

# Sequential machine learning for data assimilation with parameter estimation via hierarchical Bayesian optimization

Yang Sun<sup>1\*</sup> | Song Xi Chen<sup>2†</sup> | Yumou Qiu<sup>1\*</sup>

<sup>1</sup>School of Mathematical Sciences and and Center for Statistical Science, Peking University, Beijing, 100871, China

<sup>2</sup>Department of Statistics and Data Sciences, Tsinghua University, Beijing, 100084, China

## Correspondence

Song Xi Chen, PhD, Department of Statistics and Data Sciences, Tsinghua University, Beijing, 100084, China; Yumou Qiu, PhD, School of Mathematical Sciences and and Center for Statistical Science, Peking University, Beijing, 100871, China  
Email: sxchen@tsinghua.edu.cn; qiuyumou@math.pku.edu.cn

## Funding information

The performance of data assimilation critically depends on well-specified physical models, particularly accurate assignment of key parameters. Misspecified parameters can introduce model bias and uncertainty, degrading both forecasting and assimilation quality. This paper proposes a sequential data assimilation procedure with a hierarchical Bayesian machine learning approach for sequentially calibrating parameters in the physical model that integrates the ensemble Kalman filter. We define an objective function for parameter estimation by an optimally weighted distance between the observations and forecast states over a time segment. Based on this criterion, we develop a Bayesian learning algorithm to efficiently estimate the unknown parameters throughout the assimilation window. Statistical properties of the proposed method are developed. Numerical experiments on the Lorenz'96 system and a two-dimensional shallow water model demonstrate that the proposed method outperforms conventional approaches in terms of estimation accuracy and stability.

## KEYWORDS

Data assimilation, ensemble Kalman filter, hierarchical Bayesian optimization, nonlinear system, machine learning, parameter estimation

## 1 | INTRODUCTION

In large-scale geophysical models such as those used in atmospheric and oceanic systems, key physical parameters exert significant influence on dynamic modeling and forecast accuracy. Proper specification of these parameters is essential to ensure realistic dynamical behavior that reflects the underlying physics, as misspecified parameters can introduce systematic errors, degrade the data assimilation performance, and ultimately reduce forecast accuracy. Therefore, estimating the key parameters of the models is a fundamental task in data assimilation, aiming to improve both the state and the underlying model representation.

Data assimilation integrates observations with a dynamic numerical model to provide an estimator of the unknown time-evolving state [1, 2, 3]. As a critical tool in diverse fields such as meteorology [3, 4, 5], oceanography [6, 7], and automatic control [8], data assimilation enhances the accuracy of forecasting systems by assimilating model forecasts with observations, leading to more accurate forecasts and imputation of the state variables. The Kalman Filter (KF) [9] was a pioneering method for linear system assimilation with Gaussian noise, by providing optimal recursive solutions to the state variables. It propagates the mean and covariance of the state distribution through time, updating the state estimate whenever new observations become available. For nonlinear physical models, the ensemble Kalman Filter (EnKF) [2, 10] uses the forecast ensembles to estimate the forecast error covariance matrix and updates the system state via the analysis formula of the KF. It has become a popular data assimilation method.

There are a set of methods for implementing the EnKF. The inflation method was proposed to increase the spread of the forecast error covariance and to counter the potential model error, thereby reducing filter divergence [11, 12, 13]. To overcome the challenge in high-dimension covariance estimation, researchers developed the localization method to counter the deterioration of the ensemble forecast error covariance estimation, a technique that imposes sparsity by tapering or banding the covariance between two far away locations to zero, see [14, 15, 16] for details of the techniques.

Despite progress in data assimilation algorithms, parameter uncertainty remains a major limitation. Many geophysical models, such as oceanic or atmospheric systems, contain empirical estimates of parameters (e.g., vertical diffusivity, mixing and drag coefficients, optical depths) that strongly affect model dynamics and forecasting accuracy [17]. These parameters are typically pre-chosen based on prior knowledge, but any misspecification would lead to biased forecasts and inaccurate analysis states, regardless of the assimilation method used. There have been growing efforts devoted to incorporating parameter estimation within data assimilation systems. Early strategies include state augmentation, where parameters are augmented as state variables and are updated as part of the EnKF. Augmented state assimilation treats uncertain model parameters as an additional state variable and estimates it simultaneously with the dynamic state [4, 18]. Recently, machine learning methods have been incorporated in augmented data assimilation procedures [19].

Although the augmented assimilation method has a relatively low computational cost and has been widely used in fields such as the atmosphere and the ocean, it is sometimes difficult to ensure algorithmic convergence due to the unstable results of parameter estimation. The quasi maximal likelihood estimation (QMLE) based on the expectation-maximization algorithm was used for parameter estimation in nonlinear state space models [20, 21]. However, it relies on the ensemble Kalman smoother (EnKS), and requires re-running of the entire sequence of the EnKS for every candidate value of parameters, which incurs a heavy computation burden and may not be suitable for high-dimensional nonlinear systems. Meanwhile, [22, 23] proposed a nonlinear least squares (NLS) estimation for unknown parameters based Kalman filters. Similar as the QMLE method, it also requires re-running the entire data assimilation procedure each time the model parameters are updated during the optimization process, which inevitably leads to a substantial computational burden. Both QMLE and NLS methods do not support sequential parameter updates, and hence, can

not handle time-varying parameters.

To overcome the limitations of the existing methods, this paper aims to develop an efficient sequential machine learning algorithm coupled with the EnKF for data assimilation and estimating the underlying model parameter with a statistical guarantee. The proposed machine learning-based data assimilation algorithm meets the following four objectives: (i) sequential updating of both state variables and model parameters, (ii) computationally efficient without repeatedly recalculation of the forecast error covariance matrix in the optimization procedure for the model parameter, (iii) a more statistically efficient estimator for parameter estimation than the existing approaches, (iv) a statistical guarantee for the consistency and asymptotic distribution of the parameter estimate.

To achieve the aforementioned objectives, our method combines modern hierarchical Bayesian machine learning algorithms [24] with the EnKF framework, which leverages the adaptability of data-driven learning while preserving the predictive strengths of ensemble-based filtering through a sequential learning mechanism. We develop a data assimilation framework that integrates parameter estimation with state updating. Specifically, first, the entire time interval is divided into several segments. Parameter estimation and state updating are performed sequentially within each segment. Second, the estimated parameters in the current segment are used for an initial data assimilation for the next segment to compute the forecast error covariance matrices. This approach avoids repeated computation of the Kalman gain matrices during parameter optimization, thereby reducing the computational burden. A statistically efficient weighting matrix is selected to construct the empirical loss function for parameter estimation, which achieves minimum asymptotic variance. Finally, a hierarchical expected improvement (HEI) algorithm is employed to search for the minimum of the loss function. The HEI is an iterative algorithm, which searches for the global optimum with faster convergence by strengthening the exploration of areas with high prediction uncertainty.

To evaluate the performance of our proposed method, we conduct numerical simulations based on two benchmark models: the Lorenz'96 model with a single parameter and the shallow water model with two parameters. Numerical experiments show that our approach consistently outperformed traditional EnKF-based methods without parameter estimation, augmented data assimilation methods and existing parameter estimation methods in terms of the accuracy of both data assimilation and parameter estimation.

This paper is organized as follows. Section 2 presents the problem setting, introduces the identification conditions for parameter estimation in state-space models, and discusses the challenges associated with existing methods. Section 3 constructs the proposed sequential machine learning method for data assimilation and parameter estimation and establishes its asymptotic properties. Section 4 presents numerical studies on the Lorenz-96 model and the two-dimensional shallow water equations. Section 5 concludes the paper with discussions on future topics. The proofs of the theorems and additional numerical results are relegated to the supplementary material (SM).

## 2 | PROBLEM SETTING AND BACKGROUND

We consider the data assimilation (DA) setting with unknown parameters in the state-space models. Suppose the underlying state variable  $\mathbf{X}_t \in \mathcal{X} \subset \mathbb{R}^p$  and its observation  $\mathbf{Y}_t \in \mathbb{R}^q$  follow

$$\mathbf{X}_t = \mathcal{M}_t(\mathbf{X}_{t-1}, \boldsymbol{\theta}^*) + \boldsymbol{\eta}_t, \quad (1)$$

$$\mathbf{Y}_t = H_t \mathbf{X}_t + \mathbf{e}_t \quad (2)$$

for  $t = 1, \dots, T$ , where  $\boldsymbol{\eta}_t$  and  $\mathbf{e}_t$  are the model and observation errors with zero-mean and covariance matrices  $Q_t$  and  $R_t$ , respectively. This is the state-space model where  $\{\mathbf{X}_t\}_{t=1}^T$  may be unobservable while  $\{\mathbf{Y}_t\}_{t=1}^T$  are observable. It is

assumed that both  $\eta_t$  and  $e_t$  are independent of  $\{\mathbf{X}_1, \dots, \mathbf{X}_{t-1}\}$ , respectively. In this study, we do not impose specific distributional assumptions (e.g., Gaussianity) on  $\eta_t$  and  $e_t$ , allowing them to follow unknown distributions. Moreover,  $\mathcal{M}_t(\cdot, \cdot) : \mathbb{R}^p \times \Theta \rightarrow \mathbb{R}^p$  is a nonlinear physical model operator,  $H_t$  is a  $q \times p$  known linear observational matrix, and  $\boldsymbol{\theta}^*$  is an interior point of  $\Theta \subset \mathbb{R}^d$ , which represents the unknown true value of the key model parameters.

The purpose of data assimilation is to impute the state  $\mathbf{X}_t$  from observation  $\mathbf{Y}_t$  based on the physical model, which depends on choosing a quality estimate of  $\boldsymbol{\theta}^*$ . Misspecification of the parameter  $\boldsymbol{\theta}$  can result in systematic bias in the physical model and the data assimilation procedure [25]. In particular, when the physical model  $\mathcal{M}_t(\cdot, \boldsymbol{\theta})$  is highly sensitive to  $\boldsymbol{\theta}$ , such as the general circulation models (GCMs) [26], the resulting bias may degrade the performance of data assimilation and forecast, and in severe cases, lead to filter divergence or failure. The goal of this paper is to propose a joint estimation procedure for both states and parameters with a guarantee on statistical consistency, which can reduce model bias and maintain the quality of the analysis state.

Let  $\mathbf{X}_{t-1,j}^a$  be the  $j$ -th perturbed analysis ensemble at  $t-1$  for  $j = 1, \dots, n$ . Let  $\mathbf{X}_{t,j}^f(\boldsymbol{\theta}^*) = \mathcal{M}_t(\mathbf{X}_{t-1,j}^a, \boldsymbol{\theta}^*)$  be the one-step ensemble forecast of  $\mathbf{X}_t$  under the true parameter. Denote the  $j$ -th forecast error at  $\boldsymbol{\theta}^*$  as  $\boldsymbol{\delta}_{t,j} = \mathbf{X}_t - \mathbf{X}_{t,j}^f(\boldsymbol{\theta}^*)$ . Then, Equation (2) can be expressed as

$$\mathbf{Y}_t = H_t \mathcal{M}_t(\mathbf{X}_{t-1,j}^a, \boldsymbol{\theta}^*) + H_t \boldsymbol{\delta}_{t,j} + e_t. \quad (3)$$

Let  $\|\cdot\|$  denote the Euclidean norm of vectors. We make the following assumptions for the statistical identification of  $\boldsymbol{\theta}^*$  based on the observed data  $\mathbf{Y}_1, \dots, \mathbf{Y}_t$  via the expression in (3).

**(S1)** The analysis ensemble randomness among  $\{\mathbf{X}_{t-1,j}^a\}$  is independent with  $\{e_t\}_{t=1}^T$ .

**(S2)** The forecast error satisfies  $\mathbb{E}\{\mathbb{E}_{\text{ens}}(\boldsymbol{\delta}_{t,j}) \mid \mathbf{Y}_{1:t-1}\} = 0$ , where  $\mathbb{E}_{\text{ens}}$  represents taking expectation over the ensembles.

**(S3)** For every  $\epsilon > 0$ , we have  $\sup_{\|\boldsymbol{\theta} - \boldsymbol{\theta}^*\| \geq \epsilon} q^{-1} \mathbb{E}\|H_t \mathbb{E}_{\text{ens}}(\mathcal{M}_t(\mathbf{X}_{t-1,j}^a, \boldsymbol{\theta}) - \mathcal{M}_t(\mathbf{X}_{t-1,j}^a, \boldsymbol{\theta}^*))\|^2 > 0$ .

Assumption S1 is a natural condition as the analysis ensembles are determined by the observations  $\mathbf{Y}_{1:t-1}$  before time  $t$ . Assumption S2 assumes the unbiasedness of the average ensemble forecast of  $\mathbf{X}_t^f$  at the true parameter value  $\boldsymbol{\theta}^*$ . In applications, violations of S1 may arise from the misspecification of the model or observation operators, in which case bias modeling is required [18]. Assumption S3 imposes distinguishability of the parameter  $\boldsymbol{\theta}$  after applying the observation operator  $H_t$ . It ensures sufficient curvature of the population loss function at the true value  $\boldsymbol{\theta}^*$  so that the observed data can identify the target parameter.

For a  $q \times q$  positive definite matrix  $W$ , let  $\|\mathbf{a}\|_W = \mathbf{a}' W \mathbf{a}$  for  $\mathbf{a} \in \mathbb{R}^q$ . The following Lemma provides an identification criterion for  $\boldsymbol{\theta}^*$ .

**Lemma 1** Under the model in (1) and (2), if Assumptions S1–S3 hold at  $t$ , then for any  $q \times q$  positive definite matrix  $W$ , the true parameter  $\boldsymbol{\theta}^*$  is the unique minimum of the optimization problem:

$$\boldsymbol{\theta}^* = \arg \min_{\boldsymbol{\theta} \in \Theta} \mathbb{E} \|\mathbf{Y}_t - H_t \mathbb{E}_{\text{ens}} \{\mathcal{M}_t(\mathbf{X}_{t-1,j}^a, \boldsymbol{\theta})\}\|_W^2. \quad (4)$$

The lemma provides a theoretical basis for estimating  $\boldsymbol{\theta}^*$  by minimizing the squared residual between the observation and forecast states at a single time point. To conduct parameter estimation along with data assimilation as time progresses, we introduce the main steps of EnKF in the following.

For a given parameter value  $\boldsymbol{\theta}$  and the analysis ensembles  $\{\mathbf{X}_{t-1,j}^a\}_{j=1}^n$  at time  $t-1$ , the forecast ensembles are obtained by  $\mathbf{X}_{t,j}^f(\boldsymbol{\theta}) = \mathcal{M}_t(\mathbf{X}_{t-1,j}^a, \boldsymbol{\theta})$ . Let  $\Sigma_t^f = \Sigma_t^f(\boldsymbol{\theta}^*) = \text{Var}(\mathbf{X}_{t,j}^f(\boldsymbol{\theta}^*))$  be the forecast covariance matrix at the true

114 parameter value. Let  $\hat{\Sigma}_t^f(\boldsymbol{\theta})$  be the sample covariance of  $\{\mathbf{X}_{t,j}^f(\boldsymbol{\theta})\}_{j=1}^n$  and  $\hat{K}_t(\boldsymbol{\theta}) = \hat{\Sigma}_t^f(\boldsymbol{\theta})H_t'(H_t\hat{\Sigma}_t^f(\boldsymbol{\theta})H_t' + R_t)^{-1}$  be  
 115 the Kalman gain matrix, which highlights the role of  $\boldsymbol{\theta}$  on these quantities. Let  $e_{t,j}$  be an independent copy of the  
 116 observation error  $e_t$ , with mean 0 and covariance  $R_t$ . The EnKF updates the analysis ensembles at  $t$  via

$$\mathbf{X}_{t,j}^a(\hat{K}_t(\boldsymbol{\theta}), \boldsymbol{\theta}) = \mathbf{X}_{t,j}^f(\boldsymbol{\theta}) + \hat{K}_t(\boldsymbol{\theta})(\mathbf{Y}_t - H_t\mathbf{X}_{t,j}^f(\boldsymbol{\theta}) + e_{t,j}) \quad (5)$$

117 for  $j = 1, \dots, n$ , where  $\mathbf{X}_{t,j}^a(\hat{K}_t(\boldsymbol{\theta}), \boldsymbol{\theta})$  depends on  $\boldsymbol{\theta}$  through both the Kalman gain  $\hat{K}_t(\boldsymbol{\theta})$  and the forecast ensembles  
 118 via  $\mathcal{M}_t(\cdot, \boldsymbol{\theta})$ . Therefore, for each change of  $\boldsymbol{\theta}$ , we need to re-run the ensembles and re-estimate the Kalman gain  
 119 matrix to obtain an updated value of  $\mathbf{X}_{t,j}^a(\hat{K}_t(\boldsymbol{\theta}), \boldsymbol{\theta})$ . This would make parameter estimation very computationally  
 120 expensive for dynamic state-space models.

121 To overcome the challenge, we consider a sequential updating approach by partitioning the entire time horizon  
 122 evenly into  $K$  segments:  $\{1, \dots, T\} = \mathcal{T}_1 \cup \dots \cup \mathcal{T}_K = \cup_{k=1}^K \mathcal{T}_k$ , where  $L = T/K$  and the  $k$ -th segment  $\mathcal{T}_k = \{(k -$   
 123  $1)L + 1, \dots, kL\}$ . Our strategy is to update estimates of  $\boldsymbol{\theta}^*$  iteratively over the segments. Starting from an initial value,  
 124 say  $\boldsymbol{\theta}^{(0)}$ , we run an initial data assimilation procedure to obtain the Kalman gain matrices  $\{\hat{K}_t(\boldsymbol{\theta}^{(0)})\}$  using  $\boldsymbol{\theta} = \boldsymbol{\theta}^{(0)}$   
 125 for  $t \in \mathcal{T}_1$ . Fix those Kalman gain matrices, we calculate the analysis and forecast ensembles  $\mathbf{X}_{t,j}^a(\hat{K}_t(\boldsymbol{\theta}^{(0)}), \boldsymbol{\theta})$  and  
 126  $\mathbf{X}_{t,j}^f(\boldsymbol{\theta})$  in this segment by varying the value of  $\boldsymbol{\theta}$  in the model operator  $\mathcal{M}_t(\cdot, \boldsymbol{\theta})$ . We obtain an improved estimate  
 127  $\hat{\boldsymbol{\theta}}^{(1)}$  from the first segment  $\mathcal{T}_1$  by minimizing an empirical estimate of the loss function in (4). In general, given a  $\hat{\boldsymbol{\theta}}^{(k-1)}$   
 128 on the  $(k - 1)$ -th segment, a new estimate  $\hat{\boldsymbol{\theta}}^{(k)}$  is obtained based on  $\hat{\boldsymbol{\theta}}^{(k-1)}$  and the observations in the  $k$ -th segment  
 129  $\mathcal{T}_k$ . The rationale and detailed procedure of the proposed approach will be outlined in the next section.

130 Augmented-state EnKF methods (e.g., [19, 27]) incorporate unknown parameters into the state vector and apply  
 131 artificial evolution models (such as random walks) to enable joint state-parameter estimation. These approaches are  
 132 computationally straightforward and allow parameters to be updated sequentially in every time step. However, their  
 133 performance may be limited in complex or nonlinear systems, since the imposed parameter dynamics do not necessar-  
 134 ily reflect the true evolution, potentially affecting estimation stability and accuracy. On the other hand, optimization-  
 135 based methods such as the NLS estimator [23, 22] provide statistically consistent estimators by minimizing a global  
 136 loss function. However, they are computationally expensive, as evaluating each candidate parameter requires rerun-  
 137 ning the full data assimilation cycle over all time steps, resulting in substantial computational burden, high memory  
 138 consumption, and limited algorithmic scalability.

139 Our proposed framework bridges these two approaches by combining their respective strengths: it performs  
 140 parameter updates sequentially over time segments, like the augmented EnKF, while ensuring statistical consistency  
 141 through segment-wise optimization, as in the NLS. This hybrid framework enables scalable and accurate parameter  
 142 estimation in both high-dimensional and nonlinear state-space models.

### 143 | 3 | METHODOLOGY

144 We elaborate the detailed method for updating  $\hat{\boldsymbol{\theta}}^{(k-1)}$  to  $\hat{\boldsymbol{\theta}}^{(k)}$  using the forecast states and the observations in the  
 145  $k$ -th segment  $\mathcal{T}_k$ . Specifically, for each segment  $\mathcal{T}_k$  given the parameter estimate from the previous segment  $\hat{\boldsymbol{\theta}}^{(k-1)}$ ,  
 146 we build an objective function of  $\boldsymbol{\theta}$  using the weighted residuals between the forecast states and the observations  
 147 in this period by recycling a common Kalman gain matrix and the forecast error covariance over the segment. This  
 148 avoids repeatedly estimating the forecast error covariance matrices and the Kalman gain matrices while the value of  
 149  $\boldsymbol{\theta}$  changes. Motivated by this, we propose a sequential estimation framework for jointly estimating the states and  
 150 parameters of the state-space model via hierarchical Bayesian optimization [24].

### 3.1 | Sequential objective function

Given the estimate  $\hat{\theta}^{(k-1)}$  from the  $(k-1)$ -th segment, we run the EnKF using  $\hat{\theta}^{(k-1)}$  to obtain the initial analysis ensembles for every time point  $t$  in the  $k$ -th segment  $\mathcal{T}_k$  as in (5). We denote  $\mathbf{X}_{t,j}^a(\hat{\theta}^{(k-1)}|\hat{\theta}^{(k-1)}) = \mathbf{X}_{t,j}^a(\hat{K}_t(\hat{\theta}^{(k-1)}), \hat{\theta}^{(k-1)})$  for notation simplicity, where the two arguments of  $\theta$  in  $\mathbf{X}_{t,j}^a(\theta_1|\theta_2)$  indicate that  $\theta_1$  is used to run the model operator  $\mathcal{M}_{t-1}(\cdot, \theta_1)$  at  $t-1$  and  $\theta_2$  is used to estimate the Kalman gain matrix. The initial forecast ensembles are obtained by  $\{\mathbf{X}_{t,j}^f(\hat{\theta}^{(k-1)}) = \mathcal{M}_t(\mathbf{X}_{t-1,j}^a(\hat{\theta}^{(k-1)}|\hat{\theta}^{(k-1)}), \hat{\theta}^{(k-1)})\}_{j=1}^n$ , where  $\hat{\theta}^{(k-1)}$  is used in  $\mathcal{M}_t(\cdot, \theta)$  for all  $t \in \mathcal{T}_k$ . The Kalman gain matrices  $\hat{K}_t(\hat{\theta}^{(k-1)}) = \hat{\Sigma}_t^f(\hat{\theta}^{(k-1)})H_t'(H_t\hat{\Sigma}_t^f(\hat{\theta}^{(k-1)})H_t' + R_t)^{-1}$  for  $\mathcal{T}_k$ , where  $\hat{\Sigma}_t^f(\hat{\theta}^{(k-1)})$  is the sample covariance of the forecast ensembles  $\{\mathbf{X}_{t,j}^f(\hat{\theta}^{(k-1)})\}_{j=1}^n$ .

By keeping the Kalman gain matrices fixed at the parameter value  $\hat{\theta}^{(k-1)}$  estimated in the previous segment over the  $k$ -th segment, we find an updated estimate of  $\theta^*$  by evaluating forecast states in the  $k$ -th segment with changing values of  $\theta$ . Specifically, given  $\{\hat{K}_t(\theta^{(k-1)})\}_{t \in \mathcal{T}_k}$  and using the analysis ensembles as the end of the  $(k-1)$ -th segment as the initial states for the  $k$ -th segment, we obtain the analysis ensembles at a  $\theta$  and  $t \in \mathcal{T}_k$  as

$$\mathbf{X}_{t,j}^a(\theta|\hat{\theta}^{(k-1)}) = \mathbf{X}_{t,j}^f(\theta) + \hat{K}_t(\hat{\theta}^{(k-1)})(\mathbf{Y}_t - H_t\mathbf{X}_{t,j}^f(\theta) + e_{t,j}) \quad (6)$$

for  $j = 1, \dots, n$ , where  $\mathbf{X}_{t,j}^f(\theta) = \mathcal{M}_t(\mathbf{X}_{t-1,j}^a(\theta|\hat{\theta}^{(k-1)}), \theta)$ ,  $e_{t,j} \sim \mathcal{N}(0, R_t)$  and  $\mathcal{N}(0, R_t)$  denotes the normal distribution with mean zero and covariance  $R_t$ . Different from  $\mathbf{X}_{t,j}^a(\hat{K}_t(\theta), \theta)$  in (5) and the initial analysis ensembles  $\{\mathbf{X}_{t,j}^a(\hat{\theta}^{(k-1)}|\hat{\theta}^{(k-1)})\}$ , the sequentially updated analysis ensembles in (6) use different values of  $\theta$  for obtaining the Kalman gain matrices and computing the model forecasts, as reflected by the notation  $\mathbf{X}_{t,j}^a(\theta|\hat{\theta}^{(k-1)})$ . This allows us to find a  $\theta$  that minimizes weighted a distance between  $\mathbf{Y}_t$  and  $H_t\mathbf{X}_{t,j}^f(\theta)$  for  $t \in \mathcal{T}_k$  without re-computing the Kalman gain matrix at each time point of the  $k$ -th segment. We consider the following objective function in the  $k$ -th segment  $\mathcal{T}_k$  to minimize with respect to  $\theta$ :

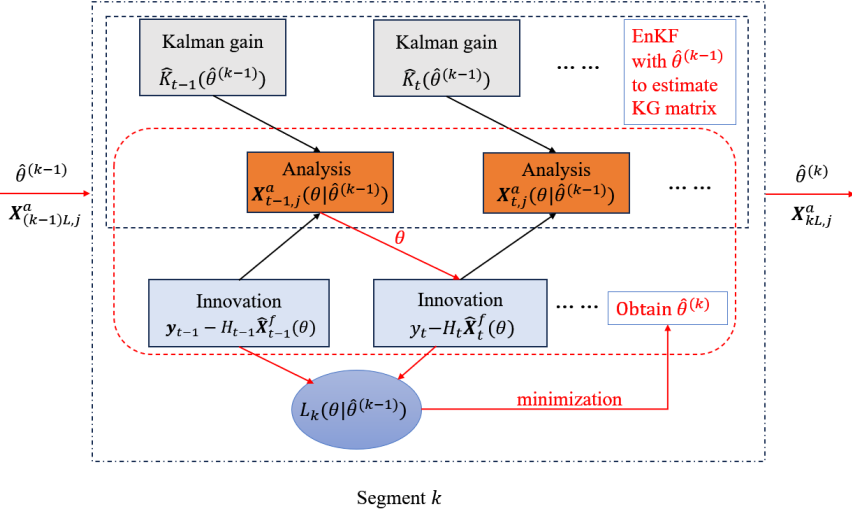
$$L_k(\theta|\hat{\theta}^{(k-1)}) := \sum_{t \in \mathcal{T}_k} (\mathbf{Y}_t - H_t\hat{\mathbf{X}}_t^f(\theta))' W_t(\hat{\theta}^{(k-1)}) (\mathbf{Y}_t - H_t\hat{\mathbf{X}}_t^f(\theta)), \quad (7)$$

where  $W_t(\hat{\theta}^{(k-1)}) = (H_t\hat{\Sigma}_t^f(\hat{\theta}^{(k-1)})H_t' + R_t)^{-1}$  and  $\hat{\mathbf{X}}_t^f(\theta) = n^{-1} \sum_{j=1}^n \mathbf{X}_{t,j}^f(\theta)$ . The weighting matrix  $W_t(\hat{\theta}^{(k-1)})$  is introduced to improve the efficiency of parameter estimation in the presence of heteroscedastic forecast errors. See the detailed explanation in the following paragraphs and the theoretical comparison with the NLS estimator without weighting in Theorem 4. Thus, we can obtain an updated estimate  $\hat{\theta}^{(k)}$  of  $\theta$  in the  $k$ -th segment as

$$\hat{\theta}^{(k)} = \arg \min_{\theta \in \Theta} L_k(\theta|\hat{\theta}^{(k-1)}). \quad (8)$$

In evaluating the objective function, obtaining the inverse of the  $q \times q$  matrix  $(H_t\hat{\Sigma}_t^f(\hat{\theta}^{(k-1)})H_t' + R_t)$  can be computationally realized by using the Sherman–Morrison–Woodbury formula if  $R_t$  is a spherical matrix [28]. Hence, the objective function, despite being weighted by  $W_t(\hat{\theta}^{(k-1)})$ , does not incur additional computational cost compared to the standard EnKF. Figure 1 illustrates the main steps of the proposed sequential updating and data assimilation process in the  $k$ -th segment.

It is noted that the NLS estimator minimizes the squared Euclidean distance  $\|\mathbf{Y}_t - H_t\hat{\mathbf{X}}_t^f(\theta)\|^2$  between  $\mathbf{Y}_t$  and  $\hat{\mathbf{X}}_t^f(\theta)$  without weighting, which is used in estimating the parameter for state-space models. However, the sequentially updating procedure was not developed. Evaluating the NLS loss function over the entire time period requires repeated re-runs of the whole data assimilation procedure, which is computationally heavy. Furthermore, even if the



**FIGURE 1** The entire time span is divided into  $L$  segments, and parameter estimation is performed within each segment using a machine learning-based optimization scheme introduced in Section 3.2. This plot shows the workflow for the  $k$ -th segment  $\mathcal{T}_k$ , where the parameter  $\hat{\theta}^{(k-1)}$  from the previous segment is used in an initial EnKF procedure to compute the Kalman gain  $\hat{K}_t(\hat{\theta}^{(k-1)})$  for  $t \in \mathcal{T}_k$ . Keeping those Kalman gain matrices unchanged, the analysis states are calculated with varying  $\theta$  in the model  $\mathcal{M}_t(\cdot, \theta)$  to evaluate the forecast error at each time step. The objective function  $L_k(\theta|\hat{\theta}^{(k-1)})$  of the weighted forecast errors is minimized to obtain the updated estimate  $\hat{\theta}^{(k)}$ . The red arrow indicates that the unknown parameter  $\theta$  is propagated through the model, enabling a sequential update of parameters and states.

183 proposed sequential estimation procedure is applied, the NLS estimator

$$\hat{\theta}_{\text{NLS}}^{(k)} = \arg \min_{\theta \in \Theta} \sum_{t \in \mathcal{T}_k} \|\mathbf{Y}_t - H_t \hat{\mathbf{X}}_t^f(\theta)\|^2 \quad (9)$$

184 without weighting on the  $k$ -th segment is not statistically efficient due to the heteroscedastic variances and correla-  
 185 tions of the forecast errors  $\mathbf{Y}_t - H_t \hat{\mathbf{X}}_t^f(\theta)$ . Namely,  $\hat{\theta}_{\text{NLS}}^{(k)}$  does not achieve the lowest possible estimation variance as  
 186 the sample size increases.

187 Note that  $\text{Var}(\mathbf{Y}_t | \mathbf{Y}_{1:t-1}) = H_t \Sigma_t^f(\theta^*) H_t' + R_t = (W_t^*)^{-1}$  if  $\Sigma_t^f(\theta^*) = \text{Var}(\mathbf{X}_t | \mathbf{Y}_{1:t-1})$ . The application of the  
 188 weighting matrix  $W_t(\hat{\theta}^{(k-1)})$  in our loss function  $L_k(\theta|\hat{\theta}^{(k-1)})$  leads to whitened forecast errors. This is in the same  
 189 spirit as the weighted least squares estimation in linear regression. Theorem 4 shows that the proposed estimator  
 190  $\hat{\theta}^{(k)}$  in (8) is more efficient than the NLS estimator  $\hat{\theta}_{\text{NLS}}^{(k)}$  under suitable conditions. In practical terms,  $\hat{\theta}_{\text{NLS}}^{(k)}$  tends to  
 191 be more sensitive to the variability of forecast errors and thus may produce a less stable parameter estimate. The  
 192 proposed weighted least squares approach enhances the stability of parameter estimation by weighting the forecast  
 193 errors, and thus, reduces the asymptotic variance of the estimate.

### 3.2 | Hierarchical expected optimization algorithm

Because the physical model  $\mathcal{M}_t(\cdot, \boldsymbol{\theta})$  is computationally expensive, particularly for high-resolution geophysical models, the optimization problem in (8) effectively becomes a black-box optimization task with substantial computational cost. To reduce the number of computationally expensive forward model evaluations, we adopt the hierarchical Bayesian optimization framework in [24], which uses a hierarchical Bayesian surrogate model to approximate the loss function  $L_k(\boldsymbol{\theta}|\hat{\boldsymbol{\theta}}^{(k-1)})$ . At each iteration, instead of numerically computing the derivative of  $\mathcal{M}_t(\cdot, \boldsymbol{\theta})$  with respect to  $\boldsymbol{\theta}$ , our method selects the most informative parameter value by maximizing an expected improvement function, called hierarchical expected improvement (HEI). Compared to the Bayesian optimization procedure using expected improvement in [23], the proposed HEI procedure provides a more efficient and stable search for  $\hat{\boldsymbol{\theta}}^{(k)}$  in (8) by balancing the exploitation of the surrogate model and exploration of the areas with high prediction uncertainty.

Let  $\{g(\boldsymbol{\theta})\} \sim \mathcal{GP}(0, \sigma^2 V_\tau)$  for  $\boldsymbol{\theta} \in \Theta$  denote a zero-mean Gaussian process (GP) on  $\boldsymbol{\theta}$  with marginal variance  $\sigma^2$  and correlation function  $V_\tau(\cdot, \cdot)$ , where  $V_\tau(\boldsymbol{\theta}, \boldsymbol{\theta}) = 1$  for any value of  $\boldsymbol{\theta}$  and  $\tau$  is the parameter. The correlation function  $V_\tau$  encodes the dependence of this process over  $\boldsymbol{\theta} \in \Theta$ . For this GP,  $(g(\boldsymbol{\theta}_1), \dots, g(\boldsymbol{\theta}_m))^T$  follows a multivariate normal distribution with mean zero and covariance  $\sigma^2 (V_\tau(\boldsymbol{\theta}_{i_1}, \boldsymbol{\theta}_{i_2}))_{m \times m}$  for any  $\boldsymbol{\theta}_1, \dots, \boldsymbol{\theta}_m \in \Theta$ . To construct a surrogate model for the objective function  $L_k(\boldsymbol{\theta}|\hat{\boldsymbol{\theta}}^{(k-1)})$ , we introduce a hierarchical GP process model  $\tilde{L}_k(\boldsymbol{\theta}; \boldsymbol{\beta}_k, \sigma^2, \tau)$  as

$$\tilde{L}_k(\boldsymbol{\theta}; \boldsymbol{\beta}_k, \sigma^2, \tau) \sim \mathcal{GP}(\mathbf{b}(\boldsymbol{\theta})' \boldsymbol{\beta}_k, \sigma^2 V_\tau), \quad (10)$$

$$\boldsymbol{\beta}_k \propto \mathbf{1} \text{ and } \sigma^2 \sim \text{IG}(\alpha_{k,1}, \alpha_{k,2}), \quad (11)$$

where  $\mathbf{b}(\boldsymbol{\theta}) = (\mathbf{b}_1(\boldsymbol{\theta}), \dots, \mathbf{b}_r(\boldsymbol{\theta}))'$  are  $r$  basis functions,  $\boldsymbol{\beta}_k \in \mathbb{R}^r$  are their coefficients,  $\boldsymbol{\beta}_k \propto \mathbf{1}$  denotes a flat (non-informative) prior, meaning that all values of  $\boldsymbol{\beta}_k$  are assigned equal prior probability, and  $\text{IG}(\alpha_{k,1}, \alpha_{k,2})$  represents the inverse Gamma distribution with parameters  $\alpha_{k,1}, \alpha_{k,2}$ . Under this hierarchical Bayesian model in (10)–(11), the mean  $\mathbf{b}(\boldsymbol{\theta})' \boldsymbol{\beta}_k$  of the GP fits the overall trend of the objective function  $L_k(\boldsymbol{\theta}|\hat{\boldsymbol{\theta}}^{(k-1)})$ , while its variance part  $\sigma^2 V_\tau$  fits the local structures and provides a credible interval for the fitted curve. Here, the prior distributions specified in (11) are designed to increase the error variance  $\sigma^2$  of the model, and hence, to widen the length of the credible intervals. This helps faster exploration of the objective function  $L_k(\boldsymbol{\theta}|\hat{\boldsymbol{\theta}}^{(k-1)})$  over the range of  $\boldsymbol{\theta}$ .

The GP model in (10) with basis function regression  $\boldsymbol{\phi}(\boldsymbol{\theta})' \boldsymbol{\beta}$  is a common statistical model used in atmospheric sciences [29, 30, 31]. The covariance function  $V_\tau$  encodes spatial-temporal dependence, analogous to error covariance structures in numerical weather prediction and data assimilation [3]. The flat prior  $\boldsymbol{\beta}_k \propto \mathbf{1}$  indicates that no strong prior knowledge is imposed on the coefficients, allowing the data to determine the deterministic trend. Similar approaches have been successfully applied in climate and weather prediction [32, 33].

The HEI optimization procedure is conducted in an iterative manner. Let  $\boldsymbol{\theta}_{k,1}, \dots, \boldsymbol{\theta}_{k,m}$  be  $m$  pre-specified design points, where we evaluate the objective function  $L_k(\boldsymbol{\theta}|\hat{\boldsymbol{\theta}}^{(k-1)})$ . To build a surrogate model, we assume the GP  $\tilde{L}_k(\boldsymbol{\theta}; \boldsymbol{\beta}_k, \sigma^2, \tau)$  takes the same value as  $L_k(\boldsymbol{\theta}|\hat{\boldsymbol{\theta}}^{(k-1)})$  on those design points  $\{\boldsymbol{\theta}_{k,i}\}_{i=1}^m$ . That is, we use the set

$$\mathcal{D}_{k,0} = \{\boldsymbol{\theta}_{k,i}, L_k(\boldsymbol{\theta}_{k,i}|\hat{\boldsymbol{\theta}}^{(k-1)})\}_{i=1}^m,$$

as the training data for the GP, where  $\tilde{L}_k(\boldsymbol{\theta}_{k,i}; \boldsymbol{\beta}_k, \sigma^2, \tau) = L_k(\boldsymbol{\theta}_{k,i}|\hat{\boldsymbol{\theta}}^{(k-1)})$  on  $\mathcal{D}_{k,0}$ . The posterior distribution of  $\tilde{L}_k(\boldsymbol{\theta}; \boldsymbol{\beta}_k, \sigma^2, \tau)$  given  $\mathcal{D}_{k,0}$  provides the predictive mean and variance, which are considered as a surrogate function of  $L_k(\boldsymbol{\theta}|\hat{\boldsymbol{\theta}}^{(k-1)})$  and its uncertainty quantification, respectively. This posterior distribution is then utilized to construct the HEI updating function for Bayesian optimization.

Let  $T_\nu(a, b^2)$  be the location-scale  $t$  distribution with the location parameter  $a$ , scale parameter  $b$  and degree of

226 freedom  $\nu$ , where  $T_\nu(a, b^2) = a + bt_\nu$ , and  $t_\nu$  is the Student's  $t$  distribution with degree of freedom  $\nu$ . Given the training  
 227 set  $\mathcal{D}_{k,0}$  and the parameter  $\tau$  in the correlation function  $V_\tau$ , let  $\hat{\boldsymbol{\beta}}_k$  and  $\hat{\sigma}_k^2$  be the maximum likelihood estimators  
 228 (MLEs) of  $\boldsymbol{\beta}_k$  and  $\sigma^2$  for the model in (10), respectively. Let  $\hat{\sigma}_k^2 = \hat{\alpha}_{k,2}/\hat{\alpha}_{k,1}$ , where  $\hat{\alpha}_{k,1} = \alpha_{k,1} + (m-r)/2$  and  
 229  $\hat{\alpha}_{k,2} = \alpha_{k,2} + m\hat{\sigma}_k^2/2$ . Note that all those estimates depend on the training set  $\mathcal{D}_{k,0}$  in the updating procedure of HEI.  
 230 Here, we ignore this dependence in those notations for simplicity of the presentation. Under the hierarchical model  
 231 in (10)-(11), the marginal posterior distribution  $\tilde{L}_k(\boldsymbol{\theta})|\mathcal{D}_{k,0}$  conditioned on the training set  $\mathcal{D}_{k,0}$  that integrates over  
 232 the prior distributions of  $\boldsymbol{\beta}_k$  and  $\sigma^2$  follows a location-scale  $t$ -distribution as

$$\tilde{L}_k(\boldsymbol{\theta})|\mathcal{D}_{k,0} \sim T_{2\hat{\alpha}_{k,1}}(\hat{L}_k(\boldsymbol{\theta}), \hat{\sigma}_k^2 s_k^2(\boldsymbol{\theta})), \quad (12)$$

233 where  $\hat{L}_k(\boldsymbol{\theta}) = \mathbf{b}(\boldsymbol{\theta})'\hat{\boldsymbol{\beta}}_k + \mathbf{v}_k(\boldsymbol{\theta})'\mathbf{V}_k^{-1}(\mathbf{L}_k - \mathbf{P}_k\hat{\boldsymbol{\beta}}_k)$  and  $s_k^2(\boldsymbol{\theta}) = 1 - \mathbf{v}_k(\boldsymbol{\theta})'\mathbf{V}_k^{-1}\mathbf{v}_k(\boldsymbol{\theta}) + \mathbf{h}_k(\boldsymbol{\theta})'(\mathbf{P}_k'\mathbf{V}_k^{-1}\mathbf{P}_k)^{-1}\mathbf{h}_k(\boldsymbol{\theta})$  with  
 234  $\mathbf{v}_k(\boldsymbol{\theta}) = (V_\tau(\boldsymbol{\theta}, \boldsymbol{\theta}_{k,1}), \dots, V_\tau(\boldsymbol{\theta}, \boldsymbol{\theta}_{k,m}))'$ ,  $\mathbf{V}_k = (V_\tau(\boldsymbol{\theta}_{k,i_1}, \boldsymbol{\theta}_{k,i_2}))_{i_1, i_2=1}^m$ ,  $\mathbf{L}_k = (L_k(\boldsymbol{\theta}_{k,1}|\hat{\boldsymbol{\theta}}^{(k-1)}), \dots, L_k(\boldsymbol{\theta}_{k,m}|\hat{\boldsymbol{\theta}}^{(k-1)}))'$ ,  
 235  $\mathbf{P}_k = (\mathbf{b}(\boldsymbol{\theta}_{k,1}), \dots, \mathbf{b}(\boldsymbol{\theta}_{k,m}))'$  and  $\mathbf{h}_k(\boldsymbol{\theta}) = \mathbf{b}(\boldsymbol{\theta}) - \mathbf{P}_k'\mathbf{V}_k^{-1}\mathbf{v}_k(\boldsymbol{\theta})$ .

236 The posterior distribution  $\tilde{L}_k(\boldsymbol{\theta})|\mathcal{D}_{k,0}$  serves as a probabilistic surrogate representation of the objective function  
 237  $L_k(\boldsymbol{\theta}|\hat{\boldsymbol{\theta}}^{(k-1)})$ , allowing uncertainty quantification through the GP covariance function and the prior distribution on  
 238  $\sigma^2$ . In the posterior distribution,  $\hat{L}_k(\boldsymbol{\theta})$  is the estimated conditional mean of  $\tilde{L}_k(\boldsymbol{\theta})$  given the set  $\mathcal{D}_{k,0}$  of training data.  
 239 Its first term,  $\mathbf{b}(\boldsymbol{\theta})'\hat{\boldsymbol{\beta}}_k$ , gives the overall trend prediction by the regression basis functions, and its the second term  
 240 provides a covariance-based correction, which exactly reproduces the same values of  $L_k(\boldsymbol{\theta}|\hat{\boldsymbol{\theta}}^{(k-1)})$  at the evaluation  
 241 points  $\{\boldsymbol{\theta}_{k,i}\}_{i=1}^m$ , while smoothly interpolates in between. Figure 2 evaluates the fitted curve of  $\hat{L}_k(\boldsymbol{\theta})$  in a Lorenz-96  
 242 system. The detailed simulation settings are introduced in Section 4.1. From this figure, the fitted curve closely aligned  
 243 with the true objective function and the leave one out cross validation errors were small over all simulation repetitions,  
 244 demonstrating that the GP model offers a good surrogate for our loss function of  $\boldsymbol{\theta}$ .

245 Note that the expression of the estimated conditional mean  $\hat{L}_k(\boldsymbol{\theta})$  is the same as that of the universal kriging  
 246 or Bayesian optimization based on a Gaussian process with a fixed covariance function [34]. Universal kriging and  
 247 Bayesian optimization are widely used in geoscience to develop statistical emulators for physical model operators. See,  
 248 for example [23] and [31]. However, our posterior variance is different due to adding a prior on the error variance  $\sigma^2$   
 249 in (11). This results in a  $t$ -distribution which has a larger variance than the normal distributed posterior in the universal  
 250 kriging and Bayesian optimization. A larger posterior variance helps for faster exploring outside the training points  
 251  $\{\boldsymbol{\theta}_{k,i}\}_{i=1}^m$  in the optimization procedure introduced in the following.

Using the posterior distribution  $\tilde{L}_k(\boldsymbol{\theta})|\mathcal{D}_{k,0}$  in (12), we now introduce an iterative procedure to minimize the  
 objective function  $L_k(\boldsymbol{\theta}|\hat{\boldsymbol{\theta}}^{(k-1)})$ . Let  $L_{k,\min} = \min\{L_k(\boldsymbol{\theta}_{k,i}|\hat{\boldsymbol{\theta}}^{(k-1)}) : i = 1, \dots, m\}$ . Based on the current training set  
 $\mathcal{D}_{k,0}$ , we aim to find a new evaluation point of  $\boldsymbol{\theta}$  that balances the exploitation and exploration. Here, exploitation  
 means finding the minimum value of the current fitted curve  $\hat{L}_k(\boldsymbol{\theta})$ , and exploration means finding the value of  $\boldsymbol{\theta}$   
 outside  $\mathcal{D}_{k,0}$  where the fitted curve is most uncertain, namely, the locations with large posterior variances. This can  
 be achieved by maximizing the expected improvement function, defined as

$$\text{HEI}_0(\boldsymbol{\theta}) = \mathbb{E}_{\tilde{L}_k(\boldsymbol{\theta})|\mathcal{D}_{k,0}}\{L_{k,\min} - \tilde{L}_k(\boldsymbol{\theta})\}_+,$$

252 where  $\{x\}_+ = \max\{x, 0\}$  is the positive truncation function. The improvement function  $\{L_{k,\min} - \tilde{L}_k(\boldsymbol{\theta})\}_+$  measures  
 253 the potential gain at  $\boldsymbol{\theta}$  if  $\tilde{L}_k(\boldsymbol{\theta})$  is smaller than the current minimum value  $L_{k,\min}$  in the set  $\mathcal{D}_{k,0}$ . Based on the result  
 254 in (12), the HEI updating function can be calculated as

$$\text{HEI}_0(\boldsymbol{\theta}) = I_k(\boldsymbol{\theta})\Phi_{2\hat{\alpha}_{k,1}}\left(\frac{I_k(\boldsymbol{\theta})}{\hat{\sigma}_k s_k(\boldsymbol{\theta})}\right) + \nu_k \hat{\sigma} s_k(\boldsymbol{\theta})\phi_{2(\hat{\alpha}_{k,1}-1)}\left(\frac{I_k(\boldsymbol{\theta})}{\nu_k \hat{\sigma} s_k(\boldsymbol{\theta})}\right), \quad (13)$$



**FIGURE 2** Fitted curve (left) of the loss function for the Lorenz-96 system with  $T = 1000$ ,  $p = 50$  and  $q = 40$ , where the parameter of interest is the forcing  $\theta_F$  of this system and its true value is  $\theta_F^* = 8$ . For demonstration of the Gaussian process fitting, only one segment of data assimilation is considered ( $K = 1$ ). The objective function (red dots)  $L_1(\theta_F | \theta_F^{(0)})$  is derived using the standard ensemble Kalman filter with ensemble size  $n = 50$  and initial value  $\theta_F^{(0)} = 11$ . The posterior mean (black line)  $\hat{L}_k(\theta)$  is computed under the model in (10) using the training data (blue points) with size  $m = 20$ , where the design points  $\{\theta_{1,i}\}_{i=1}^m$  are uniformly selected from (4, 12). The shaded area represents the 95% credible interval. The boxplot (right) shows the distribution of the relative RMSE (RMSE normalized by the true values of objective function) based on 200 replications.

255 where  $I_k(\theta) = L_{k,\min} - \hat{L}_k(\theta)$ ,  $v_k = \{\tilde{\alpha}_{k,1}/(\tilde{\alpha}_{k,1} - 1)\}^{1/2}$ , and  $\Phi_a(\cdot)$  and  $\phi_a(\cdot)$  represent the cumulative distribution  
 256 function and the probability density function of a  $t$ -distribution with degree of freedom  $a$ , respectively. The hyperpa-  
 257 rameters  $\tau$ ,  $\alpha_{k,1}$  and  $\alpha_{k,2}$  in the model (10)–(11) can be estimated empirically, and their estimates are plugged-in to  
 258 evaluate the HEI function  $\text{HEI}_0(\theta)$ .

259 Given the expression of  $\text{HEI}_0(\theta)$  in (13), we obtain the next evaluation point  $\theta_{k,m+1}$  by maximizing  $\text{HEI}_0(\theta)$ ,  
 260 namely,

$$\theta_{k,m+1} = \underset{\theta \in \Theta}{\operatorname{argmax}} \text{HEI}_0(\theta). \quad (14)$$

261 The training set is then updated as  $\mathcal{D}_{k,1} = \mathcal{D}_{k,0} \cup \{(\theta_{k,m+1}, L_k(\theta_{k,m+1} | \hat{\theta}^{(k-1)}))\}$ . Based on this updated training dataset  
 262  $\mathcal{D}_{k,1}$ , we compute a new fitted curve and the posterior distribution  $\tilde{L}_k(\theta) | \mathcal{D}_{k,1}$  as (12). Then, we evaluate the HEI  
 263 function  $\text{HEI}_1(\theta) = \mathbb{E}_{\tilde{L}_k(\theta) | \mathcal{D}_{k,1}} \{L_{k,\min} - \tilde{L}_k(\theta)\}_+$  as (13) and calculate its maximum to obtain the subsequent eval-  
 264 uation point  $\theta_{k,m+2}$ . This iterative HEI procedure continues until the number of iterations reaches a prespecified  
 265 threshold  $s_{\max}$ , which produces a sequence of evaluation points  $\{\theta_{k,i} : i = 1, \dots, m + s_{\max}\}$  for the objective function  
 266 in the  $k$ -th segment. The estimate  $\hat{\theta}^{(k)}$  in (8) is computed as

$$\hat{\theta}^{(k)} = \underset{\theta \in \{\theta_{k,i} : i=1, \dots, m+s_{\max}\}}{\operatorname{argmin}} L_k(\theta | \hat{\theta}^{(k-1)}), \quad (15)$$

267 which is the minimum value of the objective function in the set of evaluation points selected by the HEI procedure.

268 **Remark** The HEI algorithm is a computationally efficient global optimization method in machine learning which avoids  
 269 intensive computation of the derivatives of the model operator  $\mathcal{M}_t(\cdot, \boldsymbol{\theta})$ . The HEI function in (13) is a combination of  
 270 two parts, where the first part  $L_{k,\min} - \hat{L}_k(\boldsymbol{\theta})$  reaches the maximum at the minimum value of the fitted curve  $L_k(\boldsymbol{\theta})$ ,  
 271 and the second part is determined by the posterior standard deviation  $\bar{\sigma}_{s_k}(\boldsymbol{\theta})$ , which penalizes the fitted value by  
 272 its uncertainty quantification. Specifically, in the expression of  $s_k(\boldsymbol{\theta})$ , a large value of  $h_k(\boldsymbol{\theta})$  implies deviation of  $\boldsymbol{\theta}$   
 273 from the training data, which leads to a large posterior variance. This means that  $s_k(\boldsymbol{\theta})$  could be large for the points  
 274 that are far away from the training data, as demonstrated in Figure 2. Therefore, searching a new evaluation point by  
 275 maximizing the HEI function in (14) intrinsically considers an exploitation-exploration trade-off: minimizing the fitted  
 276 curve based on the training data while exploring the regions outside the training set with large prediction uncertainty.  
 277 The posterior distribution of  $\hat{L}_k(\boldsymbol{\theta})$  integrates over the prior uncertainty of  $\sigma^2$  in (11), resulting a heavier-tailed  $t$ -  
 278 distribution rather than the normal distribution used in universal kriging. This heavier tail allows HEI to put more  
 279 weight on exploration, so that it may find the global optimizer using fewer steps than Bayesian optimization.

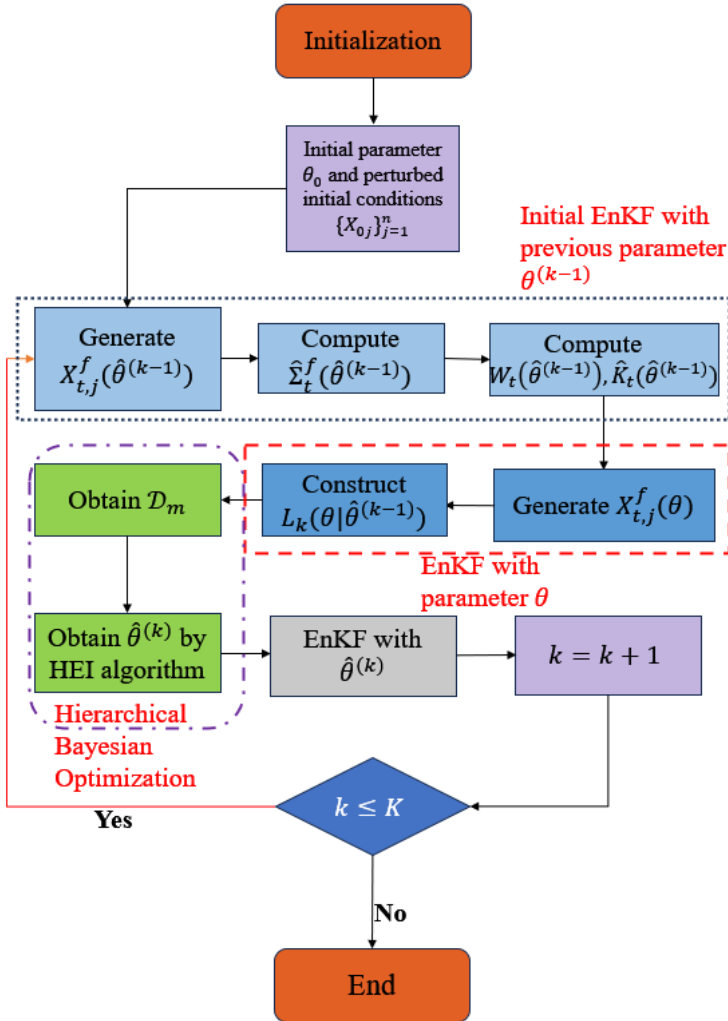
### 280 3.3 | Sequential algorithm for parameter estimation

281 Based on the sequential objective function introduced in Section 3.1 and the HEI algorithm in Section 3.2, we pro-  
 282 pose an efficient sequential machine learning framework for joint data assimilation and parameter estimation under  
 283 the models in (1) and (2). The time domain is divided into  $K$  segments  $\{\mathcal{T}_k\}_{k=1}^K$ , and the parameter estimation is per-  
 284 formed once per segment based on the accumulated observations therein. This strategy retains the adaptivity of the  
 285 augmented EnKF which updates the parameter at every assimilation step, while also enjoying the statistical properties  
 286 of minimizing a global loss function across the entire domain of data assimilation. Given the observations  $\mathbf{Y}_1, \dots, \mathbf{Y}_T$   
 287 and the initial condition  $\mathbf{X}_0$  and the covariance  $\Sigma_0$  of the background field, we develop the following algorithm for  
 288 data assimilation with parameter estimation.

- 289 • 1. **(Initialization at  $t = 0$ )** Give the initial condition  $\mathbf{X}_0$  and the background error covariance  $\Sigma_0$ , generate  $n$   
 290 ensemble members  $\{X_{0,1}, \dots, X_{0,n}\}$  from the distribution  $\mathcal{N}(\mathbf{X}_0, \Sigma_0)$  of the background field; give an initial value  
 291  $\boldsymbol{\theta}^{(0)}$  of  $\boldsymbol{\theta}$ .
- 292 • 2. **(Initial EnKF with  $\hat{\boldsymbol{\theta}}^{(k-1)}$ )** In the  $k$ -th time segment, carry forward the analysis ensemble at the end of  $\mathcal{T}_{k-1}$  as  
 293 the initial state and fix the parameter at  $\hat{\boldsymbol{\theta}}^{(k-1)}$ . Run EnKF over  $t \in \mathcal{T}_k$  using the physical model  $\mathcal{M}_t(\cdot, \hat{\boldsymbol{\theta}}^{(k-1)})$  and  
 294 obtain the initial forecast ensembles  $\{\mathbf{X}_{t,j}^f(\hat{\boldsymbol{\theta}}^{(k-1)})\}$ , then calculate their sample covariance  $\hat{\Sigma}_t^f(\hat{\boldsymbol{\theta}}^{(k-1)})$ . Using  
 295  $\hat{\Sigma}_t^f(\hat{\boldsymbol{\theta}}^{(k-1)})$ , compute the weighting matrix  $W_t(\hat{\boldsymbol{\theta}}^{(k-1)})$  and Kalman gain  $K_t(\hat{\boldsymbol{\theta}}^{(k-1)})$  as introduced in Section 3.1.
- 296 • 3. **(Loss function in  $\mathcal{T}_k$ )** With  $\{W_t(\hat{\boldsymbol{\theta}}^{(k-1)})\}$  and  $\{\hat{K}_t(\hat{\boldsymbol{\theta}}^{(k-1)})\}$  held fixed for  $t \in \mathcal{T}_k$ , we evaluate the forecast and  
 297 analysis ensembles for a candidate value  $\boldsymbol{\theta}$  over  $\mathcal{T}_k$  and construct the objective function  $L_k(\boldsymbol{\theta} \mid \hat{\boldsymbol{\theta}}^{(k-1)})$  defined  
 298 in (7). Generate  $m$  design points  $\boldsymbol{\theta}_{k,1}, \dots, \boldsymbol{\theta}_{k,m}$  uniformly sampled from  $\Theta$  by the space-filling design (Chapter 5 of  
 299 [35]). Obtain  $m$  evaluations of  $L_k(\boldsymbol{\theta} \mid \hat{\boldsymbol{\theta}}^{(k-1)})$  at the design points  $\boldsymbol{\theta}_{k,1}, \dots, \boldsymbol{\theta}_{k,m}$ .
- 300 • 4. **(HEI optimization)** Minimize the objective function  $L_k(\boldsymbol{\theta} \mid \hat{\boldsymbol{\theta}}^{(k-1)})$  and obtain an update  $\hat{\boldsymbol{\theta}}^k$  via the HEI algo-  
 301 rithm as follows:  
 302 i) construct the HEI updating function  $\text{HEI}(\boldsymbol{\theta})$  by (13);  
 303 ii) perform  $s_{\max}$  iterations of HEI, yielding  $s_{\max}$  new evaluation points  $\{\boldsymbol{\theta}_{k,m+1}, \dots, \boldsymbol{\theta}_{k,m+s_{\max}}\}$ .  
 304 iii) obtain the updated parameter  $\hat{\boldsymbol{\theta}}^{(k)}$  from all evaluation points, including the initial  $m$  design points, by (15).
- 305 • 5. **(EnKF with  $\hat{\boldsymbol{\theta}}^{(k)}$ )** Based on the updated parameter  $\hat{\boldsymbol{\theta}}^{(k)}$ , compute the updated forecast states  $\mathbf{X}_{t,j}^f(\hat{\boldsymbol{\theta}}^{(k)})$  and

- 306 analysis states  $X_{t,j}^a(\hat{\theta}^{(k)}|\hat{\theta}^{(k-1)})$  for  $j = 1, \dots, n$ .
- 307 • 6. Iterate Steps 2-5 for all the time segments  $\mathcal{T}_1, \dots, \mathcal{T}_K$ .

308 Figure 3 presents a flowchart to illustrate this sequential data assimilation algorithm with parameter estimation by  
 309 Bayesian optimization. Due to the rapid convergence of the hierarchical Bayesian optimization algorithm, the overall  
 310 computational cost can be kept low.



**FIGURE 3** Flowchart illustrating the sequential machine learning algorithm for data assimilation and parameter estimation by the hierarchical Bayesian optimization.

### 3.4 | Theoretical properties

In this subsection, we investigate the theoretical properties of the proposed estimator, which establishes its consistency, asymptotic normality, and efficiency under suitable regularity conditions. The proofs of all theorems can be found in the SM.

**Theorem 2** Under Assumptions S1'-S3' and A1-A5 in the SM,  $\hat{\boldsymbol{\theta}}^{(k)}$  converges to  $\boldsymbol{\theta}^*$  in probability for  $k = 1, \dots, K$  as  $n, L \rightarrow \infty$ .

Theorem 2 establishes the consistency of the proposed estimator of  $\boldsymbol{\theta}^*$  in each segment. Recall that  $\mathcal{N}(0, \Sigma)$  denotes the multivariate normal distribution with zero mean and covariance  $\Sigma$ . The following theorem shows the asymptotic normality of  $\hat{\boldsymbol{\theta}}^{(k)}$ .

**Theorem 3** Under Assumptions S1'-S3' and A1-A8 in the SM,  $\sqrt{L}(\hat{\boldsymbol{\theta}}^{(1)} - \boldsymbol{\theta}^*)$  converges to  $\mathcal{N}_d(0, U_0^{-1}U_{01}U_0^{-1})$  in distribution, and for  $k = 2, \dots, K$ ,  $\sqrt{L}(\hat{\boldsymbol{\theta}}^{(k)} - \boldsymbol{\theta}^*)$  converges to  $\mathcal{N}_d(0, U_k^{-1})$  in distribution, as  $n, L \rightarrow \infty$ , where

$$U_{01} = \lim_{L \rightarrow \infty} \frac{1}{L} \sum_{t=1}^L \mathbb{E} \left\{ \mathbb{E}_{\text{ens}} \{ J_t(\boldsymbol{\theta}^*) \}' H_t' (H_t \Sigma_t^f(\boldsymbol{\theta}^{(0)}) H_t' + R_t)^{-1} (H_t \Sigma_t^f H_t' + R_t) (H_t \Sigma_t^f(\boldsymbol{\theta}^{(0)}) H_t' + R_t)^{-1} H_t \mathbb{E}_{\text{ens}} \{ J_t(\boldsymbol{\theta}^*) \} \right\}$$

and

$$U_k = \lim_{L \rightarrow \infty} \frac{1}{L} \sum_{t \in \mathcal{J}_k} \mathbb{E} \left\{ \left\{ \mathbb{E}_{\text{ens}} J_t(\boldsymbol{\theta}^*) \right\}' H_t' W_t^* H_t \left\{ \mathbb{E}_{\text{ens}} J_t(\boldsymbol{\theta}^*) \right\} \right\}, k = 1, \dots, K$$

with  $J_t(\boldsymbol{\theta}^*) = \frac{\partial}{\partial \boldsymbol{\theta}} \mathcal{M}_t(\mathbf{X}_{t-1}^a, \boldsymbol{\theta}^*)$

Theorem 3 derives the asymptotic variance of the proposed estimator  $\hat{\boldsymbol{\theta}}^{(k)}$  in the  $k$ -th segment and shows its asymptotic normality property. From this theorem, the convergence rate of  $\hat{\boldsymbol{\theta}}^{(k)}$  is at the order  $L^{-1/2}$ , meaning a larger  $L$  (longer segment) leads to a more accurate estimate of  $\boldsymbol{\theta}^*$ . However, this would increase the computation burden in that segment. The  $(1 - \alpha) \times 100\%$  confidence interval of the  $\ell$ -th component of  $\boldsymbol{\theta}$  can be constructed as

$$\hat{\theta}_\ell^{(k)} \pm z_{1-\alpha/2} \sqrt{L^{-1} (\hat{U}_k^{-1})_{\ell\ell}}$$

for  $\ell = 1, \dots, d$ , where

$$\hat{U}_k = \lim_{L \rightarrow \infty} \frac{1}{L} \sum_{t \in \mathcal{J}_k} \left\{ \frac{1}{n} \sum_{j=1}^n \frac{\partial}{\partial \boldsymbol{\theta}} \mathcal{M}_t(\mathbf{X}_{t-1}^a, \hat{\boldsymbol{\theta}}^{(k)}) \right\}' H_t' W_t(\hat{\boldsymbol{\theta}}^{(k)}) H_t \left\{ \frac{1}{n} \sum_{j=1}^n \frac{\partial}{\partial \boldsymbol{\theta}} \mathcal{M}_t(\mathbf{X}_{t-1}^a, \hat{\boldsymbol{\theta}}^{(k)}) \right\}$$

is an estimator of  $U_k$ ,  $\frac{\partial}{\partial \boldsymbol{\theta}} \mathcal{M}_t(\mathbf{X}_{t-1}^a, \hat{\boldsymbol{\theta}}^{(k)})$  can be obtained by numerical differentiation,  $(\hat{U}_k^{-1})_{\ell\ell}$  denotes the  $\ell$ -th diagonal element of the inverse matrix of  $\hat{U}_k$ , and  $z_{1-\alpha/2}$  refers to the  $(1 - \alpha/2)$  quantile of the standard normal distribution.

To demonstrate the advantage of the proposed estimator  $\hat{\boldsymbol{\theta}}^{(k)}$  using a weighted loss function in (7) over the non-linear least squares estimator  $\hat{\boldsymbol{\theta}}_{\text{NLS}}^{(k)}$  in (9) without weighting, we compare their asymptotic variances and theoretically show the proposed estimator is statistically more efficient.

331 **Theorem 4** Under the same Assumptions as Theorem 3, we have  $\sqrt{L}(\hat{\boldsymbol{\theta}}_{\text{NLS}}^{(k)} - \boldsymbol{\theta}^*)$  converges to  $\mathcal{N}_d(0, \tilde{U}_k^{-1} \tilde{U}_{0k} \tilde{U}_k^{-1})$  in distri-  
 332 bution, where the expressions of  $\tilde{U}_{0k}$  and  $\tilde{U}_k$  are given in formulae (S3) and (S5) in the SM. Furthermore,  $\tilde{U}_k^{-1} \tilde{U}_{0k} \tilde{U}_k^{-1} - U_k^{-1}$   
 333 is positive semidefinite for  $k = 2, \dots, K$  as  $n, L \rightarrow \infty$ .

334 Since  $\tilde{U}_k^{-1} \tilde{U}_{0k} \tilde{U}_k^{-1} - U_k^{-1}$  is positive semidefinite, the variances of the proposed estimator are smaller than those  
 335 of the NLS estimator as the ensemble size  $n$  and the length  $L$  of each segment grow large. This demonstrates that the  
 336 proposed method yields more precise parameter estimates with smaller asymptotic variance.

337 **Remark** Both the proposed estimator and the NLS estimator are constructed from the forecast ensembles. In contrast,  
 338 one may also consider an alternative estimator

$$\hat{\boldsymbol{\theta}}_a^{(k)} = \arg \min_{\boldsymbol{\theta} \in \Theta} \sum_{t \in \mathcal{I}_k} \|\mathbf{Y}_t - H_t \hat{\mathbf{X}}_t^a(\boldsymbol{\theta} | \hat{\boldsymbol{\theta}}^{(k-1)})\|^2. \quad (16)$$

339 based on the analysis ensembles, as the mean of the analysis ensembles is generally closer to the true state than that  
 340 of the forecast ensembles. However, our theoretical analysis in the SM shows that  $\hat{\boldsymbol{\theta}}_a^{(k)}$  is not necessarily superior to  
 341  $\hat{\boldsymbol{\theta}}_{\text{NLS}}^{(k)}$ . In fact, the estimator  $\hat{\boldsymbol{\theta}}_a^{(k)}$  can in fact be viewed as a special case of the weighted loss framework, where the  
 342 weighting matrix is specified as  $\tilde{W}_t = (I_q - \hat{K}_t(\hat{\boldsymbol{\theta}}^{(k-1)})H_t)'(I_q - \hat{K}_t(\hat{\boldsymbol{\theta}}^{(k-1)})H_t)$ . In the SM, we have shown that the  
 343 proposed estimator  $\hat{\boldsymbol{\theta}}^{(k)}$  is optimal within the class of estimators that minimize a weighted loss as (7). This means that  
 344 the proposed estimator also has higher precision than  $\hat{\boldsymbol{\theta}}_a^{(k)}$ .

## 345 4 | NUMERICAL STUDIES

346 To evaluate the performance of the proposed method, we conducted a series of numerical experiments using the  
 347 Lorenz'96 model and the two-dimensional shallow water equations. All simulations were implemented in R and were  
 348 based on 50 repeated runs. The proposed sequential data assimilation method (Seq\_HEI) was compared against  
 349 five existing approaches: the standard EnKF, EnKF with inflation (EnKF\_infla), augmented assimilation (Aug), aug-  
 350 mented assimilation with inflation (Aug\_infla), and the standard EnKF using the NLS method for parameter estimation  
 351 (EnKF\_NLS) without sequential data assimilation and updating, as proposed in [23].

### 352 4.1 | Simulations on Lorenz'96 model

353 We considered the Lorenz'96 model [36], a widely used nonlinear dynamical system that has been extensively studied  
 354 in the context of data assimilation. Let the true state variable at time  $t$  be denoted by  $\mathbf{X} = (X_1, \dots, X_p)'$ , which evolved  
 355 according to the following equation:

$$\frac{dX_i}{dt} = (X_{i+1} - X_{i-2})X_{i-1} - X_i + \theta_F, \quad i = 1, \dots, p, \quad (17)$$

356 where  $X_{-1} = X_{p-1}$ ,  $X_{-2} = X_{p-2}$  and  $\theta_F^* = 8$  was the true parameter. We solved the equation (17) by the fourth-order  
 357 Runge-Kutta time integration scheme with a time step of 0.05 non-dimensional time unit. We set the initial state  
 358  $\mathbf{X}_0 = (X_{01}, \dots, X_{0p})'$  with the  $j$ -th element  $X_{0i} = \theta_F^*$  for  $i \neq \lfloor p/2 \rfloor$ ,  $X_{0i} = \theta_F^* + 0.001$  for  $i = \lfloor p/2 \rfloor$ ,  $i = 1, \dots, p$ , and the  
 359 initial parameter  $\theta_F^{(0)} = 11$ . Recall that  $Q = (Q_{i_1 i_2})$  and  $R = (R_{i_1 i_2})$  are the the model error covariance and observation  
 360 error covariance in (1), respectively. We set  $Q_{i_1 i_2} = \sigma_1 \times 0.3^{|i_1 - i_2|}$  and  $R_{i_1 i_2} = \sigma_2 \times 0.1^{|i_1 - i_2|}$  for  $i_1, i_2 = 1, \dots, p$ . We  
 361 considered two settings for the simulation experiment. In the first setting, we fixed  $\sigma_1 = 0.1$  and let  $\sigma_2 = 0.2, 0.4, 0.8$ .

In the second setting, we fixed  $\sigma_2 = 0.2$  and allowed  $\sigma_1$  to take values of 0.2, 0.4, and 0.8.

For all methods, we set  $T = 2000$ ,  $p = 200, 400$ ,  $q = 0.8p$ , and the ensemble size  $n = 30, 60$ . For the proposed method, we divided the time period  $\{1, \dots, T\}$  into  $K = 20$  segments. For the standard EnKF and EnKF\_infla that don't involve parameter estimation, the initial value  $\theta_F^{(0)} = 11$  was set to obtain their prediction states. To compare the performance between different methods, we used the average RMSE over spatial locations for the differences between the observations and the analysis states, that is,

$$\text{RMSE}_t = \sqrt{\frac{1}{q} \left\| \mathbf{y}_t - \frac{1}{n} \sum_{j=1}^n H_t \mathbf{x}_{t,j}^a(\hat{\boldsymbol{\theta}}^{(k)} | \hat{\boldsymbol{\theta}}^{(k-1)}) \right\|^2}, \text{ for } t \in \mathcal{T}_k, k = 1, \dots, K. \quad (18)$$

To carry out the proposed sequential algorithm introduced in Section 3.3 based on Lorenz'96 model, we used orthogonal polynomial basis functions to fit the loss function under model (10), and the polynomial order was determined according to the Bayesian Information Criterion rule. Due to space limits, some of the simulation results were provided in the SM.

Figure 4 together with Figures S1–S2 in the SM report the RMSEs under observation variances  $\sigma_2 = 0.2, 0.4, 0.8$ , respectively. From those figures, Seq\_HEI consistently achieved the lowest RMSE. At  $\sigma_2 = 0.2$  with  $n = 30$ , its RMSE remained below 1.5, while standard EnKF and EnKF\_infla exceeded 2.5. Increasing  $n$  to 60 reduced Seq\_HEI's RMSE to around 1.0, whereas other methods improved only marginally. As  $\sigma_2$  grew, all methods' errors increased, yet Seq\_HEI retained its lead, especially at small ensemble sizes, demonstrating superior performance under limited ensembles.

Figure 5 and Figures S3–S4 in the SM report the RMSEs under model variance  $\sigma_1 = 0.2, 0.4, 0.8$ , respectively. At  $\sigma_1 = 0.2$  and  $n = 30$ , Seq\_HEI stabilized around 2.0, outperforming EnKF\_NLS (2.5 to 3.0), Aug (around 3.0), Aug\_infla (2.5 to 3.0), standard EnKF and EnKF\_infla (both around 3.5). When  $n$  increased to 60, its RMSE fell to around 1.5 while others remained above 2.5. Although all methods' RMSEs rose with the increase of the model variance  $\sigma_1$ , Seq\_HEI consistently yielded the most accurate data assimilation results.

We also summarize the mean and standard deviation of the RMSEs for the analysis state over simulation repetitions in Tables 1 and 2, which also demonstrated that Seq\_HEI consistently attained the lowest RMSE with the smallest variability, and its advantage became more pronounced as  $p$  grew. Table 3 presents the average and standard deviation of the estimates of  $\theta_F^*$  by four methods, the proposed Seq\_HEI, EnKF\_NLS, Aug and Aug\_infla, over simulation repetitions. For the proposed method, we reported the estimate  $\hat{\boldsymbol{\theta}}^{(K)}$  from the last segment  $\mathcal{T}_K$ . Consistent with previous results, Seq\_HEI provided the most accurate and stable estimates as shown in Table 3. These results confirmed that our method provided both superior data assimilation accuracy and reliable parameter estimation.

## 4.2 | Simulations on shallow water equations

The shallow water equations (SWE) are simplified PDEs for a thin fluid layer (e.g., oceans, rivers, lakes) where horizontal scales greatly exceed the vertical depth, so vertical accelerations are negligible relative to horizontal accelerations. In atmospheric and oceanic studies, SWE are also widely used to assess the performance of data assimilation methods

**TABLE 1** The average (standard derivation) of RMSEs of the data assimilation results by six different methods, Seq\_HEI (proposed), EnKF\_NLS, EnKF, EnKF\_infla, Aug, Aug\_infla, for the the Lorenz'96 model with  $p = 200, 400, q = 160, 320$  and  $T = 2000, n = 30$ .

$n = 30, p = 200, q = 160$						
Method	$\sigma_1 = 0.1$			$\sigma_2 = 0.2$		
	$\sigma_2 = 0.2$	$\sigma_2 = 0.4$	$\sigma_2 = 0.8$	$\sigma_1 = 0.2$	$\sigma_1 = 0.4$	$\sigma_1 = 0.8$
EnKF	2.84(0.170)	2.99(0.231)	3.30(0.250)	2.68(0.203)	2.83(0.278)	3.27(0.282)
EnKF_infla	2.80(0.215)	2.97(0.232)	3.27(0.218)	2.64(0.203)	2.84(0.260)	3.26(0.223)
Seq_HEI	1.71(0.160)	1.84(0.151)	1.99(0.133)	1.71(0.182)	2.05(0.186)	2.02(0.174)
EnKF_NLS	2.22(0.215)	2.43(0.142)	2.52(0.215)	2.11(0.209)	2.38(0.203)	2.72(0.342)
Aug	2.43(0.521)	2.50(0.384)	2.69(0.447)	2.45(0.474)	2.52(0.385)	3.32(0.705)
Aug_infla	2.23(0.377)	2.47(0.279)	2.65(0.358)	2.20(0.444)	2.51(0.336)	2.82(0.543)

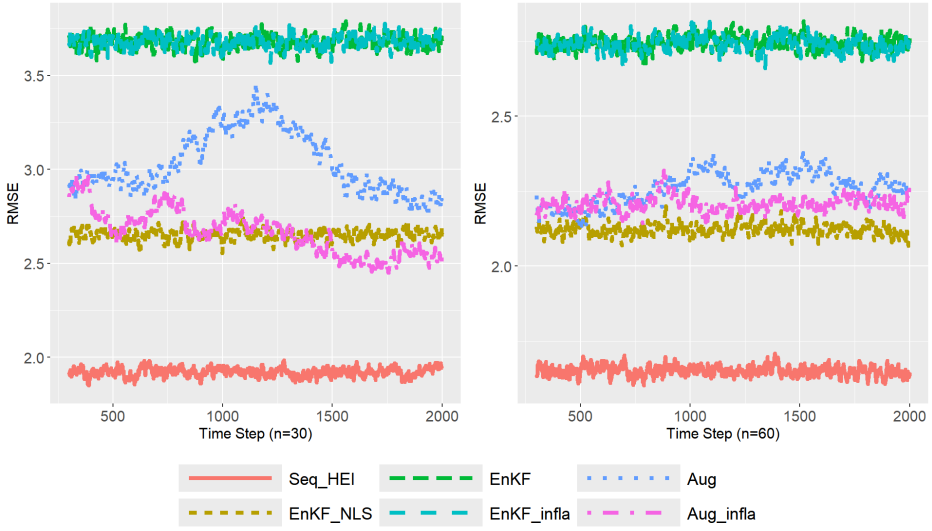
$n = 30, p = 400, q = 320$						
Method	$\sigma_1 = 0.1$			$\sigma_2 = 0.2$		
	$\sigma_2 = 0.2$	$\sigma_2 = 0.4$	$\sigma_2 = 0.8$	$\sigma_1 = 0.2$	$\sigma_1 = 0.4$	$\sigma_1 = 0.8$
EnKF	3.69(0.126)	3.87(0.139)	4.06(0.163)	3.63(0.168)	3.69(0.189)	3.74(0.207)
EnKF_infla	3.68(0.118)	3.83(0.093)	4.01(0.143)	3.62(0.189)	3.60(0.148)	3.68(0.182)
Seq_HEI	1.94(0.089)	2.05(0.111)	2.24(0.120)	2.01(0.116)	2.05(0.116)	2.21(0.116)
EnKF_NLS	2.65(0.352)	2.69(0.449)	2.97(0.350)	2.59(0.394)	2.60(0.344)	2.79(0.307)
Aug	2.84(1.224)	3.13(0.871)	2.89(0.487)	2.93(0.700)	2.96(0.660)	3.21(0.643)
Aug_infla	2.51(0.442)	2.73(0.545)	2.88(0.361)	2.62(0.659)	2.83(0.664)	3.02(0.444)

**TABLE 2** The average (standard derivation) of RMSEs of the data assimilation results by six different methods, Seq\_HEI (proposed), EnKF\_NLS, EnKF, EnKF\_infla, Aug, Aug\_infla, for the the Lorenz'96 model with  $p = 200, 400, q = 160, 320$  and  $T = 2000, n = 60$ .

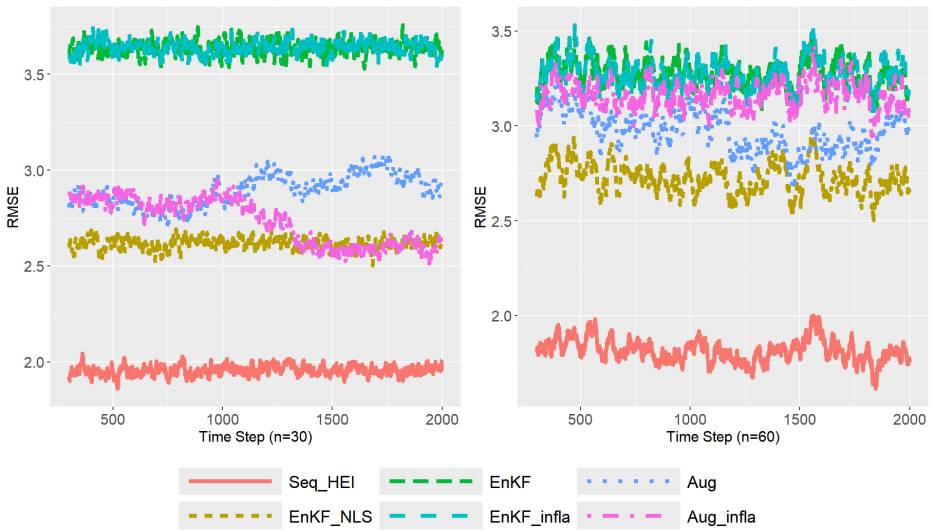
$n = 60, p = 200, q = 160$						
Method	$\sigma_1 = 0.1$			$\sigma_2 = 0.2$		
	$\sigma_2 = 0.2$	$\sigma_2 = 0.4$	$\sigma_2 = 0.8$	$\sigma_1 = 0.2$	$\sigma_1 = 0.4$	$\sigma_1 = 0.8$
EnKF	1.72(0.280)	2.18(0.252)	2.65(0.396)	1.58(0.193)	1.82(0.247)	1.98(0.240)
EnKF_infla	1.68(0.226)	2.17(0.253)	2.62(0.425)	1.57(0.218)	1.71(0.224)	1.98(0.237)
Seq_HEI	0.68(0.069)	0.87(0.076)	1.22(0.143)	0.85(0.094)	1.08(0.124)	1.36(0.146)
EnKF_NLS	1.08(0.181)	1.37(0.169)	1.97(0.391)	1.21(0.418)	1.48(0.175)	1.78(0.224)
Aug	1.16(0.259)	1.51(0.318)	1.94(0.375)	1.53(0.442)	1.60(0.165)	1.90(0.310)
Aug_infla	1.16(0.242)	1.54(0.308)	1.95(0.315)	1.32(0.132)	1.60(0.094)	1.89(0.306)

$n = 60, p = 400, q = 320$						
Method	$\sigma_1 = 0.1$			$\sigma_2 = 0.2$		
	$\sigma_2 = 0.2$	$\sigma_2 = 0.4$	$\sigma_2 = 0.8$	$\sigma_1 = 0.2$	$\sigma_1 = 0.4$	$\sigma_1 = 0.8$
EnKF	2.76(0.122)	2.92(0.186)	3.26(0.199)	3.16(0.246)	3.27(0.215)	3.50(0.241)
EnKF_infla	2.74(0.140)	2.91(0.161)	3.19(0.164)	3.13(0.249)	3.24(0.213)	3.34(0.129)
Seq_HEI	1.64(0.096)	1.75(0.097)	1.98(0.123)	1.76(0.185)	1.97(0.162)	2.32(0.172)
EnKF_NLS	2.10(0.178)	2.19(0.152)	2.50(0.141)	2.65(0.165)	2.80(0.232)	2.96(0.223)
Aug	2.27(0.416)	2.45(0.266)	2.57(0.292)	3.08(0.715)	3.12(0.546)	3.28(0.800)
Aug_infla	2.24(0.372)	2.29(0.198)	2.52(0.185)	2.96(0.476)	3.11(0.529)	3.27(0.458)



**FIGURE 4** RMSE of the data assimilation results by six different methods, Seq\_HEI (proposed), EnKF\_NLS, EnKF, EnKF\_infla, Aug, Aug\_infla, for the the Lorenz'96 model with  $(p, q, T) = (400, 320, 2000)$  and  $(\sigma_1, \sigma_2) = (0.1, 0.2)$ .



**FIGURE 5** RMSE of the data assimilation results by six different methods, Seq\_HEI (proposed), EnKF\_NLS, EnKF, EnKF\_infla, Aug, Aug\_infla, for the the Lorenz'96 model with  $(p, q, T) = (400, 320, 2000)$  and  $(\sigma_1, \sigma_2) = (0.2, 0.2)$ .

**TABLE 3** The average (standard derivation) of estimated values of the unknown parameter in the Lorenz'96 model by four different methods, Seq\_HEI (proposed), EnKF\_NLS, Aug, Aug\_infla, with  $n = 60$ ,  $p = 200, 400$ ,  $q = 160, 320$  and  $\sigma_1, \sigma_2 = 0.2, 0.4, 0.8$ .

$n = 60, p = 200, q = 160$						
Method	$\sigma_1 = 0.1$			$\sigma_2 = 0.2$		
	$\sigma_2 = 0.2$	$\sigma_2 = 0.4$	$\sigma_2 = 0.8$	$\sigma_1 = 0.2$	$\sigma_1 = 0.4$	$\sigma_1 = 0.8$
Seq_HEI	8.00(0.039)	8.00(0.048)	7.97(0.043)	8.00(0.045)	8.04(0.054)	8.05(0.303)
EnKF_NLS	7.90(0.193)	7.91(0.172)	7.88(0.243)	7.87(0.216)	7.88(0.128)	7.82(0.399)
Aug	8.11(1.744)	8.04(1.692)	7.66(1.778)	8.95(3.549)	6.99(1.673)	7.79(4.031)
Aug_infla	7.85(1.353)	8.06(1.692)	7.23(1.926)	7.62(2.028)	5.95(4.109)	7.45(3.897)

$n = 60, p = 400, q = 320$						
Method	$\sigma_1 = 0.1$			$\sigma_2 = 0.2$		
	$\sigma_2 = 0.2$	$\sigma_2 = 0.4$	$\sigma_2 = 0.8$	$\sigma_1 = 0.2$	$\sigma_1 = 0.4$	$\sigma_1 = 0.8$
Seq_HEI	7.72(0.454)	7.69(0.406)	7.53(0.216)	8.03(0.354)	7.95(0.224)	7.95(0.065)
EnKF_NLS	7.69(0.620)	7.10(1.304)	6.59(0.215)	7.76(0.442)	7.65(0.526)	7.48(0.674)
Aug	7.71(2.578)	6.99(1.376)	6.82(1.432)	7.41(4.579)	8.21(4.454)	9.18(3.332)
Aug_infla	7.37(2.221)	8.66(1.502)	8.75(1.312)	8.59(4.877)	8.35(4.641)	9.15(5.034)

393 [25, 37, 38]. Referring to [25], the continuous fields  $u(z_1, z_2, t)$ ,  $v(z_1, z_2, t)$  and  $h(z_1, z_2, t)$  satisfy

$$\begin{aligned}
 \frac{\partial u}{\partial t} + u \frac{\partial u}{\partial z_1} + v \frac{\partial u}{\partial z_2} - f v &= -g \frac{\partial h}{\partial z_1} + k \nabla^2 u, \\
 \frac{\partial v}{\partial t} + u \frac{\partial v}{\partial z_1} + v \frac{\partial v}{\partial z_2} + f u &= -g \frac{\partial h}{\partial z_2} + k \nabla^2 v, \\
 \frac{\partial h}{\partial t} + \frac{\partial(uh)}{\partial z_1} + \frac{\partial(vh)}{\partial z_2} &= k \nabla^2 h,
 \end{aligned}$$

394 where  $u, v$  and  $h$  are simplified notations for  $u(z_1, z_2, t)$ ,  $v(z_1, z_2, t)$  and  $h(z_1, z_2, t)$ , respectively,  $z_1 \in [0, L]$ ,  $z_2 \in$   
 395  $[0, D]$ ,  $g = 9.8 \text{ m s}^{-2}$ ,  $k = 5\theta_1 \times 10^4 \text{ m}^2 \text{ s}^{-1}$  is the diffusion coefficient, and  $f = \theta_2 \times 10^{-4} \text{ s}^{-1}$  is the Coriolis parameter. We  
 396 uniformly discretized  $[0, L] \times [0, D]$  to a  $p_1 \times p_2$  grid. Let  $\mathbf{u}_t = \{u(z_1, z_2, t)\}$ ,  $\mathbf{v}_t = \{v(z_1, z_2, t)\}$  and  $\mathbf{h}_t = \{h(z_1, z_2, t)\}$   
 397 be the grid-stacked vectors of state variables. We applied the proposed method to estimate the parameters  $\boldsymbol{\theta} =$   
 398  $(\theta_1, \theta_2)'$  in the Coriolis and diffusion coefficients, and recovery the true state vector  $\mathbf{X}_t = (\mathbf{u}'_t, \mathbf{v}'_t, \mathbf{h}'_t)'$ . Note that the  
 399 dimension of  $\mathbf{X}_t$  is  $p = 3\bar{p}$ , where  $\bar{p} = p_1 p_2$ .

400 We set  $L = 500 \text{ km}$  and  $D = 300 \text{ km}$ , took the true parameter  $\boldsymbol{\theta}^* = (\theta_1^*, \theta_2^*)' = (1, 1)'$ , and set the observation  
 401 dimension to  $\bar{q} = 0.2\bar{p}$ . For the proposed method, we divided the time period  $\{1, \dots, T\}$  into  $K = 12$  segments. The  
 402 observation error covariance was taken constant over time as

$$R_t \equiv R = \text{diag}(\sigma_u^2 I_{\bar{q}}, \sigma_v^2 I_{\bar{q}}, \sigma_h^2 I_{\bar{q}}),$$

403 where  $\sigma_u^2 = \sigma_v^2 = 0.5$  and  $\sigma_h^2 = 1.0$ . The initial condition  $\mathbf{X}_0 = (\mathbf{u}'_0, \mathbf{v}'_0, \mathbf{h}'_0)'$  was specified as follows. First, the initial  
 404 height (depth) field was prescribed by

$$h(z_1, z_2) = H_0 + H_1 \tanh\left(\frac{9(D/2 - z_2)}{2D}\right) + H_2 \text{sech}^2\left(\frac{9(D/2 - z_2)}{D}\right) \sin\left(\frac{2\pi z_1}{L}\right),$$

405 with  $H_0 = 50 \text{ m}$ ,  $H_1 = 5.5 \text{ m}$ , and  $H_2 = 3.325 \text{ m}$ . Then, the initial velocity fields  $(u_0, v_0)$  were then obtained from  $h_0$

using the geostrophic balance. We integrated the SWE for 24.5 hours using 30-second time steps, for a total of 2940



**FIGURE 6** RMSE of the data assimilation results by six different methods, Seq\_HEI (proposed), EnKF\_NLS, EnKF, EnKF\_infla, Aug, and Aug\_infla, for the shallow water model with  $(T, n) = (3000, 25)$ .

406 steps, after a 60-step spin-up. Observations were assimilated every 60 steps (i.e., every 30 minutes). The initial value  
 407 of  $\theta$  was set as  $\theta^{(0)} = (0.8, 0.8)'$ . For the standard EnKF and EnKF\_infla that don't involve parameter estimation, we  
 408 used this initial value to run the SWE.  
 409

410 Figure 6 reports the data assimilation RMSEs between the true state and analysis state by the proposed method  
 411 (Seq\_HEI) and five existing methods considered in Section 4.1 under two settings of model dimensions  $(p, q)$ . In both  
 412 settings, the proposed method Seq\_HEI achieved the lowest RMSE around 0.4, demonstrating its estimation accuracy  
 413 and robustness to the increase of dimensions. The EnKF\_NLS method performed slightly worse, but still outperformed  
 414 the other methods. The augmented data assimilation approaches, Aug and Aug\_infla, showed moderate performance,  
 415 achieving lower RMSE than the standard EnKF and EnKF\_infla, which suffered from large initial spikes due to the  
 416 misspecification of  $\theta$ .

417 Overall, the results from both simulation studies confirm that the proposed method can achieve better data as-  
 418 similation accuracy by a sequentially updated estimation of the model parameters. The improvement is particularly  
 419 evident when compared to traditional EnKF methods and their inflation-based variants that don't consider model  
 420 calibration or parameter estimation.

## 421 5 | CONCLUSION AND DISCUSSION

422 The unknown parameter in state-space models can introduce significant bias that even advanced data assimilation  
 423 techniques cannot fully correct. Existing methods for parameter estimation often break down under strong non-

424 linearity, high dimensionality, or large error variances, and they incur heavy computational costs. To address these  
425 challenges, this paper proposes a sequential machine learning algorithm that couples the ensemble Kalman filter with  
426 hierarchical Bayesian optimization. This method not only provides effective sequential estimates of the parameters  
427 but also achieves an efficient “estimating-while-assimilating” computational paradigm. The asymptotic properties in-  
428 cluding the consistency and the asymptotic normality of the proposed estimator have been established in this paper.  
429 Numerical experiments on the Lorenz-96 and the shallow water models demonstrate the accuracy and efficiency of  
430 the proposed method in high-dimensional settings with large variance in both observation and model errors.

431 Since this paper focuses primarily on parameter estimation, it overlooks certain data assimilation techniques, for  
432 example, the localization methods in high-dimensional data assimilation. In atmospheric or oceanic data assimilation,  
433 the dimension of state variables can be much higher than the time step and the ensemble size. In such cases, local-  
434 ization can be incorporated into our sequential data assimilation procedure to improve the accuracy of forecast error  
435 covariance estimation. The loss function in (7) and parameter estimation by the proposed HEI algorithm can be simi-  
436 larly constructed with localized data assimilation methods. Our study also assumes that the forecast error is unbiased  
437 conditional on the historical observations (Assumption S2). If the data assimilation and prediction results are biased,  
438 a debiasing step [4] is necessary before constructing the loss function for parameter estimation. We will leave those  
439 two extensions of the proposed method in future works.

## 440 Conflict of interest

441 The authors declare that they have no conflict of interest.

## 442 references

- 443 [1] Lorenc AC. Analysis methods for numerical weather prediction. *Quarterly Journal of the Royal Meteorological Society*  
444 1986;112(474):1177–1194.
- 445 [2] Evensen G. Sequential data assimilation with a nonlinear quasi-geostrophic model using Monte Carlo methods to fore-  
446 cast error statistics. *Journal of Geophysical Research: Oceans* 1994;99(C5):10143–10162.
- 447 [3] Carrassi A, Bocquet M, Bertino L, Evensen G. Data assimilation in the geosciences: An overview of methods, issues, and  
448 perspectives. *Wiley Interdisciplinary Reviews: Climate Change* 2018;9:e535.
- 449 [4] Dee DP, Da Silva AM. Data assimilation in the presence of forecast bias. *Quarterly Journal of the Royal Meteorological*  
450 *Society* 1998;124(545):269–295.
- 451 [5] Zheng X. An adaptive estimation of forecast error covariance parameters for Kalman filtering data assimilation. *Advances*  
452 *in Atmospheric Sciences* 2009;26:154–160.
- 453 [6] Oke PR, Brassington GB, Griffin DA, Schiller A. The Bluelink ocean data assimilation system (BODAS). *Ocean Modelling*  
454 2008;21(1):46–70.
- 455 [7] Sakov P, Counillon F, Bertino L, Lisæter KA, Oke PR, Korabelv A. TOPAZ4: an ocean-sea ice data assimilation system for  
456 the North Atlantic and Arctic. *Ocean Science* 2012;8(4):633–656.
- 457 [8] Friedland B. Treatment of bias in recursive filtering. *IEEE Transactions on Automatic Control* 1969;14(4):359–367.
- 458 [9] Kalman RE. A New Approach to Linear Filtering and Prediction Problems. *Journal of Basic Engineering* 1960;82(1):35–  
459 45.

- 460 [10] Evensen G. The ensemble Kalman filter: Theoretical formulation and practical implementation. *Ocean dynamics*  
461 2003;53:343–367.
- 462 [11] Anderson JL, Anderson SL. A Monte Carlo implementation of the nonlinear filtering problem to produce ensemble  
463 assimilations and forecasts. *Monthly Weather Review* 1999;127(12):2741–2758.
- 464 [12] Liang X, Zheng XG, Zhang SP, Wu GC, Dai YJ, Li Y. Maximum likelihood estimation of inflation factors on error covariance  
465 matrices for ensemble Kalman filter assimilation. *Quarterly Journal of the Royal Meteorological Society* 2012;138:263–  
466 273.
- 467 [13] Raanes PN, Bocquet M, Carrassi A. Adaptive covariance inflation in the ensemble Kalman filter by Gaussian scale  
468 mixtures. *Quarterly Journal of the Royal Meteorological Society* 2018;145:53–75.
- 469 [14] Greybush SJ, Kalnay E, Miyoshi T, Ide K, Hunt BR. Balance and ensemble Kalman filter localization techniques. *Monthly*  
470 *Weather Review* 2011;139(2):511–522.
- 471 [15] Perianez A, Reich H, Potthast R. Optimal localization for ensemble Kalman filter systems. *Journal of the Meteorological*  
472 *Society of Japan Ser II* 2014;92(6):585–597.
- 473 [16] Sun HX, Wang S, Zheng X, Chen SX. High-dimensional Ensemble Kalman Filter with Localization, Inflation and Iterative  
474 Updates. *Quarterly Journal of the Royal Meteorological Society* 2024;150:4870–4884.
- 475 [17] Mellor GL, Yamada T. Development of a turbulence closure model for geophysical fluid problems. *Review of Geophysics*  
476 1982;20(4):851–875.
- 477 [18] Dee DP. Bias and data assimilation. *Quarterly Journal of the Royal Meteorological Society* 2005;131(613):3323–3343.
- 478 [19] Sawada Y. An efficient estimation of time-varying parameters of dynamic models by combining offline batch optimization  
479 and online data assimilation. *Journal of Advances in Modeling Earth Systems* 2022;14:e2021MS002882.
- 480 [20] Yang Y, Mémin E. Estimation of physical parameters under location uncertainty using an ensemble<sup>2</sup> expecta-  
481 tion–maximization algorithm. *Quarterly Journal of the Royal Meteorological Society* 2019;145:418–433.
- 482 [21] Lucini MM, Van Leeuwen PJ, Pulido M. Model error estimation using the expectation maximization algorithm and a  
483 particle flow filter. *SIAM/ASA Journal on Uncertainty Quantification* 2021;9(2):681–707.
- 484 [22] Ahn KW, Chan KS. Approximate conditional least squares estimation of a nonlinear state-space model via an unscented  
485 Kalman filter. *Computational Statistics and Data Analysis* 2014;69.
- 486 [23] Lunderman S, Morzfeld M, Posselt DJ. Using global Bayesian optimization in ensemble data assimilation: parameter  
487 estimation, tuning localization and inflation, or all of the above. *Tellus A: Dynamic Meteorology and Oceanography*  
488 2021;73(1):1–16.
- 489 [24] Chen ZH, Mak S, Wu CFJ. A Hierarchical Expected Improvement Method for Bayesian Optimization. *Journal of the*  
490 *American Statistical Association* 2024;119(546):1619–1632.
- 491 [25] Wang S, Chen SX, Sun HX. High Dimensional Ensemble Kalman Filter. arXiv preprint arXiv:250500283v3 2025;.
- 492 [26] Zhang GJ, McFarlane NA. Sensitivity of climate simulations to the parameterization of cumulus convection in the Cana-  
493 dian Climate Centre general circulation model. *Atmosphere-ocean* 1995;33(3):407–446.
- 494 [27] Sawada Y, Duc L. An efficient and robust estimation of spatio-temporally distributed parameters in dynamic models by  
495 an ensemble Kalman filter. *Journal of Advance in Modeling Earth Systems* 2024;16:e2023MS003821.
- 496 [28] Hager WW. Updating the inverse of a matrix, vol. 31. SIAM; 1989.

- 497 [29] Castruccio S, McInerney DJ, Stein ML, Liu Crouch F, Jacob RL, Moyer EJ. Statistical Emulation of Climate Model Projec-  
498 tions Based on Precomputed GCM Runs. *Journal of Climate* 2014;27(5):1829–1844.
- 499 [30] Girard S, Mallet V, Korsakissok I, Mathieu A. Emulation and Sobol' Sensitivity Analysis of an Atmospheric Dispersion  
500 Model Applied to the Fukushima Nuclear Accident. *Journal of Geophysical Research: Atmospheres* 2016;121(7):3484–  
501 3496.
- 502 [31] Song Y, Khalid Z, Genton MG. Efficient stochastic generators with spherical harmonic transformation for high-resolution  
503 global climate simulations from CESM2-LENS2. *Journal of the American Statistical Association* 2024;119(548):2493–  
504 2507.
- 505 [32] Rougier J, Goldstein M, House L. Second-order exchangeability analysis for multimodel ensembles. *Journal of the*  
506 *American Statistical Association* 2013;108(503):852–863.
- 507 [33] Sansom PG, Ferro CA, Stephenson DB, Goddard L, Mason SJ. Best practices for postprocessing ensemble climate  
508 forecasts. Part I: Selecting appropriate recalibration methods. *Journal of Climate* 2016;29(20):7247–7264.
- 509 [34] Cressie NAC. *Statistics for spatial data*. John Wiley & Sons; 1993.
- 510 [35] Santner T, Williams B, Notz W. *The Design and Analysis of Computer Experiments*. 2nd ed. New York: Springer; 2018.
- 511 [36] Papamichail C, Bouzebda S, Limnios N. Predictability: a problem partly solved. *Proceedings on Seminar on Predictability*  
512 1996;1:1–18.
- 513 [37] Greybush SJ, Kalnay E, Miyoshi T, Ide K, Hunt BR. Balance and ensemble Kalman filter localization techniques. *Monthly*  
514 *Weather Review* 2011;139:511–522.
- 515 [38] Lei L, Stauffer DR, Deng A. A hybrid nudging-ensemble Kalman filter approach to data assimilation. Part II: application  
516 in a shallow-water model. *Tellus A: Dynamic Meteorology and Oceanography* 2012;64:18485.