninski, 2005), in practice it is difficult 220 to evaluate, especially in high-dimensional problems (e.g., Huan & Marzouk, 2013). 221

222

# 2.3 Targeted uncertainty quantification algorithm

The combined algorithm for targeted uncertainty quantification consists of two stages: first, finding an optimal design point  $\tilde{k}$  in a design stage and, second, evaluating parameter uncertainty with data from  $\tilde{k}$  in an uncertainty quantification stage. Let D be the finite index set for the set of design points, and define  $W_k$  to be the restriction map for any  $k \in D$ . The two stages then are as follows:

1. The design stage consists of the following steps:

229

230

231

232

234

235

243

(a) Generate a sample of GCM simulated data  $\boldsymbol{y} = \mathcal{G}_T(\boldsymbol{\theta}^*; \boldsymbol{v}^{(0)})$ , and estimate the internal variability covariance matrix  $\Sigma(\boldsymbol{\theta}^*)$ . We approximate  $\Sigma(\boldsymbol{\theta})$  as  $\Sigma(\boldsymbol{\theta}^*)$ .

(b) For each  $k \in D$ , solve (5), in parallel, for the posterior of  $(\boldsymbol{\theta} \mid W_k \boldsymbol{y})$ , using the CES-type algorithm described in Appendix A.

(c) For each  $k \in D$ , calculate the *D*-utility  $U(W_k)$  from (6) and choose

$$\tilde{k} = \arg \max_{k \in D} U(W_k)$$

233 2. The uncertainty quantification stage consists of the following steps:

- (a) At the optimal design point  $\tilde{k}$ , obtain a sample  $\boldsymbol{z}_{\tilde{k}}$ .
  - (b) Solve the inverse problem (2) for the posterior distribution of  $(\boldsymbol{\theta} \mid \boldsymbol{z}_{\tilde{k}})$ .

The complexity of the first stage grows linearly with the candidate design points k because we only consider a point at a time. However, if one wishes to choose a design composed of K simultaneous points from a set D, a combinatorial problem arises, with complexity growing like |D|!/((|D| - |K|)!|K|!). This will become prohibitively costly to solve by brute force, even in parallel. We focus on the algorithm for single design points k for now, addressing scaling questions in the discussion section.

# <sup>242</sup> 3 Idealized GCM and Experimental Setup

## 3.1 Idealized GCM, parameters, and priors

To demonstrate the algorithm in a simplified setting, we use the idealized aqua-244 planet GCM described by D. M. W. Frierson et al. (2006) and O'Gorman and Schnei-245 der (2008b). The aquaplanet is a climate model with atmosphere and a simplified slab 246 ocean covering the entire planet surface. Without topography, it exhibits symmetries in 247 the longitudinal directions. The aquaplanet can produce statistically stationary climates 248 by prescribing fixed insolation. It can also cyclostationary statistics over seasons through 249 seasonally varying insolation (Bordoni & Schneider, 2008a; Howland et al., 2021). It has 250 been shown in Dunbar et al. (2021); Howland et al. (2021) that the parameters of a sim-251 ple quasi-equilibrium moist convection parameterization can be calibrated within this 252 GCM in the stationary and cyclostationary regimes. The quasi-equilibrium moist con-253 vection scheme relaxes temperature and specific humidity toward moist-adiabatic ref-254 erence profiles with a fixed relative humidity RH (D. M. W. Frierson, 2007). The timescale 255 with which the temperature and specific humidity relax to their respective reference pro-256 files is given by the parameter  $\tau$ . The two parameters RH and  $\tau$  are the key parameters 257 to be calibrated and whose uncertainties we want to determine and minimize. 258

The priors for these parameters are taken to be logit-normal and lognormal distributions, RH ~ Logitnormal(0, 1) and  $\tau$  ~ Lognormal(12 h, (12 h)<sup>2</sup>). That is, we define the invertible transformation

$$\mathcal{T}(\mathrm{RH}, \tau) = \left(\mathrm{logit}(\mathrm{RH}), \ \ln\left(\frac{\tau}{1\ \mathrm{s}}\right)\right),$$

which transforms each parameter to values along the real axis. We label the transformed (or computational) parameters as  $\boldsymbol{\theta} = \mathcal{T}(\mathrm{RH}, \tau)$ , and the untransformed (or physical) parameters (relative humidity and timescale) are uniquely defined by  $\mathcal{T}^{-1}(\boldsymbol{\theta})$ . We apply our calibration methods in the space of the transformed parameters  $\boldsymbol{\theta}$ , where priors are unit-free, normally distributed, and unbounded; meanwhile, the climate model uses the physical parameters  $\mathcal{T}^{-1}(\boldsymbol{\theta})$ , with  $\mathrm{RH} \in [0, 1]$  and  $\tau \in [0, \infty)$ . In this way, the prior distributions enforce physical constraints on the parameters.

# 3.2 Objective function for parameter learning

266

We learn from climate statistics that are known to be sensitive to the parameters. 267 We have knowledge about these sensitivities from a body of previous studies of large-268 scale atmosphere dynamics and mechanisms of climate changes which used this ideal-269 ized GCM (e.g., O'Gorman & Schneider, 2008b, 2008a; Bordoni & Schneider, 2008b; O'Gorman 270 & Schneider, 2009b; Schneider et al., 2010; Merlis & Schneider, 2011; O'Gorman, 2011; 271 Kaspi & Schneider, 2011, 2013; X. Levine & Schneider, 2015; Bischoff & Schneider, 2014; 272 Wills et al., 2017; Wei & Bordoni, 2018). We know, for example, that the convection scheme 273 274 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 275 humidity parameter RH in the convection scheme controls the humidity of the tropical 276 free troposphere but has a weaker effect on the humidity of the extratropical free tro-277 posphere (O'Gorman et al., 2011). Thus, we expect tropical circulation statistics to be 278 especially informative about the parameters in the convection scheme. However, convec-279 tion plays a central role in intense precipitation events at all latitudes (O'Gorman & Schnei-280 der, 2009b, 2009a), so we expect statistics of precipitation intensity to be informative 281 about convective parameters, and in particular to contain information about the relax-282 ation timescale  $\tau$ . 283

As statistics to learn from, we therefore choose averages of the free-tropospheric 284 relative humidity, of the precipitation rate, and of a measure of the frequency of intense 285 precipitation. We use averages over T = 30 days in statistically stationary simulations 286 (Dunbar et al., 2021) and over T = 90 days in simulations of the seasonal cycle (Howland 287 et al., 2021). We exploit the symmetry in the GCM by taking zonal averages in addi-288 tion to the time averages. The relative humidity data are evaluated at  $\sigma = 0.5$  (where 289  $\sigma = p/p_s$  is pressure p normalized by the local surface pressure  $p_s$ ), the precipitation 290 rate is taken daily, and as a measure of the frequency of intense precipitation, we use the 291 frequency with which daily precipitation exceeds the latitude-dependent 90th percentile 292 of precipitation rates in a long (18000 days) control simulation. We run the GCM at the 293 coarse horizontal spectral resolution of T21, implying 32 discrete latitudes on the spec-294 tral transform grid. Hence, we have 3 statistics, each a function of 32 latitude points, 295 resulting in a 96-dimensional processed output, defined as  $\mathcal{H}_T$ . In the statistically sta-296 tionary case, we take the forward map  $\mathcal{G}_T = \mathcal{H}_T$ . 297

For the simulations with a seasonal cycle,  $\mathcal{H}_T$  is not statistically stationary but is cyclostationary over multiples of a year. The year length in the GCM is 360 days. We stack four 90-day seasons of data together (Howland et al., 2021) and define the forward map

$$\mathcal{G}_T(\boldsymbol{\theta}; \boldsymbol{v}^{(0)}) = [\mathcal{H}_T(\boldsymbol{\theta}; \boldsymbol{v}^{(0)}), \dots, \mathcal{H}_T(\boldsymbol{\theta}; \boldsymbol{v}^{(3)})],$$

over a one year cycle (360 days), where  $v^{(i)}$  is the model state at the beginning of each 90-day long season labelled i = 0, 1, 2, 3. With this batching, we have now constructed stationary statistics for the stacked data. The theory of Section 2 applies, and our inverse problems can be formulated in the seasonally varying case.

#### 302 3.3 Design choices

In the stationary GCM setting, we aggregate statistics temporally and zonally. Thus, a local design implies a restriction to certain latitudes. Recall our discretization has 32 discrete latitudes. We therefore choose a design space that contains sets of  $\ell$  consecutive discrete latitudes, indexed from south to north poles with the design points k = 1, ..., 32 - $(\ell-1)$ . In the stationary experiments, we choose  $\ell = 3$ , indexing designs k = 1, ..., 30, unless otherwise specified. The choice of  $\ell$  is discussed in Section 4.1.

In the seasonal GCM setting, we still aggregate in time and longitudinal directions, but we also stack the seasons in a vector. We define a local design by indexing both a restriction to a season and a restriction to certain latitudes. We choose a design space that contains sets of  $\ell$  consecutive discrete latitudes, collected season by season, indexed from south to north poles as  $1, \ldots, 32 - (\ell - 1)$ , and from spring to winter as  $0, \ldots, 3$ , all collected as k = (season, latitude). In the seasonal experiments, we choose  $\ell = 1$ , which indexes the designs  $k = (0, 1), \ldots, (3, 32)$ .

316

#### 3.4 Synthetic data and noise

To generate surrogates of locally available data from high-resolution simulations, we generate data with the idealized GCM itself at a fixed parameter vector  $\boldsymbol{\theta}^{\dagger}$ , adding Gaussian noise  $\delta$  with zero mean and covariance matrix  $\Delta$  as in (2). The implication is that we generate  $\boldsymbol{z}_k$  with the restricted idealized GCM  $\mathcal{S}_T(\boldsymbol{\theta}^{\dagger};k)$ , corrupted by noise to describe model error (Kennedy & O'Hagan, 2001b). In this way, the inverse problem (2) can be written in the compact form

$$\boldsymbol{z}_k = \mathcal{S}_{\infty}(\boldsymbol{\theta}; k) + \gamma_k, \qquad \gamma_k \sim N(0, W_k(\Sigma(\boldsymbol{\theta}) + \Delta)W_k^T).$$
 (7)

We construct the measurement error covariance matrix  $\Delta$  to be diagonal with entries  $d_i^2 = \Delta_{ii} > 0$ , where *i* indexes over data type (three observed quantities) and over the number of discrete latitudes,

$$\Sigma + \operatorname{diag}(d_i^2) = \Sigma + \Delta. \tag{8}$$

We choose  $d_i$  so that it is proportional to the mean  $\mu_i$  of the variable in question, with a proportionality factor  $C_{\text{max}} = 0.1$ . To prevent the noise from becoming so large that the variables can cross a physical boundary  $\partial \Omega_i$  (e.g., relative humidity becoming negative), we limit the noise standard deviation to a factor C = 0.2 times the distance between the approximate 95% noise confidence interval and the physical boundary:

$$d_{i} = \min\left(C\min\left(\operatorname{dist}(\mu_{i} + 2\sqrt{\Sigma_{ii}}, \partial\Omega_{i}), \operatorname{dist}(\mu_{i} - 2\sqrt{\Sigma_{ii}}, \partial\Omega_{i})\right), C_{\max}\mu_{i}\right)$$

We carry out a set of control simulations, with the parameters fixed to standard 317 values  $\boldsymbol{\theta}^{\dagger}$ , where  $\mathcal{T}^{-1}(\boldsymbol{\theta}^{\dagger}) = (0.7, 2 \text{ h})$  are standard values used in previous studies (O'Gorman 318 & Schneider, 2008b). We use this set of control simulations to estimate the restricted 319 covariance matrix  $W_k \Sigma(\boldsymbol{\theta}) W_k^T \approx W_k \Sigma(\boldsymbol{\theta}^{\dagger}) W_k^T$  for performing uncertainty quantifica-320 tion with local data  $z_k$  (stage 2 in Section 2.3). In the statistically stationary case, we 321 carry out control simulations for 650 windows of length T = 30 days, discarding the 322 first 50 months for spin-up, and calculate the sample covariance matrix  $\Sigma(\theta^{\dagger})$  from the 323 latter 600 samples. Here,  $W_k \Sigma(\boldsymbol{\theta}^{\dagger}) W_k^T$  is a symmetric matrix whose size depends on the 324 design space; it represents noise from internal variability in 30-day time averages. In the 325 seasonally varying case, we carry out a control simulation for 150 years, discarding the 326 first 4 years for spin-up, and obtain the sample covariance matrix  $\Sigma(\theta^{\dagger})$  from the stacked 327 seasonal (T = 90 days) averages. In the seasonal case, it is a symmetric matrix whose 328 size depends on 4 times the design space and represents noise from internal variability 329 in the 90-day time averages. In practical implementations of this methodology, good es-330 timates of the local variability that we represent with  $W_k \Sigma(\boldsymbol{\theta}^{\dagger}) W_k^T$  can be made from 331 the observed climatology of the statistics of interest. 332

For the design stage (stage 1 in Section 2.3) we estimate  $\Sigma(\theta^*)$  from a second set 333 of control simulations of the GCM in which we fix the parameters to the prior mean  $\theta^*$ . 334 equivalent to the physical values  $\mathcal{T}^{-1}(\boldsymbol{\theta}^*) = (0.5, 12 \text{ h})$ . In the stationary case, the 3 335 latitude-dependent fields evaluated at 32 latitude points produce a  $96 \times 96$  symmetric 336 matrix  $\Sigma(\boldsymbol{\theta}^*)$ , representing noise from internal variability in 30-day averages; in the sea-337 sonal case, the stacked statistics produce a  $384 \times 384$  symmetric matrix  $\Sigma(\theta^*)$ , repre-338 senting noise from variability of 90-day averages. In either case, we take  $\Sigma(\boldsymbol{\theta}) = \Sigma(\boldsymbol{\theta}^*)$ 339 in the optimal design stage of the algorithm. 340

The mean and 95% confidence interval of the data at  $\theta^*$ , with covariance constructed from  $\Sigma(\theta^*)$ , are shown in Figure 1 for the statistically stationary case and in Figure 2



Figure 1. Aggregated climate statistics in the statistically stationary control simulation, with parameters set to the mean of the prior  $\theta^*$ . The mean (grey lines) and 95% confidence intervals (shading) of the data are plotted against latitude. One realization of the data is shown (black line). No noise is added here.



Figure 2. Aggregated climate statistics in the seasonally varying control simulation, with parameters set to the mean of the prior  $\theta^*$ . The mean (solid lines) and 95% confidence intervals (shading) of the data are plotted against latitude, with the colors indicating the seasons. No noise is added here.

for the seasonally varying case. The black (stationary) and colored (seasonal) solid lines illustrate a realization of the data for one initial condition. Similarly, the mean and 95% confidence interval of the data at  $\theta^{\dagger}$ , with noise added with covariance matrix  $\Delta + \Sigma(\theta^{\dagger})$ (over all designs for illustration), are shown in Figure 3 for the stationary and in Figure 4 for the seasonal case.

#### 348 4 Results

349

#### 4.1 Stationary statistics

We first apply the optimal design algorithm to the statistically stationary GCM. The logarithm of the utility function is shown in Figure 5, with four representative samples shown by the colored discs (specifically, these are the design points k = 15, 14, 20,and 3, in decreasing order of utility). The extent to which hemispheric symmetry of the statistics is broken in Figure 5 is an indication of sampling variability, as the infinite-time GCM statistics are hemispherically symmetric.

The distribution of the inflated climate statistics produced at the true parameters  $\theta^{\dagger}$  are represented by the mean and 95% confidence interval in grey in Figure 3, which also shows the data samples for each three-latitude design stencil as colored discs for four representative design locations. We apply the uncertainty quantification stage in Sec-



**Figure 3.** Aggregated climate statistics in the statistically stationary control simulation using the ground truth parameters. Mean (grey lines) and 95% confidence intervals (shading) of the data are plotted against latitude. Additional inflation noise is added. Each set of colored discs represents a 30-day realization of inflated GCM data coming from a different 3-latitude design used in the experiment.



**Figure 4.** Aggregated climate statistics produced from the seasonally varying control simulation using the ground truth parameters, and with additional inflation. The mean (solid lines) and 95% confidence intervals (shading) of the data are plotted against latitude, with the colors indicating the seasons. The blue vertical line indicates the location and season (northern winter) in which we observe the data for uncertainty quantification; the specific 90-day realization of inflated GCM data for the 1-latitude design is given by the blue disc.



Figure 5. Logarithm of the data utility as a function of latitude, with designs represented by a node at the center of each stencil (comprised of three neighboring latitudes). The colored discs signify the four representative designs indicated in Fig. 3, which are used in the uncertainty quantification experiment.



**Figure 6.** Posterior distributions for convection parameters learned from data restricted to different design points. The drawn contours bound 50%, 75% and 99% of the distribution. Panels (a, b, c, d) correspond to designs k = (51, 14, 20, 3), ordered to express learning from data at decreasingly informative design points (i.e., points of decreasing utility). The true parameter values in the control simulation are given by the blue circle. The parameters found to be optimal in the calibration scheme (given a single random realization of data) are given by the red star in each case (in panel (d) this is outside the plotting region).



**Figure 7.** Performance for different optimal design selections at smaller stencil sizes. The contours bound 50%, 75%, and 99% of the distribution. Panels (a,c) is a two-latitude stencil, and panels (b,d) is a one-latitude stencil. The top row displays the logarithm of the utility plot, and the bottom row the corresponding posterior from a sample at the optimal latitude, marked by a disc at the top.

tion 2.3 to each location, resulting in the posterior distributions shown in Figure 6. Each 360 panel shows the density contours bounding 50%, 75%, and 99% of the posterior distri-361 bution, shaded dark to light; the priors are largely uninformative and have been excluded 362 from the plots. The panels are ordered (a - d) by decreasing D-utility, a predictor of in-363 formation content based on uncertainty at the prior mean parameter  $\theta^*$ . We see this mono-364 tonicity is preserved when considering data produced from the true parameter  $\theta^{\dagger}$  in this 365 example. In particular, this implies that the design optimizing the chosen utility pro-366 duces minimal uncertainty in the uncertainty quantification stage. As observed in other 367 investigations (Dunbar et al., 2021), the posterior distributions are subject to variabil-368 ity due to the finite-time sampling and the inflation. However, all distributions capture 369 the true parameter values within 99% of the posterior mass. 370

For the statistically stationary case, we investigate the choice of  $\ell$ , a measure of the 371 design sparsity. To this end we repeat the experiment, choosing  $\ell = 2$  or 1 in Section 3.3). 372 For each, Figure 7 shows the utility function against the latitude at the center of the sten-373 cil and the posterior distribution at the respective optimal designs. We see that in both 374 cases, the optimal design remains robust, coinciding with the three-stencil case. Peak 375 utility is consistently at a design near the equator. The posterior distributions are seen 376 to be far broader than in the three-latitude case, offering only marginal improvement over 377 the prior distribution in the one-latitude stencil case. They are non-Gaussian and mul-378 timodal but nevertheless capture the true parameters (blue disc) with high probability. 379 They provide insight into the correlation structure between the parameters at the op-380 timal design location. We observe that for these sparser designs, non-identifiability (mul-381



**Figure 8.** Logarithm of the data utility plotted against latitude (1 design per latitude and season). The blue disc signifies that a latitude in northern winter maximizes the utility function.

timodality) appears only at data from  $\theta^{\dagger}$ , but not at  $\theta^*$ . As a result, the optimal uncertainty is not guaranteed to be found at the location of optimal utility. This is remedied by having a better initial guess through the prior, or (as demonstrated in the threelatitude set) a less sparse data set from which the parameters are more identifiable.

# 386 4.2 Seasonally varying statistics

In the seasonally varying case, we choose the optimal design with the algorithm 387 in Section 2.3 applied to the stacked data. Figure 8 shows the logarithm of the utility 388 function. Hemispheric and seasonal asymmetries are evident here. In northern winter, 389 latitudes just south of the the equator (k = (3, 16)) optimize the design, in the vicin-390 ity of the ITCZ. Conversely, in northern summer, latitudes just north of the equator (k =391 (1,17)) optimize the design, again in the vicinity of the seasonally migrating ITCZ; ad-392 ditional peaks at around 30 degrees can be seen. The equinox seasons have less utility 393 at the optimal designs (k = (0, 17) and k = (2, 16)). Because the equinoctial Hadley 394 cells and ascent regions in the ITCZ are less pronounced than the solstitial Hadley cells 395 (Schneider et al., 2010), utility is more spread out across the latitudes. 396

We solve the analogue inverse problem (7) as in the nonseasonal case with a data 397 sample at the optimal spatial design location for each season. The posterior distributions 398 are collected in Figure 9, colored by season. In general, the true parameter values lie in 399 regions of high posterior density in each case. We see qualitatively that the utility of the 400 different designs predicts the size of support of the corresponding posterior distri-401 bution, in particular the design with highest utility (northern winter) also has the small-402 est support. This indicates that the utility is still a good predictor of data quality for 403 learning the convection parameters in the cyclo-stationary settings. 404

# **5** Conclusions and Discussion

We have presented a novel framework for automated optimal placement of highresolution simulations embedded in lower-resolution models. The framework can be used with computationally expensive and chaotic (noisy) low-resolution models, whose derivatives may not be available. Given low-resolution simulations, we use parameter uncertainty information provided by the CES algorithm to guide our choice of design. We have demonstrated the efficacy of the algorithm for finding optimally informative locations in perfect-model settings in which we generated surrogates of embedded high-resolution



Figure 9. Posterior distribution obtained from using data at the optimal latitude from each season. Contours bound 50%, 75% and 99% of the distribution (in decreasing color saturation). Panels (a, b, c, d) correspond to designs k = ((0, 17), (1, 17), (2, 16), (3, 16)), ordered by season. The true parameter values in the control simulation are given by the blue circle. The parameters found to be optimal in the calibration scheme (given a single random realization of data) are given by the red star in each case.

simulations directly with an idealized GCM, with statistically stationary or seasonally
varying statistics. In these settings, we have demonstrated learning about parameters
in a convection parameterization, exploring both spatial and spatio-temporal designs.
Our design framework can also be used more broadly, to automate selection of optimally
informative climate statistics from libraries of high-resolution simulations (Shen et al.,
2021).

With the idealized GCM, we showed how to optimally target a location at which 419 additional data will produce parameter estimates that minimize uncertainty. In our proof-420 421 of-concept in which we calibrated parameters in a simple convection scheme, the automatically targeted optimal location for new data was consistently near the equator, in 422 the vicinity of the seasonally migrating ITCZ. This is consistent with the fact that the 423 convection scheme in the idealized GCM is most important near the ITCZ (O'Gorman 424 & Schneider, 2008). We showed that the optimal targeting is limited in its effectiveness 425 in settings of very sparse data, where parameter posteriors can be multimodal. However, 426 with access to the posterior distributions of the parameters, the behavior is both diag-427 nosable a posteriori, and actionable with successive iterations of the optimal design pro-428 cess (simply take the current posterior as the prior for a subsequent iteration with ad-429 ditional data). 430

The design algorithm is very efficient with respect to evaluations of the GCM and 431 the high-resolution model. Due to our integration of the design framework within the 432 CES algorithm (detailed in Appendix A), only a relatively modest  $\mathcal{O}(100)$  forward model 433 evaluations of the GCM are required for the design selection process; no evaluations of 434 the high-resolution model (or, in our proof-of-concept, of the surrogate for it) are required. 435 The calculation of the utility function can be performed in an embarrassingly parallel 436 fashion. Thus, for moderately sized design spaces, the computational cost is dominated 437 by the cost of evaluation of the GCM. 438

Despite being efficient, the current algorithm relies on evaluating utilities naively 439 at all design points. In a practical climate model application, where we may have  $10^2$ 440 LES that are computationally affordable to be placed optimally within  $10^6$  or more pos-441 sible locations, such naive approaches are infeasible. Instead, one can use more sophis-442 ticated optimization algorithms. For determinant based (i.e., D-optimal) utilities, this 443 typically requires accelerating the determinant evaluation (and its gradients). Various 444 methods have been developed to do so, e.g., using Laplace approximations (Long et al., 445 2013; Beck et al., 2018; Rue et al., 2009), polynomial chaos surrogates (Huan & Mar-446 zouk, 2014), optimization of criteria bounds (Tsilifis et al., 2017), fast random determi-447 nant approximation (Alexanderian et al., 2014; Alexanderian & Saibaba, 2018), and Gaus-448 sian process surrogates (Buathong et al., 2020; Paglia et al., 2020). The latter, kernel-449 based approaches are particularly amenable to our setting, as they give sparse represen-450 tations of the utility function that are independent of the underlying computational grid. 451 They may offer a way forward in the climate modeling setting. 452

# 453 Acknowledgments

<sup>54</sup> We gratefully acknowledge the generous support of Eric and Wendy Schmidt (by rec-

- <sup>455</sup> ommendation of Schmidt Futures) and the National Science Foundation (grant AGS-1835860).
- 456 The simulations were performed on Caltech's High Performance Cluster, which is par-
- tially supported by a grant from the Gordon and Betty Moore Foundation. AMS is also
- supported by the Office of Naval Research (grant N00014-17-1-2079).
- 459 Data Availability. All computer code used in this paper is open source. The code for
  460 the idealized GCM, the Julia code for the optimal design algorithm, the plot tools, and
  461 the slurm/bash scripts to run both GCM and design algorithms are available at:
- 462 https://doi.org/10.5281/zenodo.5835269.

463	References
464	Alexanderian, A., Petra, N., Stadler, G., & Ghattas, O. (2014). A-optimal design
465	of experiments for infinite-dimensional bayesian linear inverse problems with
466	regularized \ell_0-sparsification. SIAM Journal on Scientific Computing, 36(5),
467	A2122–A2148.
468	Alexanderian, A., & Saibaba, A. K. (2018). Efficient d-optimal design of experi-
469	ments for infinite-dimensional bayesian linear inverse problems. SIAM Journal
470	on Scientific Computing, $40(5)$ , A2956–A2985.
471	Beck, J., Dia, B. M., Espath, L. F., Long, Q., & Tempone, R. (2018). Fast bayesian
472	experimental design: Laplace-based importance sampling for the expected in-
473	formation gain. Computer Methods in Applied Mechanics and Engineering,
474	334, 523 - 553.
475	Bischoff, T., & Schneider, T. (2014). Energetic constraints on the position of the In-
476	tertropical Convergence Zone. J. Climate, 27, 4937–4951. doi: 10.1175/JCLI-D
477	-13-00650.1
478	Bishop, C. H., & Toth, Z. (1999). Ensemble transformation and adaptive observa-
479	tions. Journal of the atmospheric sciences, 56(11), 1748–1765.
480	Bony, S., Colman, R., Kattsov, V. M., Allan, R. P., Bretherton, C. S., Dufresne,
481	JL Webb, M. J. (2006). How well do we understand and evalu-
482	ate climate change feedback processes? J. Climate, 19, 3445–3482. doi:
483	10.1175/JCLI3819.1
484	Bony, S., & Dufresne, J. L. (2005). Marine boundary layer clouds at the heart of
485	tropical cloud feedback uncertainties in climate models. <i>Geophys. Res. Lett.</i> .
486	<i>32</i> , L20806.
487	Bordoni, S., & Schneider, T. (2008a). Monsoons as eddy-mediated regime transitions
488	of the tropical overturning circulation. <i>Nature Geoscience</i> , 1(8), 515–519.
480	Bordoni S & Schneider T (2008b) Monsoons as eddy-mediated regime transi-
490	tions of the tropical overturning circulation. <i>Nature Geosci.</i> 1, 515–519. doi:
491	10.1038/ngeo248
492	Brient, F., & Schneider, T. (2016). Constraints on climate sensitivity from space-
493	based measurements of low-cloud reflection. J. Climate, 29, 5821–5835. doi: 10
494	.1175/JCLI-D-15-0897.1
495	Buathong, P., Ginsbourger, D., & Kritvakierne, T. (2020). Kernels over sets of finite
496	sets using RKHS embeddings, with application tobBayesian (combinatorial)
497	optimization. In International conference on artificial intelligence and statistics
498	(pp. 2731–2741).
499	Campin, JM., Hill, C., Jones, H., & Marshall, J. (2011). Super-parameterization in
500	ocean modeling: Application to deep convection. Ocean Modelling. 36(1), 90 -
501	101. doi: https://doi.org/10.1016/j.ocemod.2010.10.003
502	Cess, R. D., Potter, G., Blanchet, J., Boer, G., Ghan, S., Kiehl, J., others
503	(1989). Interpretation of cloud-climate feedback as produced by 14 atmo-
504	spheric general circulation models. <i>Science</i> , 245, 513–516.
505	Cess, R. D., Potter, G. L., Blanchet, J. P., Boer, G. J., Del Genio, A. D., Déqué,
506	M Zhang, MH. (1990). Intercomparison and interpretation of climate
507	feedback processes in 19 atmospheric general circulation models. J. Geophys.
508	<i>Res.</i> 95, 16601–16615. doi: $10.1029/JD095iD10p16601$
509	Chaloner, K., & Verdinelli, I. (1995). Bayesian experimental design: A review. Sta-
510	tistical Science, $10(3)$ , $273-304$ .
511	Chen, Y., & Oliver, D. S. (2012). Ensemble randomized maximum likelihood method
512	as an iterative ensemble smoother. Mathematical Geosciences, $44(1)$ , 1–26.

Cleary, E., Garbuno-Inigo, A., Lan, S., Schneider, T., & Stuart, A. M. (2021a). 513 Journal of Computational Physics, 424, 109716. Calibrate, emulate, sample. 514 Retrieved from https://www.sciencedirect.com/science/article/pii/ 515  $\texttt{S0021999120304903} \hspace{0.1 cm} \text{doi: } \hspace{0.1 cm} \text{https://doi.org/10.1016/j.jcp.2020.109716}$ 516

Cleary, E., Garbuno-Inigo, A., Lan, S., Schneider, T., & Stuart, A. M. (2021b). Cal-517 ibrate, emulate, sample. J. Comp. Phys., 424, 109716. Retrieved from 10 518 .1016/j.jcp.2020.109716 519 Cook, A. R., Gibson, G. J., & Gilligan, C. A. (2008). Optimal observation times in 520 experimental epidemic processes. *Biometrics*, 64(3), 860-868. 521 Couvreux, F., Hourdin, F., Williamson, D., Roehrig, R., Volodina, V., Villefranque, 522 N., ... others (2021). Process-based climate model development harnessing 523 machine learning: I. A calibration tool for parameterization improvement. J. 524 Adv. Model. Earth Sys., 13, e2020MS002217. doi: 10.1029/2020MS002217 525 Dashti, M., & Stuart, A. M. (2013).The Bayesian approach to inverse problems. 526 arXiv preprint arXiv:1302.6989. 527 de Rooy, W. C., Bechtold, P., Fröhlich, K., Hohenegger, C., Jonker, H., Mironov, D., 528 ... Yano, J.-I. (2013). Entrainment and detrainment in cumulus convection: 529 an overview. Quart. J. Roy. Meteor. Soc., 139, 1–19. 530 Dunbar, O. R. A., Garbuno-Inigo, A., Schneider, T., & Stuart, A. M. (2021).531 532 Calibration and uncertainty quantification of convective parameters in an idealized GCM. J. Adv. Model. Earth Sys., 13, e2020MS002454. doi: 533 10.1029/2020MS002454 534 Duncan, A. B., Stuart, A. M., & Wolfram, M.-T. (2021).Ensemble inference 535 methods for models with noisy and expensive likelihoods. arXiv preprint 536 arXiv:2104.03384. 537 Emanuel, K., Raymond, D., Betts, A., Bosart, L., Bretherton, C., Droegemeier, K., 538 Report of the first prospectus development team of the ... others (1995).539 us weather research program to noaa and the nsf. Bulletin of the American 540 Meteorological Society, 1194–1208. 541 Emerick, A. A., & Reynolds, A. C. (2013). Ensemble smoother with multiple data 542 assimilation. Computers & Geosciences, 55, 3–15. 543 Fox-Kemper, B., & Menemenlis, D. (2013). Can large eddy simulation techniques 544 improve mesoscale rich ocean models? In Ocean modeling in an eddying regime 545 (p. 319-337). American Geophysical Union (AGU). doi: 10.1029/177GM19 546 (2007).The dynamics of idealized convection schemes and their Frierson, D. M. 547 effect on the zonally averaged tropical circulation. Journal of the Atmospheric 548 Sciences, 64(6), 1959-1976. 549 Frierson, D. M. W. (2007). The dynamics of idealized convection schemes and their 550 effect on the zonally averaged tropical circulation. J. Atmos. Sci., 64, 1959-551 1976. 552 Frierson, D. M. W., Held, I. M., & Zurita-Gotor, P. (2006). A gray-radiation aqua-553 planet moist GCM. Part I: Static stability and eddy scale. J. Atmos. Sci., 63, 554 2548 - 2566.555 Geyer, C. J. (2011).Introduction to markov chain monte carlo. In S. Brooks, 556 A. Gelman, G. L. Jones, & X.-L. Meng (Eds.), Handbook of markov chain 557 monte carlo (pp. 3–48). Chapman and Hall/CRC. 558 Hohenegger, C., & Bretherton, C. S. (2011). Simulating deep convection with a shal-559 low convection scheme. Atmos. Chem. Phys., 11, 10389–10406. doi: 10.5194/ 560 acp-11-10389-2011 561 Hourdin, F., Williamson, D., Rio, C., Couvreux, F., Roehrig, R., Villefranque, N., 562 ... Volodina, V. (2021). Process-based climate model development harnessing 563 machine learning: II. model calibration from single column to global. J. Adv. 564 Model. Earth Sys., 13, e2020MS002225. doi: 10.1029/2020MS002225 565 Howland, M. F., Dunbar, O. R. A., & Schneider, T. (2021).Parameter uncer-566 tainty quantification in an idealized gcm with a seasonal cycle. arXiv preprint 567 arXiv:2108.00827. 568 Huan, X., & Marzouk, Y. (2014). Gradient-based stochastic optimization methods in 569 bayesian experimental design. International Journal for Uncertainty Quantifi-570 cation, 4(6). 571

- Huan, X., & Marzouk, Y. M. (2013, 1). Simulation-based optimal bayesian experimental design for nonlinear systems. *Journal of Computational Physics*, 232(1).
- Iglesias, M. A., Law, K. J., & Stuart, A. M. (2013). Ensemble kalman methods for inverse problems. *Inverse Problems*, 29(4), 045001.
- Kaipio, J., & Somersalo, E. (2006). Statistical and computational inverse problems
   (Vol. 160). Springer Science & Business Media.
- Kalnay, E. (2003). Atmospheric modeling, data assimilation and predictability. Cambridge, UK: Cambridge Univ. Press.
- Kaspi, Y., & Schneider, T. (2011). Winter cold of eastern continental boundaries in duced by warm ocean waters. *Nature*, 471, 621–624.
- Kaspi, Y., & Schneider, T. (2013). The role of stationary eddies in shaping midlati tude storm tracks. J. Atmos. Sci., 70, 2596–2613.
- Kennedy, M. C., & O'Hagan, A. (2000). Predicting the output from a complex computer code when fast approximations are available. *Biometrika*, 87(1), 1–13.
- Kennedy, M. C., & O'Hagan, A. (2001a). Bayesian calibration of computer mod els. Journal of the Royal Statistical Society: Series B (Statistical Methodology),
   63(3), 425-464.
- Kennedy, M. C., & O'Hagan, A. (2001b). Bayesian calibration of computer models.
   J. Roy. Statist. Soc. B, 63, 425–464. doi: 10.1111/1467-9868.00294
- Khairoutdinov, M. F., Krueger, S. K., Moeng, C.-H., Bogenschutz, P. A., & Randall,
  D. A. (2009). Large-eddy simulation of maritime deep tropical convection. J. Adv. Model. Earth Sys., 1, Art. #15, 13 pp. doi: 10.3894/JAMES.2009.1.15
  - Kim, W., Pitt, M. A., Lu, Z.-L., Steyvers, M., & Myung, J. I. (2014). A hierarchical adaptive approach to optimal experimental design. *Neural Computation*, 26(11), 2465-2492.
  - Levine, M. E., & Stuart, A. M. (2021). A framework for machine learning of model error in dynamical systems.
    - (https://arxiv.org/abs/2107.06658)

595

596

597

598

599

600

601

602

603

604

605

606

- Levine, X., & Schneider, T. (2015). Baroclinic eddies and the extent of the Hadley circulation: An idealized GCM study. J. Atmos. Sci., 72, 2744–2761. doi: 10 .1175/JAS-D-14-0152.1
- Li, Q., & Fox-Kemper, B. (2017). Assessing the effects of langmuir turbulence on the entrainment buoyancy flux in the ocean surface boundary layer. Journal of Physical Oceanography, 47(12), 2863–2886.
- Li, Q., Reichl, B. G., Fox-Kemper, B., Adcroft, A. J., Belcher, S. E., Danabasoglu,
   G., ... others (2019). Comparing ocean surface boundary vertical mixing
   schemes including langmuir turbulence. Journal of Advances in Modeling Earth
   Systems, 11(11), 3545–3592.
- Liu, C., Moncrieff, M. W., & Grabowski, W. W. (2001). Hierarchical modelling of tropical convective systems using explicit and parametrized approaches. *Quart. J. Roy. Meteor. Soc.*, 127, 493–515.
- Long, Q., Scavino, M., Tempone, R., & Wang, S. (2013). Fast estimation of expected information gains for bayesian experimental designs based on laplace
   approximations. Computer Methods in Applied Mechanics and Engineering, 259, 24–39.
- Lorenz, E. N., & Emanuel, K. A. (1998). Optimal sites for supplementary weather observations: Simulation with a small model. J. Atmos. Sci., 55, 399–414. doi: 10.1175/1520-0469(1998)055(0399:OSFSWO)2.0.CO;2
- Matheou, G., & Chung, D. (2014). Large-eddy simulation of stratified turbulence.
   Part II: Application of the stretched-vortex model to the atmospheric bound ary layer. J. Atmos. Sci., 71, 4439–4460. doi: 10.1175/JAS-D-13-0306.1
- Merlis, T. M., & Schneider, T. (2011). Changes in zonal surface temperature gradi ents and walker circulations in a wide range of climates. J. Climate, 24, 4757–
   4768.

Notz, W. I., Santner, T. J., & Williams, B. J. (2018).The design and analysis of 627 computer experiments (2nd ed. ed.). Springer. 628 (2011).The effective static stability experienced by eddies in a O'Gorman, P. A. 629 moist atmosphere. J. Atmos. Sci., 68, 75–90. 630 O'Gorman, P. A., Lamquin, N., Schneider, T., & Singh, M. S. (2011). The relative 631 humidity in an isentropic advection-condensation model: Limited poleward in-632 fluence and properties of subtropical minima. J. Atmos. Sci., 68, 3079–3093. 633 O'Gorman, P. A., & Schneider, T. (2008a). Energy of midlatitude transient eddies 634 in idealized simulations of changed climates. J. Climate, 21, 5797-5806. 635 O'Gorman, P. A., & Schneider, T. (2008b). The hydrological cycle over a wide range 636 of climates simulated with an idealized GCM. J. Climate, 21, 3815–3832. 637 O'Gorman, P. A., & Schneider, T. (2009a). The physical basis for increases in pre-638 cipitation extremes in simulations of 21st-century climate change. Proc. Natl. 639 Acad. Sci., 106, 14773-14777. 640 O'Gorman, P. A., & Schneider, T. (2009b). Scaling of precipitation extremes over 641 a wide range of climates simulated with an idealized GCM. J. Climate, 22, 642 5676-5685. 643 Oliver, D. S., Reynolds, A. C., & Liu, N. (2008). Inverse theory for petroleum reser-644 voir characterization and history matching. Cambridge Univ. Press. 645 O'Gorman, P. A., & Schneider, T. (2008). The hydrological cycle over a wide range 646 of climates simulated with an idealized gcm. Journal of Climate, 21(15), 3815-647 3832. doi: 10.1175/2007JCLI2065.1 648 Paglia, J., Eidsvik, J., & Karvanen, J. (2020). Efficient spatial designs using haus-649 dorff distances and bayesian optimisation. Statistical modeling for safer drilling 650 operations, 77. 651 Paninski, L. (2005). Asymptotic theory of information-theoretic experimental design. 652 Neural Computation, 17(7), 1480-1507. 653 Pressel, K. G., Kaul, C. M., Schneider, T., Tan, Z., & Mishra, S. (2015). Large-eddy 654 simulation in an anelastic framework with closed water and entropy balances. 655 J. Adv. Model. Earth Sys., 7, 1425–1456. doi: 10.1002/2015MS000496 656 Pressel, K. G., Mishra, S., Schneider, T., Kaul, C. M., & Tan, Z. (2017). Numerics 657 and subgrid-scale modeling in large eddy simulations of stratocumulus clouds. 658 J. Adv. Model. Earth Sys., 9, 1342-1365. doi: 10.1002/2016MS000778 659 Reich, S. (2011). A dynamical systems framework for intermittent data assimilation. 660 BIT Numerical Mathematics, 51(1), 235–249. 661 Reichl, B. G., & Hallberg, R. (2018). A simplified energetics based planetary bound-662 ary layer (epbl) approach for ocean climate simulations. Ocean Modelling, 132, 663 112 - 129. doi: https://doi.org/10.1016/j.ocemod.2018.10.004 664 Reichl, B. G., Wang, D., Hara, T., Ginis, I., & Kukulka, T. (2016). Langmuir tur-665 bulence parameterization in tropical cyclone conditions. Journal of Physical 666 Oceanography, 46(3), 863-886.667 Romps, D. M. (2016). The Stochastic Parcel Model: A deterministic parameteriza-668 tion of stochastically entraining convection. J. Adv. Model. Earth Sys., 8, 319-669 344. doi: 10.1002/2015MS000537 670 Rue, H., Martino, S., & Chopin, N. (2009).Approximate bayesian inference for 671 latent gaussian models by using integrated nested laplace approximations. 672 Journal of the Royal Statistical Society: Series B (Statistical Methodology), 673 71(2), 319-392.674 Ryan, E. G., Drovandi, C. C., McGree, J. M., & Pettitt, A. N. (2016).A review 675 of modern computational algorithms for bayesian optimal design. International 676 Statistical Review, 84(1), 128-154. 677 Ryan, E. G., Drovandi, C. C., Thompson, M. H., & Pettitt, A. N. (2014). Towards 678 bayesian experimental design for nonlinear models that require a large number 679 of sampling times. Computational Statistics & Data Analysis, 70, 45 - 60. 680

681 682	Schalkwijk, J., Jonker, H. J. J., Siebesma, A. P., & Van Meijgaard, E. (2015). Weather forecasting using GPU-based large-eddy simulations. Bull. Amer.
683	Meteor. Soc., 96, $(15-723. \text{ doi: } 10.1175/\text{BAMS-D-14-00114.1})$
685	verse problems. SIAM Journal on Numerical Analysis, 55(3), 1264–1290.
686	Schneider, T., Lan, S., Stuart, A., & Teixeira, J. (2017). Earth system model-
687	ing 2.0: A blueprint for models that learn from observations and targeted
688	high-resolution simulations. Geophys. Res. Lett., 44, 12396–12417. doi:
689	10.1002/2017 GL076101
690 691	Schneider, T., & O'Gorman, P. A. (2008). Moist convection and the thermal stratifi- cation of the extratropical troposphere. J. Atmos. Sci., 65, 3571–3583.
692	Schneider, T., O'Gorman, P. A., & Levine, X. J. (2010). Water vapor
693	and the dynamics of climate changes. <i>Rev. Geophys.</i> , 48, RG3001.
694	(doi:10.1029/2009 RG000302)
695	Schneider, T., Stuart, A. M., & Wu, J. (2021). Imposing sparsity within ensemble
696	Kalman inversion.
697	Schneider, T., Teixeira, J., Bretherton, C. S., Brient, F., Pressel, K. G., Schär, C.,
698	& Siebesma, A. P. (2017). Climate goals and computing the future of clouds.
699	Nature Climate Change, 7, 3–5. doi: 10.1038/nclimate3190
700	Shen, Z., Pressel, K. G., Tan, Z., & Schneider, T. (2020). Statistically steady
701	state large-eddy simulations forced by an idealized GCM: 1. forcing frame-
702	work and simulation characteristics. J. Adv. Model. Earth Sys., 12. doi: 10.1020/2010MC001914
703	10.1029/2019MS001814
704	Snen, Z., Sridnar, A., Tan, Z., Jaruga, A., & Schneider, T. (2021). A li-
705	https://essoar.org_doi: https://doi.org/10.1002/essoar.10507112.1
700	Siebesma A P Bretherton C S Brown A Chlond A Cuxart J Duvnkerke
708	P. G Stevens, D. E. (2003). A large eddy simulation intercomparison
709	study of shallow cumulus convection. J. Atmos. Sci., 60, 1201–1219.
710	Siebesma, A. P., Soares, P. M. M., & Teixeira, J. (2007). A combined eddy-
711	diffusivity mass-flux approach for the convective boundary layer. J. Atmos.
712	Sci., $64$ , 1230–1248. doi: 10.1175/JAS3888.1
713	Smalley, M., Suselj, K., Lebsock, M., & Teixeira1, J. (2019). A novel framework
714	for evaluating and improving parameterized subtropical marine boundary layer
715	cloudiness. Mon. Wea. Rev., 147, 3241–3260.
716	Souza, A. N., Wagner, G. L., Ramadhan, A., Allen, B., Churavy, V., Schloss, J.,
717	Ferrari, R. (2020). Uncertainty quantification of ocean parameteriza-
718	tions: Application to the k-profile-parameterization for penetrative convection.
719	Journal of Advances in Modeling Earth Systems, 12(12), e2020MIS002108.
720	$(e^{2020}MS002108)$ ( $e^{2020}MS002108$ 10 1029/2020MS002108) doi:
722	https://doi.org/10.1029/2020MS002108
723	Stephens, G. L. (2005). Cloud feedbacks in the climate system: A critical review. J.
724	Climate, 18, 237–273. doi: 10.1175/JCLI-3243.1
725	Stevens, B., Moeng, CH., Ackerman, A. S., Bretherton, C. S., Chlond, A., de
726	Roode, S., Zhu, P. (2005). Evaluation of large-eddy simulations via obser-
727	vations of nocturnal marine stratocumulus. Mon. Wea. Rev., 133, 1443–1462.
728	doi: 10.1175/MWR2930.1
729	Stuart, A. M. (2010). Inverse problems: a Bayesian perspective. Acta Numerica, 19,
730	401–009. The 7 Keyl C M Dressel K C Calar V Calar 1. $\Pi$ ( $\Pi$ ) $\Pi$ : I (2010)
731	Ian, Z., Kaul, C. M., Pressel, K. G., Conen, Y., Schneider, I., & Teixeira, J. (2018).
733	subgrid-scale turbulence and convection. J. Adv. Model. Earth Sys., 10, 770-

800. doi:  $10.1002/2017 {\rm MS001162}$ 

734



Figure A1. Procedure of the uncertainty quantification framework (blue), to produce output (pink). A restriction operator  $W_1$  extracting a subset of the GCM output (yellow); the subsequent emulate and sample stages may be performed in parallel for different  $W_i$ .

- Tarantola, A. (2005). Inverse problem theory and methods for model parameter esti mation (Vol. 89). siam.
- Tsilifis, P., Ghanem, R. G., & Hajali, P. (2017). Efficient bayesian experimentation
   using an expected information gain lower bound. SIAM/ASA Journal on Uncertainty Quantification, 5(1), 30–62.
- Van Roekel, L., Adcroft, A. J., Danabasoglu, G., Griffies, S. M., Kauffman, B.,
- Large, W., ... Schmidt, M. (2018). The kpp boundary layer scheme for the ocean: Revisiting its formulation and benchmarking one-dimensional simulations relative to les. *Journal of Advances in Modeling Earth Systems*, 10(11), 2647–2685.
- Vial, J., Dufresne, J.-L., & Bony, S. (2013). On the interpretation of inter-model
  spread in CMIP5 climate sensitivity estimates. *Clim. Dyn.*, 41, 3339–3362. doi: 10.1007/s00382-013-1725-9
- Wang, D., Large, W. G., & McWilliams, J. C. (1996). Large-eddy simulation of the equatorial ocean boundary layer: Diurnal cycling, eddy viscosity, and horizon-tal rotation. *Journal of Geophysical Research: Oceans*, 101(C2), 3649–3662.
- Webb, M. J., Lambert, F. H., & Gregory, J. M. (2013). Origins of differences in cli mate sensitivity, forcing and feedback in climate models. *Clim. Dyn.*, 40, 677–
   707. doi: 10.1007/s00382-012-1336-x
- Wei, H.-H., & Bordoni, S. (2018). Energetic constraints on the ITCZ position in ide alized simulations with a seasonal cycle. J. Adv. Model. Earth Sys., 10. doi: 10
   .1029/2018MS001313
- Williams, C. K., & Rasmussen, C. E. (2006). Gaussian processes for machine learn *ing* (Vol. 2) (No. 3). MIT press Cambridge, MA.
- Wills, R. C., Levine, X. J., & Schneider, T. (2017). Local energetic constraints on
   Walker circulation strength. J. Atmos. Sci., 74, 1907–1922. doi: 10.1175/JAS
   -D-16-0219.1
- Zhang, F., Sun, Y. Q., Magnusson, L., Buizza, R., Lin, S.-J., Chen, J.-H., &
   Emanuel, K. (2019). What is the predictability limit of midlatitude
   weather? Journal of the Atmospheric Sciences, 76(4), 1077 1091. Re trieved from https://journals.ametsoc.org/view/journals/atsc/76/4/
   jas-d-18-0269.1.xml doi: 10.1175/JAS-D-18-0269.1
- Zhang, M., Bretherton, C. S., Blossey, P. N., Austin, P. H., Bacmeister, J. T., Bony,
  S., ... others (2013). CGILS: Results from the first phase of an international project to understand the physical mechanisms of low cloud feedbacks
  in general circulation models. J. Adv. Model. Earth Sys., 5, 826–842. doi:
  10.1002/2013MS000246

# Appendix A Calibrate-Emulate-Sample with design

773 One fundamental aspect of this work, is the ability to efficiently calculate the the 774 posterior distribution (in particular the covariance), which is needed to calculate the util-775 ity function (6) at all designs. We present a methodology: calibrate-extract-emulate-sample, 776 (CEES) which allows for the calculation of posterior covariance for all designs with just 777  $\mathcal{O}(100)$  evaluations of our forward model.

The methodology is based on the calibrate-emulate-sample (CES) algorithm, for 778 full details of the individual stages see Cleary et al. (2021a); Dunbar et al. (2021), here 779 we present an overview and motivation. The core purpose of CES is to form a compu-780 tationally cheap statistical emulator of  $\mathcal{G}_{\infty}$  from intelligently chosen samples of  $\mathcal{G}_T$ ; then 781 one is able to solve the Bayesian inverse problem for the emulated  $\mathcal{G}_{\infty}$  with a sampling 782 method. We achieve this by using Gaussian process emulators, trained on the samples 783 of the (noisy and expensive) forward map. The Gaussian process mean function is nat-784 urally smoother than the data it is trained on (Kennedy & O'Hagan, 2001a; Notz et al., 785 2018), and is capable of representing the the noise of the forward model within the co-786 variance function, leading to a smooth likelihood function that is quick to evaluate. The 787 training points for the Gaussian Process are given by applying an optimization scheme, 788 EKI (Ensemble Kalman Inversion), (Iglesias et al., 2013; Schillings & Stuart, 2017) to 789 the inverse problem in its finite-time averaged form (3). Theoretical work shows that noisy 790 continuous-time versions of EKI exhibit an averaging effect that skips over fluctuations 791 superimposed onto the ergodic averaged forward model (Duncan et al., 2021), and sim-792 ilar effects are observed in practice for EKI, thus it is highly suited to optimization of 793 parameters coming from a noisy, expensive model without derivatives available. Ensem-794 ble Kalman methods are scalable to very high dimensional problems (Kalnay, 2003; Oliver 795 et al., 2008) with use of localization and regularization. 796

Let *D* index a finite space of designs. Given a time T > 0, and prior on  $\boldsymbol{\theta}$  with prior mean  $\boldsymbol{\theta}^*$ . Draw a sample  $\boldsymbol{y} = \mathcal{G}_T(\boldsymbol{\theta}^*, \boldsymbol{v}^{(0)})$ , from any initial condition  $\boldsymbol{v}^{(0)}$ :

1. Calibrate: We solve (3) with  $\boldsymbol{y}$  using evaluations of  $\mathcal{G}_T$  in an optimization sense, where we minimize the functional.

$$\Phi_T(\boldsymbol{\theta}, \boldsymbol{y}) = \|\boldsymbol{y} - \mathcal{G}_T(\boldsymbol{\theta}; \boldsymbol{v}^{(0)})\|_{2\Sigma}^2.$$
(A1)

The notation  $\|\cdot\|_{\Sigma} = \|\Sigma^{-\frac{1}{2}}\cdot\|_2$  is the Mahalanobis distance. The weight  $2\Sigma$  is 799 the sum of internal variability of  $\mathcal{G}_T$  and of y. The optimization is performed us-800 ing several iterations the Ensemble Kalman Inversion algorithm. This leads to  $\{\boldsymbol{\theta}_i, \mathcal{G}_j(\boldsymbol{\theta}_j)\}_{i=1}^J$ 801 of input-output pairs that are localized around the optimal parameter value. 802 2. Extract: For each design  $k \in D$ , we apply the restriction mapping  $W_k$  to the 803 forward map,  $\{\boldsymbol{\theta}_j, W_k \mathcal{G}_T(\boldsymbol{\theta}_j)\}_{i=1}^J$ , and apply the following **Emulate(k)** and **Sam-**804 ple(k) stages. 805 3. Emulate(k): We decorrelate the data space with an SVD on the internal vari-806

3. Emulate(k): We decorrelate the data space with an SVD on the internal variability covariance  $\Sigma$ , yielding a change-of-basis matrix V. We train Gaussian process emulators, on the pairs  $\{\boldsymbol{\theta}_j, VW_k\mathcal{G}_T(\boldsymbol{\theta}_j)\}_{j=1}^J$ , yielding  $(\mathcal{G}_{\mathrm{GP}}(\boldsymbol{\theta}), \Sigma_{\mathrm{GP}}(\boldsymbol{\theta}))$ , where  $\mathcal{G}_{\mathrm{GP}} \approx VW_k\mathcal{G}_{\infty}(\boldsymbol{\theta})$  (crucially  $\mathcal{G}_{\infty}$  and not  $\mathcal{G}_T$ ) and  $\Sigma_{\mathrm{GP}}(\boldsymbol{\theta}) \approx VW_k\Sigma W_k^T V^T$ .

807

808

809

4. Sample(k): We now solve the inverse problem (5), This is feasible as the emulator provides us with an approximation of  $\mathcal{G}_{\infty}$  (not just  $\mathcal{G}_T$ ). The posterior distribution associated with (5) is proportional to a product of prior and likelihood contribution from Bayes theorem. Explicitly, for a Gaussian prior  $N(\boldsymbol{m}, C)$  on the computational parameters, and the likelihood dependent on the emulator, we write the MCMC objective function (also known as the log-posterior) as

$$\begin{split} \Phi_{\mathrm{MCMC}}(\boldsymbol{\theta}, VW_k \boldsymbol{y}) = & \frac{1}{2} \| VW_k \boldsymbol{y} - \mathcal{G}_{\mathrm{GP}}(\boldsymbol{\theta}) \|_{\Sigma_{\mathrm{GP}}(\boldsymbol{\theta})}^2 + \frac{1}{2} \log \det \Sigma_{\mathrm{GP}}(\boldsymbol{\theta}) \\ &+ \frac{1}{2} \| \boldsymbol{\theta} - \boldsymbol{m} \|_C^2 \,. \end{split}$$

The posterior is then given by

810

 $\mathbb{P}(\boldsymbol{\theta} \mid VW_k \boldsymbol{y}) \propto \exp(-\Phi_{MCMC}(\boldsymbol{\theta}, VW_k \boldsymbol{y})).$ 

This can be sampled with a standard random walk metropolis sampling algorithm.

We then collect the posterior distributions  $\{\boldsymbol{\theta} \mid W_k \boldsymbol{y}\}_k$ ,  $\forall k \in D$  and calculate the utility function using (6). This CEES algorithm in Figure A1. In particular, note that the subsampling occurs after the J model evaluations, therefore all posterior distributions can be performed in an embarassingly parallel fashion, and all use the same forward model evaluations.

The CEES algorithm is also used to solve (2) at a given design  $\tilde{k}$ , with this algorithm using the model  $S_T(\cdot; \tilde{k}, \cdot)$ , and data sample  $\boldsymbol{z}_k = S_T(\boldsymbol{\theta}; \tilde{k}, \boldsymbol{v}^{(0)}) + \delta$ , and weighting the data misfit norm with the additional contribution from  $\delta$ . We then perform **Emulate** $(\tilde{k})$ , and **Sample** $(\tilde{k})$  at the chosen design.