

# UC San Diego

## UC San Diego Previously Published Works

### Title

High-Dimensional Covariance Estimation From a Small Number of Samples

### Permalink

<https://escholarship.org/uc/item/7qp2z12v>

### Journal

Journal of Advances in Modeling Earth Systems, 16(9)

### ISSN

1942-2466

### Authors

Vishny, David

Morzfeld, Matthias

Gwirtz, Kyle

et al.

### Publication Date

2024-09-01

### DOI

10.1029/2024ms004417

### Copyright Information

This work is made available under the terms of a Creative Commons Attribution-NonCommercial-NoDerivatives License, available at

<https://creativecommons.org/licenses/by-nc-nd/4.0/>

Peer reviewed

# High-Dimensional Covariance Estimation From a Small Number of Samples

David Vishny<sup>1</sup>, Matthias Morzfeld<sup>1</sup>, Kyle Gwartz<sup>2</sup>, Eviatar Bach<sup>3,4</sup>, Oliver  
R.A. Dunbar<sup>3</sup>, Daniel Hodyss<sup>5</sup>

<sup>1</sup>Scripps Institution of Oceanography, University of California, San Diego, CA, USA

<sup>2</sup>NASA Goddard Space Flight Center, Greenbelt, MD, USA

<sup>3</sup>California Institute of Technology, Pasadena, CA, USA

<sup>4</sup>University of Reading, Reading, UK

<sup>5</sup>Remote Sensing Division, Naval Research Laboratory, Washington DC, USA

## Key Points:

- We introduce several methods of covariance matrix estimation that adaptively select regularization parameters based on estimates of sampling error.
- One method, Noise-Informed Covariance Estimation (NICE), stands out because it guarantees a positive semi-definite estimator at a low computational cost.
- All new covariance estimation methods perform well on a large variety of test problems.

**Abstract**

We synthesize knowledge from numerical weather prediction, inverse theory, and statistics to address the problem of estimating a high-dimensional covariance matrix from a small number of samples. This problem is fundamental in statistics, machine learning/artificial intelligence, and in modern Earth science. We create several new adaptive methods for high-dimensional covariance estimation, but one method, which we call NICE (**n**oise-**i**nformed **c**ovariance **e**stimation), stands out because it has three important properties: (i) NICE is conceptually simple and computationally efficient; (ii) NICE guarantees symmetric positive semi-definite covariance estimates; and (iii) NICE is largely tuning-free. We illustrate the use of NICE on a large set of Earth science-inspired numerical examples, including cycling data assimilation, inversion of geophysical field data, and training of feed-forward neural networks with time-averaged data from a chaotic dynamical system. Our theory, heuristics and numerical tests suggest that NICE may indeed be a viable option for high-dimensional covariance estimation in many Earth science problems.

**Plain Language Summary**

Models of physical processes must be fitted to real-world data before they are useful for prediction. In some cases, the most practical way to fit models to data is to run a set—or *ensemble*—of simulations with different physics or initial conditions. One then uses the covariances among the inputs and outputs to modify the simulations so that they fit the data better. To reduce noise in the covariances, one ideally uses an ensemble size that is larger than the number of unknown variables, but this becomes impractical when the number of unknowns is large. To improve the performance of this fitting process when the ensemble size is small, one can discount covariances between variables that are likely due to noise. We introduce several methods of covariance estimation that determine the degree to which covariances are discounted based on expected levels of noise. All new methods perform well on a series of Earth science-inspired problems, but we highlight one method that preserves a key property of covariance matrices at a low computational cost.

**1 Introduction**

We consider the problem of estimating the covariance matrix  $\mathbf{P}$  of an  $n_x$ -dimensional random variable  $\mathbf{x}$ , based on a set of  $n_e \ll n_x$  independent samples  $\mathbf{x}_i$ ,  $i = 1, \dots, n_e$ . Estimating a covariance matrix from scarce samples is a fundamental challenge in science, engineering, statistics, and in the sub-fields of machine learning and artificial intelligence (Wainwright, 2019). Our interest in covariance estimation is motivated by the problem of fitting models of Earth processes to data. As an example, consider numerical weather prediction (NWP), where the  $\mathbf{x}_i$  represent an ensemble of global weather forecasts. The dimension  $n_x$  corresponds to the number of unknowns in a global atmospheric model, and it is on the order of  $10^8$ . The number of forecasts (the ensemble size  $n_e$ ) is small because each forecast requires a simulation of Earth’s atmosphere, which is expensive. A commonly used ensemble size in NWP is on the order of  $10^2$ —six orders of magnitude smaller than the number of unknowns. A common approach to update the forecast with atmospheric data is the ensemble Kalman filter (EnKF, Evensen (1994, 2009)). The EnKF updates rely on the covariance matrix associated with the ensemble, but the *empirical covariance matrix*

$$\hat{\mathbf{P}} = \frac{1}{n_e - 1} \sum_{i=1}^{n_e} (\mathbf{x}_i - \hat{\boldsymbol{\mu}})(\mathbf{x}_i - \hat{\boldsymbol{\mu}})^T, \quad \hat{\boldsymbol{\mu}} = \frac{1}{n_e} \sum_{i=1}^{n_e} \mathbf{x}_i, \quad (1)$$

63 is generally inaccurate if  $n_e \ll n_x$  (Bickel & Levina, 2008; Wainwright, 2019). Various  
 64 strategies for improving the accuracy of the empirical estimate have been developed over  
 65 the years, and we review the relevant literature below.

66 The prevailing method of covariance estimation in NWP is called *localization* (Houtekamer  
 67 & Mitchell, 1998, 2001; Ott et al., 2004). Localization enforces on the empirical covari-  
 68 ance matrix the assumption that covariances decay with spatial distance (although the  
 69 terminology has also been used in other contexts and to refer to covariance corrections  
 70 that are not spatial, see, e.g., Morzfeld et al. (2019)). To execute the localization, one  
 71 defines an  $n_x \times n_x$ , symmetric positive semi-definite (PSD) *localization matrix*  $\mathbf{L}$ , which  
 72 encodes the spatial decay pattern of correlations (Gaspari & Cohn, 1999; Gilpin et al.,  
 73 2023). One then obtains the localized covariance estimator

$$74 \quad \hat{\mathbf{P}}_{\text{loc}} = \mathbf{L} \circ \hat{\mathbf{P}}, \quad (2)$$

75 where the open circle denotes the Hadamard (element-wise) product. Localization has  
 76 proven successful for estimating high-dimensional covariance matrices from a small set  
 77 of samples in NWP and, for that reason, localization is a standard component in oper-  
 78 ational weather forecasting systems (Hamill et al., 2009; Bannister, 2017).

79 We present a new covariance estimation method that is more broadly applicable  
 80 than classical localization because it does not require *a priori* assumptions about the cor-  
 81 relation structure (e.g., the spatial decay in covariance localization). We call our method  
 82 **Noise-Informed Covariance Estimation (NICE)**. NICE replaces assumptions about the  
 83 correlation structure with the assumption that *small to medium correlations are likely*  
 84 *caused by sampling error* and, therefore, should be damped or deleted. This assumption  
 85 is not universally true (it is easy to come up with counter examples), but it is rooted in  
 86 rigorous sampling error theory (Ménétrier et al., 2015; Morzfeld & Hodyss, 2023; Flow-  
 87 erdew, 2015; Lee, 2021b; Anderson, 2012). NICE achieves three main objectives:

- 88 1. **Adaptivity.** NICE ensures that differences between sampled and corrected cor-  
 89 relations are within an expected noise level. The noise level is determined by the  
 90 sample size and the distribution of empirical correlations so that the entire covari-  
 91 ance estimation process is adaptive and largely tuning-free.
- 92 2. **Positive semi-definiteness.** NICE guarantees a symmetric positive semi-definite  
 93 (PSD) covariance estimator. Symmetry and positive semi-definiteness are defin-  
 94 ing properties of covariance matrices, but some competing methods do not guar-  
 95 antee PSD estimates.
- 96 3. **Computational efficiency.** NICE is computationally efficient and easy to im-  
 97 plement because it avoids solving optimization problems over PSD matrices.

98 We put NICE to the test in a variety of problems with different and unknown cor-  
 99 relation structures: (i) estimation of covariance matrices from Gaussian samples; (ii) cy-  
 100 cling data assimilation problems with ensemble Kalman filters (Evensen, 2009); (iii) in-  
 101 version of geophysical data with regularized ensemble Kalman inversion (EKI, Chada  
 102 et al. (2020)); and (iv) training of a feed-forward neural network with EKI (Iglesias et  
 103 al., 2013; Kovachki & Stuart, 2019; Cleary et al., 2021). Various error metrics are used  
 104 to evaluate performance in these problems. Across all examples and all error metrics,  
 105 we find that NICE works out-of-the-box with minimal tuning.

106 Estimated noise levels can also be used to make other covariance estimation meth-  
 107 ods adaptive and largely tuning-free. We introduce *new* adaptive versions of power law  
 108 corrections (Ad.-PLC, see Lee (2021b) and Section 3.4.1), adaptive (spatial) localization  
 109 (Ad.-Loc., Section 3.4.2), adaptive soft-thresholding (Ad.-ST, see Wainwright (2019) and  
 110 Section 3.4.3) and adaptive sparse covariance estimation (ASCE, see Xue et al. (2012)  
 111 and Section 3.4.4). All new methods fall under the umbrella of noise-informed covari-  
 112 ance estimation because all of them leverage an understanding of noise in empirical cor-

113 relations. However, some do not guarantee a PSD estimator and others are more com-  
 114 putationally involved. The specific method we refer to as NICE is the *only* method that  
 115 satisfies all three of our objectives: adaptivity, PSD guarantees and computational ef-  
 116 ficiency.

117 It is important to be specific about the terms “high-dimensional” and about com-  
 118 putational efficiency. In this paper, we focus on covariance estimation methods that con-  
 119 struct the entire covariance matrix. As such, the methods are limited in their use to ma-  
 120 trices of dimension  $10^4 \times 10^4$  or smaller. Higher-dimensional problems, e.g., of the ex-  
 121 treme size of NWP ( $10^8$  or more unknowns), require that we perform computations with-  
 122 out constructing the whole covariance matrix. The methods we describe here could po-  
 123 tentially be adapted to such problems, but these adaptations are beyond the scope of  
 124 this paper. The computational efficiency of covariance estimation depends on the algo-  
 125 rithms used. We focus on algorithms that perform simple element-wise operations on the  
 126 empirical covariance matrix. Many methods in the statistical literature, however, per-  
 127 form covariance estimation by solving optimization problems over PSD matrices, which  
 128 is computationally expensive.

129 The rest of this paper is organized as follows. Section 2 reviews background ma-  
 130 terials. We first explain why covariance estimation from a small number of samples is  
 131 important in Earth science, specifically in EnKF and in EKI. We further emphasize the  
 132 importance of PSD covariance estimates in the context of EnKF or EKI. We then re-  
 133 view covariance localization in NWP and several covariance estimation methods from  
 134 the statistical literature. Finally, we briefly describe Morozov’s discrepancy principle, a  
 135 classical concept in inverse theory. The discrepancy principle is the tool we use to make  
 136 covariance estimation methods adaptive. Section 3 describes our new methodology (NICE),  
 137 and other new adaptive covariance estimation methods. We apply NICE and a large num-  
 138 ber of competing methods (new and old) in a wide variety of problems in Section 4, be-  
 139 fore ending the paper with a summary and conclusions in Section 5.

## 140 2 Background

### 141 2.1 Ensemble Kalman Filters and their Localization

142 The goal of ensemble Kalman filtering (EnKF) is to use data to update a forecast  
 143 generated by a computational model. An important example is numerical weather pre-  
 144 diction (NWP), where the forecast describes atmospheric states in the form of  $n_e$  vec-  
 145 tors  $\mathbf{x}_i$ ,  $i = 1, \dots, n_e$ , each of dimension  $n_x$ . The vectors  $\mathbf{x}_i$  are referred to as “ensem-  
 146 ble members.” Typically, the ensemble size  $n_e$  is smaller than the dimension of the en-  
 147 semble members ( $n_e \ll n_x$ ). The reason is that each ensemble member is the result of  
 148 a simulation with a computationally expensive atmospheric model, so that  $n_e$  must be  
 149 small, or else the computations are infeasible. In NWP,  $n_e$  is usually a few hundred, and  
 150  $n_x$  is in the billions.

151 The forecast is updated by an observation (data), which is an  $n_y$ -dimensional vec-  
 152 tor  $\mathbf{y}$ , where, often but not always,  $n_e \ll n_y \ll n_x$ . For ease of presentation, we as-  
 153 sume that the observation is a linear function of the forecasted variables so that

$$154 \quad \mathbf{y} = \mathbf{H}\mathbf{x} + \boldsymbol{\varepsilon}, \quad (3)$$

155 where  $\mathbf{H}$  is an  $n_y \times n_x$  matrix and  $\boldsymbol{\varepsilon}$  is a Gaussian random variable with mean zero and  
 156 covariance matrix  $\mathbf{R}$ , which we write as  $\boldsymbol{\varepsilon} \sim \mathcal{N}(\mathbf{0}, \mathbf{R})$ . The assumption of a linear ob-  
 157 servation is commonly violated; nevertheless, we demonstrate in our numerical exper-  
 158 iments that the intuition and conclusions from the linear analysis extend to the nonlin-  
 159 ear case.

160 Ensemble Kalman filtering (EnKF) is a catch-all term for a whole suite of meth-  
 161 ods that merge the observation and the forecast within a Bayesian framework. The up-

date step of a stochastic EnKF (Evensen, 1994; Burgers et al., 1998; Evensen, 2009) is

$$\mathbf{x}_i^a = \mathbf{x}_i + \hat{\mathbf{K}}(\mathbf{y} - (\mathbf{H}\mathbf{x}_i + \boldsymbol{\varepsilon}_i)), \quad (4)$$

where  $\boldsymbol{\varepsilon}_i$  is a sample drawn from  $\mathcal{N}(\mathbf{0}, \mathbf{R})$ . The Kalman gain  $\hat{\mathbf{K}}$  is computed from the ensemble as

$$\hat{\mathbf{K}} = \hat{\mathbf{P}}\mathbf{H}^T(\mathbf{H}\hat{\mathbf{P}}\mathbf{H}^T + \mathbf{R})^{-1}, \quad (5)$$

where  $\hat{\mathbf{P}}$  is the empirical covariance in (1). The Kalman gain defines how to update each ensemble member in view of the observation. Since the Kalman gain depends critically on the forecast covariance  $\hat{\mathbf{P}}$ , the EnKF update is only useful if the covariance estimate is accurate, which usually requires that  $n_e$  is larger than  $n_x$  (although the situation can be more complex, e.g., with  $n_e$  directly depending on the number of observations and their independence (Chorin & Morzfeld, 2013; Agapiou et al., 2017; Al Ghattas & Sanz-Alonso, 2022; Hodyss & Morzfeld, 2023)).

Localization is a technique that enables the use of EnKF when  $n_e \ll n_x$ . A common version of localization in the EnKF is to use Hadamard products as in (2) and to define the localization matrix by the Gaspari–Cohn covariance function (Gaspari & Cohn, 1999) or its anisotropic extensions (Gilpin et al., 2023). The localization matrix implements a spatial decay of correlation and the rate of decay can be controlled via a length scale. Different methods for adaptively selecting this length scale, or localizing in a flow-dependent manner to account for temporal variations in the correlation structure, have been proposed (Zhen & Zhang, 2014; Chevrotière & Harlim, 2017; Anderson, 2012; Bishop & Hodyss, 2009a, 2009b, 2007, 2011; Luk et al., 2024).

Other implementations of the EnKF include the ensemble adjustment Kalman filter (EAKF) (Anderson, 2001) and ensemble transform filters (ETKF) (Ott et al., 2004; Tippett et al., 2003). Localization in an EAKF is achieved by working directly with the Kalman gain, reducing the effects of an observation on elements of the Kalman gain that are far from the observation (Morzfeld & Hodyss, 2023; Hodyss & Morzfeld, 2023). Localization in an ETKF is implemented by performing a “local” analysis, so that each grid point is updated by a set of nearby observations (domain localization). Variational/hybrid data assimilation (DA) algorithms combine a classical minimization (variational) approach (Talagrand & Courtier, 1987) with an ensemble to approximate uncertainties (Hamill & Snyder, 2000; Lorenc, 2003; Zhang et al., 2009; Buehner et al., 2013; Kuhl et al., 2013; Poterjoy & Zhang, 2015). Hybrid DA also requires localization, which is usually applied using Hadamard products, but without explicitly forming the covariance matrix (Buehner, 2005). Multi-scale extensions of localization are available for hybrid DA and/or EnKFs (Buehner, 2012; Miyoshi & Kondo, 2013; Buehner & Shlyayeva, 2015; Lorenc, 2017; Harty et al., 2021).

Finally, we note that all conventional localization methods require tuning. The tuning process usually amounts to picking a length scale that defines the localization and then running a cycling EnKF over a set of training observations. This process is repeated with various length scales until one encounters a length scale that leads to an acceptable error metric.

## 2.2 Ensemble Kalman Inversion

The goal in ensemble Kalman inversion (EKI, Iglesias et al. (2013)) is to minimize the cost function

$$J(\mathbf{x}) = \left\| \mathbf{R}^{-1/2}(\mathbf{y} - \mathcal{G}(\mathbf{x})) \right\|_2^2, \quad (6)$$

where vertical bars denote the two-norm (i.e.,  $\|\mathbf{b}\|_2 = \sqrt{\mathbf{b}^T\mathbf{b}}$ ),  $\mathbf{y}$  are data,  $\mathbf{x}$  are unknown model parameters, and  $\mathcal{G}(\cdot)$  is a nonlinear model that maps the model parameters to the data; the symmetric positive definite matrix  $\mathbf{R}$  defines expected errors in the

210 data, represented by a mean-zero Gaussian random variable with covariance matrix  $\mathbf{R}$ ;  
 211  $\mathbf{R}^{-1/2}$  is the inverse of a matrix square root of  $\mathbf{R} = \mathbf{R}^{1/2} (\mathbf{R}^{1/2})^T$ .

212 EKI performs the optimization by iteratively updating an ensemble as follows. The  
 213 ensemble at iteration  $k$  are the  $n_e$  vectors  $\mathbf{x}_i^k$  and we define  $n_e$  corresponding vectors  $\mathbf{g}_i^k =$   
 214  $\mathcal{G}(\mathbf{x}_i^k)$ . Each ensemble member is updated according to

$$215 \quad \mathbf{x}_i^{k+1} = \mathbf{x}_i^k + \hat{\mathbf{C}}_{xg}^k (\hat{\mathbf{C}}_{gg}^k + \mathbf{R})^{-1} (\mathbf{y} - (\mathbf{g}_i^k + \boldsymbol{\eta}_i)), \quad (7)$$

216 where  $\hat{\mathbf{C}}_{gg}^k$  is the covariance matrix associated with the vectors  $\mathbf{g}_i^k$ ,  $\hat{\mathbf{C}}_{xg}^k$  is the covariance  
 217 between the vectors  $\mathbf{x}_i^k$  and  $\mathbf{g}_i^k$  and where  $\boldsymbol{\eta}_i$  is a draw from the Gaussian with mean zero  
 218 and covariance matrix  $\mathbf{R}$ . More specifically, if we define the matrices (ensemble pertur-  
 219 bations)

$$220 \quad \mathbf{X}^k = \frac{1}{\sqrt{n_e - 1}} (\mathbf{x}_1^k - \bar{\mathbf{x}}^k \quad \mathbf{x}_2^k - \bar{\mathbf{x}}^k \quad \dots \quad \mathbf{x}_{n_e}^k - \bar{\mathbf{x}}^k), \quad \bar{\mathbf{x}}^k = \frac{1}{n_e} \sum_{j=1}^{n_e} \mathbf{x}_j^k, \quad (8)$$

$$221 \quad \mathbf{G}^k = \frac{1}{\sqrt{n_e - 1}} (\mathbf{g}_1^k - \bar{\mathbf{g}}^k \quad \mathbf{g}_2^k - \bar{\mathbf{g}}^k \quad \dots \quad \mathbf{g}_{n_e}^k - \bar{\mathbf{g}}^k), \quad \bar{\mathbf{g}}^k = \frac{1}{n_e} \sum_{j=1}^{n_e} \mathbf{g}_j^k, \quad (9)$$

222 then the covariances are

$$223 \quad \hat{\mathbf{C}}_{xg}^k = \mathbf{X}^k \otimes \mathbf{G}^k, \quad (10)$$

$$224 \quad \hat{\mathbf{C}}_{gg}^k = \mathbf{G}^k \otimes \mathbf{G}^k, \quad (11)$$

225 where the symbol  $\otimes$  denotes the outer product  $\mathbf{A} \otimes \mathbf{B} = \mathbf{AB}^T$ , where  $\mathbf{A}$  and  $\mathbf{B}$  are  
 226 vectors or matrices of compatible sizes. Note that the EKI update equation (7) is equiv-  
 227 alent to an EnKF update in (4) because  $\hat{\mathbf{C}}_{xg}^k = \hat{\mathbf{P}}\mathbf{H}^T$  and  $\hat{\mathbf{C}}_{gg}^k = \mathbf{H}\hat{\mathbf{P}}\mathbf{H}^T$  when  $\mathcal{G}(\mathbf{x}) =$   
 228  $\mathbf{H}\mathbf{x}$  is linear. The theory around EKI tells us that the iteration (7) converges, in the sense  
 229 that the ensemble collapses onto the minimizer of the cost function, under typical as-  
 230 sumptions (Schillings & Stuart, 2018, 2017; Chada & Tong, 2022). As with EnKF, there  
 231 are several variants of EKI (Huang et al., 2022; Lee, 2021a).

232 Convergence of the EKI iteration requires that the covariance estimates  $\hat{\mathbf{C}}_{xg}^k$  and  
 233  $\hat{\mathbf{C}}_{gg}^k$  are sufficiently accurate, which usually means that the ensemble size is large. Lo-  
 234 calization can be used within an EKI to keep the ensemble size small (Tong & Morzfeld,  
 235 2023; Al Ghattas & Sanz-Alonso, 2022; Lee, 2021b).

236 EKI has found application in climate sciences (Cleary et al., 2021; Bieli et al., 2022;  
 237 Schneider et al., 2021; Dunbar et al., 2022), and Julia code for it is available (Dunbar  
 238 et al., 2022). In a climate science context, model parameters appear in sub-gridscale clo-  
 239 sures of climate models (e.g., physical constants or weights of neural networks). A promis-  
 240 ing approach to optimizing sub-gridscale closures is to formulate the cost function based  
 241 on the misfit between modeled and observed climate statistics (Schneider et al., 2024).  
 242 In this scenario, derivatives of the cost function with respect to the model parameters  
 243 are difficult or impossible to compute, making derivative-free optimization via EKI at-  
 244 tractive.

245 Ensemble algorithms that are related to EKI, and which in fact pre-date EKI, are  
 246 known as iterative ensemble Kalman filters/smoothers (Chen & Oliver, 2013, 2010, 2017;  
 247 Emerick & Reynolds, 2011; Luo et al., 2018; Bocquet & Sakov, 2014; Bocquet, 2016; Hodyss,  
 248 Bishop, & Morzfeld, 2016) or multiple data assimilation (Emerick & Reynolds, 2013).  
 249 These methods are popular in reservoir modeling, but also find applications in atmospheric  
 250 sciences. A recent, mathematical overview of how some of the methods are related is given  
 251 by Chada et al. (2021) and an NWP-focused overview is provided by Hodyss, Bishop,  
 252 and Morzfeld (2016).

### 2.3 Positive Semi-Definite Covariance Estimators

A fundamental property of covariance matrices is that they are symmetric positive semi-definite (PSD, Horn and Johnson (1991)). The practical relevance of PSD estimates of  $\hat{\mathbf{P}}$  is apparent in the EnKF, where the Kalman gain (5) requires that the matrix

$$\mathbf{H}\hat{\mathbf{P}}\mathbf{H}^T + \mathbf{R}, \quad (12)$$

is well-conditioned. Since the observation error covariance matrix  $\mathbf{R}$  is usually positive definite, a PSD estimate  $\hat{\mathbf{P}}$  guarantees that (12) is positive definite, invertible and well-conditioned. One can run into numerical trouble if  $\hat{\mathbf{P}}$  is not PSD, because the matrix in (12) may be singular or ill-conditioned. Localization via Hadamard products, as used in NWP, guarantees a PSD covariance estimate by the Schur product theorem when the localization matrix  $\mathbf{L}$  is PSD (Schur, 1911).

In general, however, the PSD constraint is not easy to satisfy during covariance estimation, and many covariance estimation methods do not guarantee a PSD estimate (Khare et al., 2019; Xue et al., 2012). When we review covariance estimation methods, we comment on their PSD guarantees.

### 2.4 Beyond Localization

It has long been recognized that the assumption of a spatial decay of correlation, which is at the core of localization, is not universally applicable. Adaptive localization methods (Anderson, 2012; Lee, 2021b; Bishop & Hodyss, 2009a, 2009b, 2007) are well established in Earth science, and recent theoretical works (Ménétrier et al., 2015; Morzfeld & Hodyss, 2023; Flowerdew, 2015) address this issue as well.

Covariance estimation is also a fundamental problem in statistics. Theoretical aspects of localization in NWP, for example, are described by Furrer and Bengtsson (2007) and Bickel and Levina (2008), and a review of various covariance estimation methods is provided by Pourahmadi (2011). The textbook by Wainwright (2019) emphasizes the difficulty of estimating a covariance matrix when the ensemble size is small. As representatives of the many statistical techniques that have been created over the years, we consider a soft-thresholding method (Wainwright, 2019), the graphical Lasso (G-Lasso, Friedman et al. (2007)), convex sparse Cholesky selection (CSCS, Khare et al. (2019)), and sparse covariance estimation (Xue et al., 2012).

#### 2.4.1 Prior Optimal Localization

The idea of optimal localization is to find a Hadamard product estimator, defined by the matrix  $\mathbf{L}$ , that minimizes the cost function

$$F_{\text{POLO}}(\mathbf{L}) = \left\| \langle \mathbf{L} \circ \hat{\mathbf{P}} - \mathbf{P}_{n_e \rightarrow \infty} \rangle \right\|_{\text{Fro}}, \quad (13)$$

where  $\mathbf{P}_{n_e \rightarrow \infty}$  is the “true” covariance matrix one would obtain from an infinite ensemble and where the brackets  $\langle \cdot \rangle$  denote an expected value over ensemble draws (Ménétrier et al., 2015; Morzfeld & Hodyss, 2023; Flowerdew, 2015);  $\|\cdot\|_{\text{Fro}}$  is the Frobenius norm, i.e., the square root of the sum of the squares of all elements of a matrix. Under Gaussian assumptions, one can solve this optimization analytically to obtain

$$[\mathbf{L}]_{ij} = \frac{\rho_{ij}^2(n_e - 1)}{1 + \rho_{ij}^2 n_e}, \quad (14)$$

which we refer to as prior optimal localization (POLO). Here,  $\rho_{ij}$  is the true correlation between the variables with indices  $i$  and  $j$ . While POLO does *not* rely on a spatial decay of correlations, it assumes that the correlations are known. POLO is, therefore, not



297 a viable algorithm but it can be used as a benchmark for practical algorithms. Empir-  
 298 ical localization functions (ELF) are closely related to optimal localization and are im-  
 299 plemented based on the idea of learning a localization matrix from training/simulation  
 300 data (Anderson & Lei, 2013).

301 POLO does *not* guarantee a PSD estimator. To see why, consider a theorem in lin-  
 302 ear algebra: If one applies a function element-wise to a PSD matrix whose elements are  
 303 in  $(0, 1)$ , the only functions that always preserve semi-definiteness have a power series  
 304 representation with non-negative coefficients (Schoenberg, 1942; Guillot & Rajaratnam,  
 305 2015). The POLO matrix in (14) does not satisfy this theorem and, hence, the matrix  
 306  $\mathbf{L}$  is not guaranteed to be PSD, which in turn implies that the POLO covariance esti-  
 307 mator is not guaranteed to be PSD. Indeed, we routinely observe non-PSD POLO esti-  
 308 mates in the numerical examples in Section 4.

### 309 2.4.2 Sampling Error Corrections and Power Law Corrections

310 Anderson (2012) introduces the terminology and methodology of *sampling error*  
 311 *correction* (SEC). SEC constructs covariance corrections quite similarly to POLO, but  
 312 the SEC corrections are based on numerical experiments with “training data” and groups  
 313 of ensembles, so that the correction depends on the sample correlation, rather than on  
 314 the true correlation (compare Figure 1(b) of this paper with Figure 1 of Anderson (2012)).

315 Lee (2021b) subsequently noticed that the corrections may be efficiently approx-  
 316 imated by a power law. Specifically, let  $\hat{\rho}$  be the empirical estimate of the ensemble *cor-*  
 317 *relations* and define the PLC estimator of the correlations by

$$318 \hat{\rho}_{\text{PLC}} = \mathbf{L}(\beta) \circ \hat{\rho}, \quad (15)$$

319 where the elements of the matrix  $\mathbf{L}(\beta)$  are given by

$$320 [\mathbf{L}(\beta)]_{ij} = |[\hat{\rho}]_{ij}|^\beta, \quad (16)$$

321 i.e., we raise the absolute values of the empirical correlations *element-wise* to the power  
 322  $\beta$ . The exponent  $\beta$  is a tunable parameter. Once we have selected a suitable  $\beta$ , we ob-  
 323 tain the covariance estimator

$$324 \hat{\mathbf{P}}_{\text{PLC}} = \hat{\mathbf{V}} \hat{\rho}_{\text{PLC}} \hat{\mathbf{V}}, \quad (17)$$

325 where  $\hat{\mathbf{V}}$  is a  $n \times n$  diagonal matrix whose diagonal elements are the ensemble standard  
 326 deviations. For the rest of this paper, we refer to this algorithm as “power law correc-  
 327 tions” (PLC).

328 PLC does not guarantee a PSD covariance estimator: one can apply the same the-  
 329 orems and logic as outlined above when discussing the PSD property in the context of  
 330 POLO. The PLC correlation estimate, however, *is* positive semi-definite if the exponent  
 331 is “large enough.” To understand why, we derive lower bounds for the eigenvalues of  $\mathbf{L}(\beta)$   
 332 using Gershgorin’s circle theorem. The theorem implies that an eigenvalue,  $\lambda$ , of  $\mathbf{L}(\beta)$   
 333 satisfies the inequalities

$$334 1 - Z_i \leq \lambda \leq 1 + Z_i, \quad (18)$$

335 where  $Z_i$  is the sum of the absolute values of the off-diagonal elements of a row (or col-  
 336 umn) of  $\mathbf{L}(\beta)$ :

$$337 Z_i = \sum_{i \neq j} |\hat{\rho}_{ij}|^\beta. \quad (19)$$

338 If we pick the exponent  $\beta$  to guarantee that  $Z_i \leq 1$  for all  $i$  (all rows of  $\mathbf{L}(\beta)$ ), then Ger-  
 339 shgorin’s theorem implies positive semi-definiteness of the matrix  $\mathbf{L}(\beta)$  and, via the Schur  
 340 product theorem, positive semi-definiteness of the PLC estimator. In our examples, and  
 341 with our adaptive strategy for choosing the exponent  $\beta$  (see Section 3.4.1), we never ran  
 342 into trouble with definiteness of the estimators, but we cannot guarantee that this is gen-  
 343 erally the case.

344 Covariance estimation using powers of ensemble correlations is also at the core of  
 345 a method called ECO-RAP (ensemble correlations raised to a power, Bishop and Hodyss  
 346 (2009a, 2009b, 2007)). In ECO-RAP, only positive, even exponents are considered, which  
 347 ensures that the ECO-RAP estimator is PSD, and that ECO-RAP, embedded within an  
 348 ensemble transform approach, is scalable to high-dimensional problems.

### 349 **2.4.3 Soft-Thresholding**

350 The idea of thresholding is to set small covariances to zero. This can be achieved  
 351 by applying the soft-thresholding function

$$352 \quad T_\lambda(s) = \begin{cases} s - \lambda \operatorname{sign}(s) & \text{if } |s| > \lambda \\ 0 & \text{otherwise} \end{cases}, \quad (20)$$

353 element-wise to the sampling covariance matrix, so that the soft-thresholding covariance  
 354 estimate is

$$355 \quad [\hat{\mathbf{P}}_{\text{ST}}]_{ij} = T_\lambda([\hat{\mathbf{P}}]_{ij}). \quad (21)$$

356 Here,  $\lambda$  is a positive scalar. Soft-thresholding has favorable asymptotic properties (Wainwright,  
 357 2019) and is computationally simple to implement, but the soft-thresholding covariance  
 358 estimator is not always PSD (Khare et al., 2019). The parameter  $\lambda$  is usually determined  
 359 via a tuning process. In Section 3.4.3, we describe how to find this parameter adaptively.

### 360 **2.4.4 Sparse Covariance Estimation**

361 Xue et al. (2012) note that soft-thresholding corresponds to the minimizer of the  
 362 cost function

$$363 \quad F_{\text{Soft Thres.}}(\mathbf{P}) = \frac{1}{2} \|\mathbf{P} - \hat{\mathbf{P}}\|_{\text{Fro}}^2 + \lambda \sum_{j \neq k} |\mathbf{P}_{jk}|, \quad (22)$$

364 where  $\hat{\mathbf{P}}$  is the empirical covariance matrix. The authors then describe an algorithm to  
 365 minimize the cost function (22) subject to the constraint that  $\mathbf{P} \geq \epsilon \mathbf{I}$  (i.e., the matrix  
 366  $\mathbf{P} - \epsilon \mathbf{I}$  is PSD), where  $\mathbf{I}$  is the identity matrix and where  $\epsilon > 0$  is a nuisance parame-  
 367 ter that can be set to a small number ( $10^{-5}$  is suggested). The constraint guarantees that  
 368 the covariance estimator is symmetric positive definite. Moreover, the estimator is sparse  
 369 because large off-diagonal elements are penalized and the 1-norm drives small covariances  
 370 to zero. This means that this technique, which we call *sparse covariance estimation*, is  
 371 most applicable in situations where one expects that most covariances should be zero.  
 372 We note that sparse covariance estimation requires tuning to find an appropriate reg-  
 373 ularization strength  $\lambda$ . In Section 3.4.4, we explain how to find the regularization strength  
 374 adaptively.

### 375 **2.4.5 Graphical Lasso**

376 Soft-thresholding and sparse covariance estimation find sparse estimates of the cov-  
 377 ariance matrix, i.e., the underlying assumption is that the majority of the covariances  
 378 are equal to zero. One can also search for a covariance matrix whose *inverse* is sparse.  
 379 The inverse of the covariance matrix is called the precision matrix,  $\Theta = \mathbf{P}^{-1}$ . The graph-  
 380 ical Lasso (G-Lasso, Friedman et al. (2007)) finds an estimator of the precision matrix  
 381  $\Theta$  by minimizing the cost function

$$382 \quad F_{\text{G-Lasso}}(\Theta) = \operatorname{tr}(\hat{\mathbf{P}}\Theta) - \log \det(\Theta) + \lambda \sum_{j,k} |\Theta_{jk}|, \quad (23)$$

383 over all PSD matrices  $\Theta$ . Here,  $\hat{\mathbf{P}}$  is the empirical covariance matrix and  $\lambda$  is a regular-  
 384 ization strength, so that large  $\lambda$  promote sparsity of the precision matrix estimate. Note  
 385 that minimizing (23) over all PSD matrices guarantees that the precision matrix esti-  
 386 mate is PSD, which in turn guarantees that the covariance matrix estimate is PSD. On

387 the other hand, a sparse precision matrix does not, in general, guarantee a sparse covari-  
 388 ance matrix, so the underlying assumptions of the G-Lasso and sparse covariance esti-  
 389 mation or soft-thresholding are quite different (Bickel & Lindner, 2012; Morzfeld et al.,  
 390 2019). The G-Lasso can be computationally expensive because (i) the optimization prob-  
 391 lem (23) is non-trivial; (ii) the method requires tuning to find an appropriate  $\lambda$ .

#### 392 **2.4.6 Convex Sparse Cholesky Selection**

393 Khare et al. (2019) describe a method called *convex sparse Cholesky selection* (CSCS),  
 394 which works with the triangular Cholesky factor  $\mathbf{A}$  of the precision matrix  $\Theta = \mathbf{A}^T \mathbf{A}$ .  
 395 Specifically, the goal is to find a sparse Cholesky factor by minimizing the cost function

$$396 F_{\text{CSCS}}(\mathbf{A}) = \text{tr}(\mathbf{A}^T \mathbf{A} \hat{\mathbf{P}}) - 2 \log \det(\mathbf{A}) + \lambda \sum_{1 \leq j < i} |\mathbf{A}_{ij}|, \quad (24)$$

397 where  $\lambda > 0$ . Due to the Cholesky factorization, the CSCS method guarantees that the  
 398 resulting estimators of the precision or covariance matrices are PSD.

### 399 **2.5 Morozov’s Discrepancy Principle**

400 Morozov’s discrepancy principle is a technique to adjust regularization parameters  
 401 in inverse problems (Morozov, 1984; Anzengruber & Ramlau, 2009). Suppose that we  
 402 are interested in solving the inverse problem whose cost function is

$$403 F_{\alpha}(\mathbf{x}) = \frac{1}{2} \|\mathbf{y} - f(\mathbf{x})\|_2^2 + \frac{\alpha}{2} \|\mathbf{x}\|_2^2, \quad (25)$$

404 where  $\mathbf{y}$  are the data,  $\mathbf{x}$  is a vector of unknowns,  $f(\cdot)$  is a nonlinear function (forward  
 405 model) and  $\alpha$  is a regularization parameter. Solving the inverse problems amounts to  
 406 minimizing the cost function. We denote the solution of the inverse problem for a given  
 407  $\alpha$  as  $\mathbf{x}_{\alpha}^*$ . The discrepancy principle determines the regularization parameter to be the  
 408 largest value of  $\alpha$  such that

$$409 \|\mathbf{y} - f(\mathbf{x}_{\alpha}^*)\|_2 \leq S, \quad (26)$$

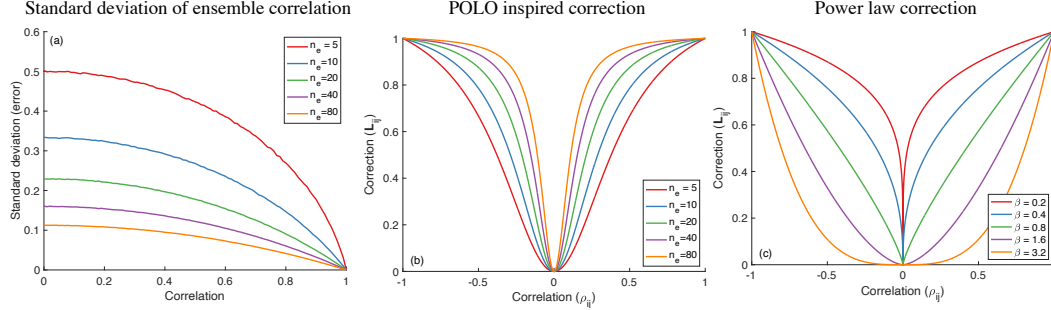
410 where the scalar  $S$  describes the “noise level” in the problem. For example, if the errors  
 411 in the data are described by Gaussian noise, then  $S$  is derived from the variances of that  
 412 noise. Application of Morozov’s discrepancy principle in practice requires that we solve  
 413 a sequence of inverse problems, parameterized by  $\alpha$ , to find the regularization param-  
 414 eter that leads to a solution that is compatible with the assumed noise level. These ideas  
 415 can also be used to obtain “regularized” covariance estimates, as we will explain below.

## 416 **3 New Methods for Noise-Informed Covariance Estimation**

417 Our goal is to design a Hadamard product estimator as in (2), which means that  
 418 we must build a correction matrix  $\mathbf{L}$ . Our design must go beyond assuming a spatial de-  
 419 cay of correlations, because this assumption is not reasonable in many cases. The de-  
 420 sign must also adapt itself to diverse situations in order to minimize tuning. We focus  
 421 on correcting correlations, and we estimate variances directly from the ensemble. This  
 422 is common in NWP (Whitaker & Hamill, 2012; Hodyss, Campbell, & Whitaker, 2016;  
 423 Gharamti et al., 2019) and perhaps intuitive because correlations are naturally scaled  
 424 to the interval  $[-1, 1]$ .

### 425 **3.1 Motivation: Damp Small Correlations More Heavily than Larger** 426 **Ones**

427 We base the design of our new method on a basic fact about estimating correla-  
 428 tions: estimating small correlations is notoriously difficult, and estimating large corre-  
 429 lations is, by comparison, easy. One way to understand this fact is to generate ensem-  
 430 bles of bivariate Gaussian random variables with varying degrees of correlation and then



**Figure 1.** (a) Standard deviation of ensemble correlation as a function of correlation (adapted from Figure 1 of Lee (2021b)). (b) POLO inspired correction factor, shown as a function of correlation for different ensemble sizes. (c) Power law correction factor, as proposed by Lee (2021b), shown as a function of correlation for different choices of the exponent  $\beta$ .

431 compute the ensemble correlation. Repeating this process many times allows us to com-  
 432 pute the standard deviation in the correlation estimate as a proxy for the error we should  
 433 expect in the correlation estimate (Lee, 2021b; Anderson, 2012). The average standard  
 434 deviation (averaged over independent ensemble draws) as a function of the “true” cor-  
 435 relation is shown in Figure 1(a), for several ensemble sizes. We note that the standard  
 436 deviation, or expected error, in the correlation estimate is large if the “true” correla-  
 437 tion is small. This means that small correlations are usually not trustworthy, unless the en-  
 438 semble size is huge. Consequently, it is natural to damp small correlations because it is  
 439 nearly impossible to distinguish “true” small correlations from sampling error. Large cor-  
 440 relations, on the other hand, are usually trustworthy, even if the ensemble size is small.  
 441 In fact a correlation equal to one should *always* be trusted—the standard deviation goes  
 442 to zero as the correlation goes to one. This simple numerical experiment thus tells us  
 443 that a reasonable correlation correction should damp small correlations more heavily than  
 444 large correlations. The larger error in estimating small correlations is a known feature  
 445 of the sampling distribution of the correlation coefficient between Gaussian random vari-  
 446 ables (Flowerdew, 2015).

447 POLO reiterates the idea that one can usually “trust” large correlations and that  
 448 small correlations should be damped. To see why, note that if we re-scale the POLO cor-  
 449 rection (14) so that correlations equal to one are uncorrected, we obtain

$$450 \quad [\mathbf{L}]_{ij} = \frac{(n_e + 1)\rho_{ij}^2}{1 + \rho_{ij}^2 n_e}. \quad (27)$$

451 This re-scaled correction factor is shown as a function of correlation in Figure 1(b), and  
 452 we see that it mimics the ideas described just above. At any ensemble size, small cor-  
 453 relations are subject to a stronger correction than large ones.

454 Power law corrections (PLC, Lee (2021b)) and “ensemble correlation raised to a  
 455 power” (ECO-RAP, Bishop and Hodyss (2009a, 2009b, 2007)) are also based on the sim-  
 456 ple fact that one should damp small correlations more severely than larger ones. This  
 457 is illustrated in Figure 1(c), where we show PLC correction factors ( $|\rho|^\beta$ ) for a few choices  
 458 of  $\beta$ . Moreover, PLC nicely resembles the SEC of Anderson (2012) (compare Figures 1(b),(c)  
 459 with Figure 1 of Anderson (2012)).

### 460 3.2 Noise-Informed Covariance Estimation

461 Our new covariance estimator is based on the simple idea that small correlations  
 462 should be reduced more heavily than large correlations, which we implement by adapt-

463 ing ideas from power law corrections. Additionally, we make use of an understanding of  
 464 sampling error (noise) in empirical correlations to make the method adaptive. The use  
 465 of uncertainties leads to the name of the method, “noise-informed covariance estimation”  
 466 (NICE).

467 NICE requires some *a priori* work that will be used to define the noise level within  
 468 the correlation estimates. Following the ideas described in Figure 1(a), we use (offline)  
 469 numerical experiments to determine a standard deviation associated with a “grid” of em-  
 470 pirical correlations (using bivariate Gaussian random variables, see Section 3.1). We then  
 471 form a lookup table so that we can assign a standard deviation to *any* empirical corre-  
 472 lation via interpolation.

473 After the offline work, the first actual step of NICE is to compute the  $n$  empirical  
 474 ensemble standard deviations, and the  $n(n-1)/2$  empirical ensemble correlations, which  
 475 we compile in a symmetric  $n \times n$  correlation matrix  $\hat{\rho}$  (with ones on the diagonal). The  
 476 sum total noise level, which we call  $S_\rho$ , is defined as follows. Using the lookup table, we  
 477 can assign a standard deviation  $\sigma_{\rho_{ij}}$  to each correlation  $\hat{\rho}_{ij}$  in the matrix  $\hat{\rho}$ , with the un-  
 478 derstanding that the standard deviation is zero if the correlation is one. The noise level  
 479  $S_\rho$  is a sum of all noises in the empirical estimate of the correlations:

$$480 \quad S_\rho = \sqrt{\sum_{i=1}^n \sum_{j=1}^n (\sigma_{\rho_{ij}})^2}. \quad (28)$$

481 In the second step, we use Morozov’s discrepancy principle, applied to the estima-  
 482 tion of correlation matrices. The “data” are the empirical estimates of the correlations  
 483  $\hat{\rho}$ , and the preliminary correlation estimate is

$$484 \quad \hat{\rho}_\gamma = \hat{\rho}^{\circ\gamma} \circ \hat{\rho}, \quad (29)$$

485 where  $\gamma$  is a positive, even integer. The elements of the matrix  $\hat{\rho}^{\circ\gamma}$  are  $[\hat{\rho}^{\circ\gamma}]_{ij} = ([\hat{\rho}]_{ij})^\gamma$ ,  
 486 i.e., we raise the empirical correlations *element-wise* to an even, positive power  $\gamma$ . Mo-  
 487 rozov’s discrepancy principle suggests to pick  $\gamma$  such that

$$488 \quad \|\hat{\rho} - \hat{\rho}_\gamma\|_{\text{Fro}} \leq \delta S_\rho, \quad (30)$$

489 where the scalar  $\delta$  is a tunable factor which we usually set to be equal to one (see nu-  
 490 merical examples in Section 4, for cross-covariances in EKI we set  $\delta = 0.5$ ). Specifically,  
 491 we pick the smallest even, positive integer  $\gamma^*$  that violates the discrepancy principle so  
 492 that

$$493 \quad \|\hat{\rho} - \hat{\rho}_{\gamma^*}\|_{\text{Fro}} \geq \delta S_\rho, \quad (31)$$

494 This procedure determines an exponent  $\gamma^*$  that leads to a correlation matrix estimate  
 495 that is PSD ( $\gamma^*$  is positive and even) and too strongly regularized according to the dis-  
 496 crepancy principle.

497 The third and final step linearly interpolates between a correction matrix that is  
 498 too strong (power  $\gamma^*$ ) and a correction with a smaller even integer (power  $\gamma^*-2$ ), which  
 499 is ostensibly “too weak”:

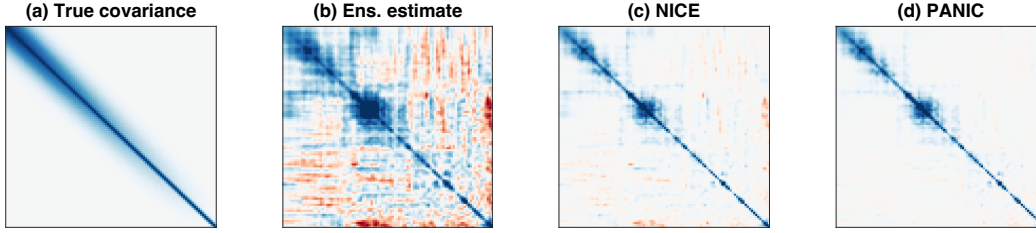
$$500 \quad \mathbf{L}(\alpha) = \alpha \hat{\rho}^{\circ\gamma^*} + (1 - \alpha) \hat{\rho}^{\circ(\gamma^*-2)}. \quad (32)$$

501 The associated correlation estimate is

$$502 \quad \hat{\rho}_\alpha = \mathbf{L}(\alpha) \circ \hat{\rho}. \quad (33)$$

503 The discrepancy principle then determines the interpolation factor  $\alpha$ . Specifically, we  
 504 find  $\alpha^*$  to be the largest  $\alpha \in [0, 1]$  such that

$$505 \quad \|\hat{\rho} - \hat{\rho}_{\alpha^*}\|_{\text{Fro}} \leq \delta S_\rho, \quad (34)$$



**Figure 2.** (a) The true covariance matrix. (b) Empirical estimate of the covariance matrix. (c) NICE approximation of the covariance matrix (Section 3.2). (d) PANIC approximation of the covariance matrix (Section 3.3). All estimation methods use  $n_e = 20$  samples. The colormap is red for  $-1$ , white for  $0$  and blue for  $1$ .

506 i.e., we determine the largest PSD correction that satisfies the discrepancy principle. The  
 507 resulting, corrected correlation estimate is

$$508 \quad \hat{\rho}_{\text{nice}} = \mathbf{L}(\alpha^*) \circ \hat{\rho}, \quad (35)$$

509 which in turn leads to the covariance estimate

$$510 \quad \hat{\mathbf{P}}_{\text{nice}} = \hat{\mathbf{V}} \hat{\rho}_{\text{nice}} \hat{\mathbf{V}}, \quad (36)$$

511 where  $\hat{\mathbf{V}}$  is a  $n \times n$  diagonal matrix whose diagonal elements are the ensemble standard  
 512 deviations.

513 We can summarize NICE in the following steps.

- 514 1. Compute the empirical correlations  $\hat{\rho}$  and empirical standard deviations.
- 515 2. Determine the noise level  $S_\rho$  via a lookup table and equation (28).
- 516 3. Determine the smallest positive, even integer  $\gamma^*$  that violates the discrepancy prin-  
 517 ciple (31).
- 518 4. Determine the largest interpolation factor  $\alpha^*$  that satisfies the discrepancy prin-  
 519 ciple (34).
- 520 5. Perform the element-wise correction of the correlation matrix in (35).
- 521 6. Use the corrected correlation matrix along with the empirical variances to com-  
 522 pute the covariance estimate via (36)

523 The effects of NICE are illustrated in Figure 2, where it is applied to estimate a  
 524  $100 \times 100$  covariance matrix, used by Bishop et al. (2017) to study localization in the  
 525 context of satellite data assimilation (compare our Figure 2(a) to Figure 2 in Bishop et  
 526 al. (2017)). We show the true covariance in Figure 2(a), the empirical estimate in Fig-  
 527 ure 2(b), and the NICE estimator in Figure 2(c). All approximations use the same en-  
 528 semble of size  $n_e = 20$ . The empirical estimate is noisy (large off-diagonal elements rep-  
 529 resent spurious covariances) and NICE improves on the empirical estimate by damping  
 530 small correlations.

### 531 **3.2.1 Implementation Details and Positive Semi-Definiteness**

532 Step 1 limits the applicability of NICE in extremely high dimensions because we  
 533 assume that *all* empirical correlations can be computed. As is, NICE can be applied to  
 534 problems with thousands of unknowns (which we demonstrate in numerical experiments),  
 535 but it may be computationally expensive if the dimension is  $10^5$  or larger. When used  
 536 in EnKF or EKI, one may be able to push these limitations further if the number of ob-

537 observations is relatively small (see Section 4.3), or if NICE is incorporated within an en-  
 538 semble transform framework (as in ECORAP, Bishop and Hodyss (2009a, 2009b)), or  
 539 serial filters (Anderson, 2001), or hybrid DA.

540 Step 2 is trivial, unless the dimension is huge (see comments above about Step 1).  
 541 For Step 3, we first try  $\gamma = 2$  and check the discrepancy principle. If it is violated, we  
 542 have found  $\gamma^*$  and move to Step 3. If not, we try  $\gamma = 4$  and so on. In the examples be-  
 543 low, a correction with  $\gamma^* = 6$  (or less) was always sufficient, meaning that we need about  
 544 three (or less) simple iterations to determine  $\gamma^*$ . Moreover, note that if  $\gamma^* = 2$  is se-  
 545 lected, then step four interpolates between the element-wise power two and the power  
 546 zero (no correction). For Step 4, we try a small  $\alpha$  and gradually increase it (line search)  
 547 until we violate the discrepancy principle, which then defines the “optimal”  $\alpha^*$  to be the  
 548 previous  $\alpha$  we just tried. Alternatively, a root-finding algorithm (e.g., the bisection method)  
 549 could be used.

550 We note that instead of a lookup table, one can also directly estimate the noise level  
 551  $S_\rho$  in (28) using the Fisher transformation. The distribution of the sample correlation  
 552 coefficient  $\hat{\rho}_{ij}$  between normally distributed variables is such that, when the Fisher trans-  
 553 formation is taken,

$$554 \quad z_{ij} = \operatorname{arctanh}(\hat{\rho}_{ij}) = \frac{1}{2} \log \left( \frac{1 + \hat{\rho}_{ij}}{1 - \hat{\rho}_{ij}} \right), \quad (37)$$

555 we have that for  $n_e > 3$ ,

$$556 \quad z_{ij} \overset{\text{approx}}{\sim} \mathcal{N} \left( \operatorname{arctanh}(\rho_{ij}), \frac{1}{n_e - 3} \right), \quad (38)$$

557 where  $\rho_{ij}$  is the true correlation (see, e.g., Flowerdew (2015)). Thus, we can estimate  
 558 the standard deviation of  $\hat{\rho}_{ij}$  as follows. Taking  $\hat{\rho}_{ij}$  as an estimate of  $\rho_{ij}$ , we draw  $m$  sam-  
 559 ples  $z_{ij}$  from the above Gaussian distribution, but replacing  $\operatorname{arctanh}(\rho_{ij})$  with  $\operatorname{arctanh}(\hat{\rho}_{ij})$   
 560 in the mean. Second, we apply the inverse Fisher transformation  $\tanh(z_{ij})$  to each of the  
 561 samples and compute their standard deviation. This strategy of computing the noise level  
 562 in the correlations is attractive because it is (i) easy; and (ii) it avoids having to pre-compute  
 563 lookup tables. The lookup tables, however, have a slight edge over the Fisher transfor-  
 564 mation approach in terms of their online cost.

565 Finally, the positive semi-definiteness of the correlation estimator,  $\hat{\boldsymbol{\rho}}_{\text{nice}}$ , follows from  
 566 basic facts about Hadamard products. Specifically, raising the elements of a PSD ma-  
 567 trix to an even power preserves definiteness, and the sum of two PSD matrices is PSD.  
 568 The positive semi-definiteness of the covariance estimator  $\mathbf{P}_{\text{nice}}$  follows from the fact that  
 569 a PSD correlation matrix leads to a PSD covariance matrix.

### 570 **3.3 Partially Adaptive Noise-Informed Covariance (PANIC)**

571 In some problems, e.g., in NWP, one may be in the situation where details of the  
 572 correlation structure are not well-understood, but one may be quite certain that corre-  
 573 lations should decay at far distances. For example, Bishop and Hodyss (2011) use a “par-  
 574 tially adaptive” method which combines an adaptive localization matrix with a tuned  
 575 (non-adaptive) localization matrix that eliminates correlations in the far-field. If the prob-  
 576 lem indeed has this structure (far-field being uncorrelated), then adding this informa-  
 577 tion should increase the accuracy of the estimator because small sampling errors in the  
 578 far-field accumulate to large errors in high-dimensions (Hodyss & Morzfeld, 2023; Morzfeld  
 579 & Hodyss, 2023).

580 One can easily combine these ideas with NICE. Since the resulting method requires  
 581 some tuning, it is “partially adaptive” (using the language in Bishop and Hodyss (2011))  
 582 and we call the method PANIC (**p**artially **a**daptive **n**oise **i**nformed **c**ovariance). PANIC  
 583 amounts to localizing the NICE estimator. Specifically, we use a localization matrix  $\mathbf{L}(\ell)$ ,

584 that depends on a length scale  $\ell$ , to obtain

$$585 \quad \hat{\rho}_{\text{panic}} = \mathbf{L}(\ell) \circ \hat{\rho}_{\text{nice}}. \quad (39)$$

586 Here, the length scale  $\ell$  is chosen a priori to be “large enough” to be certain that cor-  
587 relations beyond that length scale are physically implausible. With the correlation esti-  
588 mate we obtain the covariance matrix in the usual way via

$$589 \quad \hat{\mathbf{P}}_{\text{panic}} = \hat{\mathbf{V}} \hat{\rho}_{\text{panic}} \hat{\mathbf{V}}, \quad (40)$$

590 where  $\hat{\mathbf{V}}$  is a  $n \times n$  diagonal matrix whose diagonal elements are the ensemble standard  
591 deviations. Figure 2(d) illustrates PANIC and compares it to NICE. We note that the  
592 PANIC estimator reduces spurious correlations in the far field, but in the near field, PANIC  
593 and NICE are quite similar by construction. Moreover, the PANIC estimator is PSD be-  
594 cause NICE generates a PSD covariance estimate which is subsequently localized (Schur  
595 product with a PSD localization matrix); both steps preserve symmetry and definite-  
596 ness.

### 597 **3.4 Other New Adaptive Covariance Estimation Methods**

598 Within NICE, we combine an understanding of the noise in empirical correlations  
599 with Morozov’s discrepancy principle and, for that reason, the method is adaptive and  
600 tuning-free. This idea extends to other covariance estimation methods as well, and we  
601 now describe how to make some existing covariance estimation methods adaptive.

#### 602 **3.4.1 Adaptive Power Law Correction**

603 PLC requires that one determines the exponent  $\beta$ . In adaptive PLC (Ad.-PLC),  
604 we use the largest (but not necessarily integer)  $\beta$  that satisfies the discrepancy princi-  
605 ple

$$606 \quad \|\hat{\rho} - \mathbf{L}(\beta) \circ \hat{\rho}\|_{\text{Fro}} \leq S_{\rho}. \quad (41)$$

607 Recall that  $\mathbf{L}(\beta)$  is a matrix whose elements are the absolute values of the empirical cor-  
608 relations raised to the power  $\beta$ :  $[\mathbf{L}(\beta)]_{ij} = |[\hat{\rho}]_{ij}|^{\beta}$ . For that reason, Ad.-PLC does *not*  
609 guarantee positive semi-definiteness of the covariance estimator (just as PLC). In our  
610 numerical examples, however Ad.-PLC always leads to PSD covariance estimators, be-  
611 cause the adaptive strategy picks out exponents that are large enough to ensure that the  
612 matrix is PSD (see Section 2.4.2).

#### 613 **3.4.2 Adaptive Localization**

614 In “traditional” localization, we define a localization matrix by a length scale  $\ell$  that  
615 controls the decay of correlations. In adaptive localization (Ad.-Loc), we determine  $\ell$  to  
616 be the largest length scale that satisfies the discrepancy principle

$$617 \quad \|\hat{\rho} - \mathbf{L}(\ell) \circ \hat{\rho}\|_{\text{Fro}} \leq S_{\rho}. \quad (42)$$

618 In our numerical experiments below we use a simple line search over the length scale  $\ell$   
619 to find an optimal length scale.

#### 620 **3.4.3 Adaptive Soft-Thresholding**

621 Soft-thresholding requires that we determine the thresholding parameter  $\lambda$  in (20).  
622 In adaptive soft-thresholding (Ad.-ST), we correct *correlations* and determine the thresh-  
623 olding parameter  $\lambda^*$  to be the largest  $\lambda$  that satisfies the discrepancy principle

$$624 \quad \|\hat{\rho} - \hat{\rho}_{\lambda}\|_{\text{Fro}} \leq S_{\rho}, \quad (43)$$



625 where  $\hat{\rho}_\lambda$  is the empirical correlation matrix thresholded with parameter  $\lambda$ , i.e.,

$$626 \quad [\hat{\rho}_\lambda]_{ij} = T_\lambda([\hat{\rho}]_{ij}), \quad (44)$$

627 where  $T_\lambda(\cdot)$  is the soft-thresholding function in (20). With  $\lambda^*$  defined in this way, we ob-  
628 tain the Ad.-ST covariance estimator by

$$629 \quad \mathbf{P}_{\text{Ad.-ST}} = \hat{\mathbf{V}}\hat{\rho}(\lambda^*)\hat{\mathbf{V}}, \quad (45)$$

630 where  $\mathbf{V}$  is a diagonal matrix whose diagonal elements are the ensemble standard de-  
631 viations (as in NICE). Note that Ad.-ST, just like soft-thresholding, does not guaran-  
632 tee a PSD estimate.

### 633 **3.4.4 Adaptive Sparse Covariance Estimation**

634 The sparse covariance estimation algorithm (Xue et al., 2012), which we briefly de-  
635 scribe in Section 2.4, finds a covariance estimate by minimizing the cost function (22)  
636 subject to the constraint that the estimator satisfies  $\mathbf{P}_{\text{ASC}} \geq \epsilon \mathbf{I}$ , which guarantees that  
637 the covariance estimator is PSD. The optimization problem can be solved efficiently, but  
638 the optimization problem depends on the regularization parameter  $\lambda$ , which defines the  
639 amount of sparsity in the estimate.

640 Adaptive sparse covariance estimation (ASCE) determines the regularization pa-  
641 rameter automatically. As noted by Xue et al. (2012), sparse covariance estimation and  
642 soft-thresholding are closely related, because sparse covariance estimation solves the same  
643 optimization problem as soft thresholding does, except with an added PSD constraint.  
644 Thus, we first perform adaptive soft-thresholding to find an optimal  $\lambda^*$ , and then per-  
645 form a single optimization with this  $\lambda^*$  to find a sparse correlation estimator  $\hat{\rho}_{\text{ASCE}}$  (note  
646 that we work exclusively with correlations, not covariances). The ASCE correlation es-  
647 timator defines the ASCE covariance estimator by

$$648 \quad \mathbf{P}_{\text{ASCE}} = \mathbf{V}\hat{\rho}_{\text{ASCE}}\mathbf{V}, \quad (46)$$

649 where  $\mathbf{V}$  is, as before, a diagonal matrix whose diagonal elements are the ensemble stan-  
650 dard deviations.

## 651 **4 Numerical Illustrations**

652 We compare NICE to a variety of competing methods, some new and some old. Specif-  
653 ically, we consider the following 13 methods for covariance estimation. We introduce ab-  
654 breviations for all methods that will be used in the numerical illustrations and in the Fig-  
655 ures.

### 656 **New adaptive methods**

- 657 1. Noise informed covariance estimation (**NICE**, Section 3.2)
- 658 2. Partially adaptive noise informed covariance (**PANIC**, Section 3.3).
- 659 3. Adaptive power law corrections (**Ad.-PLC**, Section 3.4.1).
- 660 4. Adaptive localization (**Ad.-Loc**, Section 3.4.2).
- 661 5. Adaptive soft-thresholding (**Ad.-ST**, Section 3.4.3).
- 662 6. Adaptive sparse covariance estimation (**ASCE**, Section 3.4.4).

### 663 **Methods for comparison**

- 664 7. The uncorrected, empirical estimate (**Ens.**) serves as the baseline for the improve-  
665 ment a more sophisticated covariance estimation can achieve.

- 666 8. **POLO** uses the correction matrix defined in Equation (14) with the “true” cor-  
 667 relations (see Section 2.4.1). Using POLO in this way describes a best-case sce-  
 668 nario, but we remind the reader that POLO is not a practical algorithm because  
 669 the true correlations are typically unknown (except in some of our synthetic nu-  
 670 merical illustrations).
- 671 9. POLO with ensemble correlations (**Ens.-POLO**) uses the correction matrix  $\mathbf{L}$  in (14),  
 672 but the correlations  $\rho_{ij}$  are uncorrected *empirical* correlations. This is perhaps  
 673 the simplest method of increasing the accuracy of the empirical covariance ma-  
 674 trix, but we will see that NICE and other methods are superior.
- 675 10. Localization (**Loc**) is implemented via a Gaussian localization whose elements are

$$676 \quad [\mathbf{L}]_{ij} = \exp(-(d_{ij}/\ell)^2), \quad (47)$$

677 where  $d_{ij}$  is the distance between grid points  $i$  and  $j$  and where the length scale  
 678  $\ell$  is tuned (see below for details). This is an example of the commonly used Hadamard  
 679 product localization in NWP, which relies on the assumption of a spatial decay  
 680 of correlation.

- 681 11. Power law corrections (**PLC**, Section 2.4.2), with tuned (non-integer) exponent  $\beta$ .
- 682 12. The Graphical Lasso (**G-Lasso**, Section 2.4) is implemented in Matlab code that  
 683 is available on GitHub (we downloaded the code at <https://gist.github.com/samwhitehall/6422598>).  
 684 The code yields the G-Lasso estimate of the precision matrix and we subsequently  
 685 compute its inverse to obtain an estimate of the covariance matrix. We tune the  
 686 regularization parameter of the G-Lasso in the same way as we tune localization  
 687 and PLC.
- 688 13. Convex sparse Cholesky selection (**CSCS**, Khare et al. (2019), see also Section 2.4)  
 689 gives a Cholesky factor of the inverse of the covariance matrix. As in the G-Lasso,  
 690 we use matrix inversion to find the covariance matrix. We tune the regularization  
 691 parameter in CSCS in the same way as we tune localization, PLC or G-Lasso.

692 The various covariance estimation techniques and some of their properties are sum-  
 693 marized in Table 1. All techniques, except G-Lasso, CSCS and ASCE can be used on  
 694 non-square correlation matrices (and cross-covariance matrices), which will become im-  
 695 portant in examples with EKI and in the geomagnetic data assimilation example.

696 We tune the localization (length scale  $\ell$ ), PLC (exponent  $\beta$ ), G-Lasso and CSCS  
 697 (regularization parameter) as follows. We perform a (large) number of training exper-  
 698 iments in which we vary the tunable parameter (line search). We then compute an av-  
 699 erage error and declare the parameter that leads to the smallest error as optimal. The  
 700 optimal parameter is used in subsequent experiments. We repeat the tuning for each nu-  
 701 merical example because the optimal tunable parameters are problem-dependent.

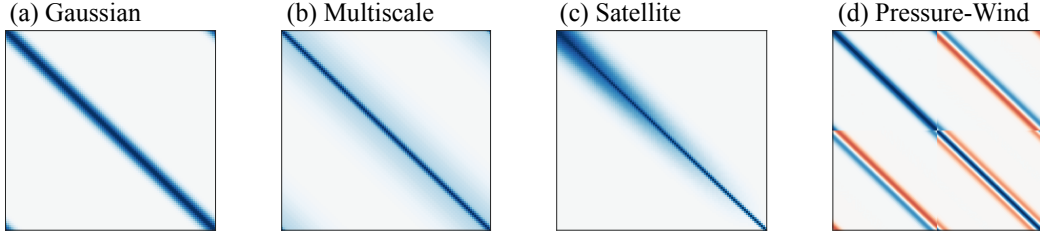
702 We use all 13 methods in our first set of numerical experiments with simple Gaus-  
 703 sians. Subsequently, we do not use methods that are computationally expensive and that  
 704 do not yield good results in simple experiments. G-Lasso and CSCS, for example, are  
 705 quite slow and do not perform well on our first set of simple tests. Other methods, e.g.,  
 706 PANIC, may not be applicable in subsequent examples because they presume a spatial  
 707 decay of correlation. Finally, POLO (with true correlations) can only be used in synthetic  
 708 scenarios where the correlations are known a priori, which is only true for our first set  
 709 of very simple experiments.

#### 710 4.1 Simple Gaussian Tests

711 We define a  $100 \times 100$  covariance matrix  $\mathbf{P}$  and draw  $n_e = 20$  ensemble mem-  
 712 bers from the corresponding Gaussian with mean zero. We then use NICE to estimate

**Table 1.** Summary of covariance estimation methods and their properties. In the assumptions column, “small corr. noisy” stands for the assumption that small correlations are noisy and, therefore reduced; “sparse cov.” stands for the assumption of a sparse covariance matrix; “known corr.” stands for the assumption that all correlations are known; “sparse inv. cov.” stands for the assumption that the inverse of the covariance matrix is sparse; “sparse Cholesky” stands for the assumption that a Cholesky factor of the covariance matrix is sparse. We assign a “low” computational cost if the technique performs simple operations on the elements of a covariance or correlation matrix. We assign a “high” computational cost if the technique solves optimization problems over (PSD) matrices. The computational cost is labeled “medium” for ASCE, which does perform an optimization over matrices but does so in a particularly speedy way (Xue et al., 2012).

Method	Adaptivity	Assumptions	PSD guarantees	Computational Cost
NICE	yes	small corr. noisy	yes	low
PANIC	yes	small corr. noisy+spatial decay of correlation	yes	low
Ad.-PLC	yes	small corr. noisy	no	low
PLC	no	small corr. noisy	no	low
Ad.-Loc	yes	spatial decay of correlation	yes	low
Loc	no	spatial decay of correlation	yes	low
Ad.-ST	yes	small corr. noisy	no	low
ASCE	yes	sparse cov.	yes	medium
POLO	-	known corr.	no	low
Ens.-POLO	-	none	no	low
Ens.	-	none	yes	low
G-Lasso	no	sparse inv. cov.	yes	high
CSCS	no	sparse Cholesky	yes	high



**Figure 3.** The covariance matrices used in Section 4.1. (a) Gaussian kernel. (b) Multi-scale kernel. (c) Covariance inspired by satellite data assimilation. (d) Covariance of two spatial fields (pressure and wind). Color indicates the matrix elements with blue corresponding to one, white to zero, and red to minus one.

713 the covariance matrix and measure the error in the estimate by

$$714 \text{ Error} = \frac{\|\hat{\mathbf{P}}_{\text{nice}} - \mathbf{P}\|_{\text{Fro}}}{\|\mathbf{P}\|_{\text{Fro}}}. \quad (48)$$

715 Since the error is random, we average over ensemble draws, and the average error indi-  
 716 cates an error we should typically expect. We use the same procedure to compute the  
 717 error of other covariance estimation methods.

718 We consider four different covariance matrices, illustrated in Figure 3.

1. *Gaussian kernel.* A covariance matrix  $\mathbf{P}$  with a Gaussian kernel is defined by the elements

$$[\mathbf{P}]_{ij} = \exp\left(-\frac{1}{2}\left(\frac{d_{ij}}{\ell}\right)^2\right),$$

719 where the length scale is  $\ell = 5$  and where  $d_{ij}$  is a periodic distance between the  
 720 grid points  $i$  and  $j$ . Note that this covariance has the same kernel function as the  
 721 localization matrix used during classical covariance localization (Loc).

2. *Multi-scale kernel.* A multi-scale covariance  $\mathbf{P}$  is defined by the superposition of two covariance matrices with Gaussian kernels and different length scales:

$$[\mathbf{P}]_{ij} = 0.7 \exp\left(-\frac{1}{2}\left(\frac{d_{ij}}{\ell_1}\right)^2\right) + 0.3 \exp\left(-\frac{1}{2}\left(\frac{d_{ij}}{\ell_2}\right)^2\right).$$

722 We chose the length scales to be  $\ell_1 = 2$  and  $\ell_2 = 20$  (Morzfeld & Hodyss, 2023;  
 723 Flowerdew, 2015).

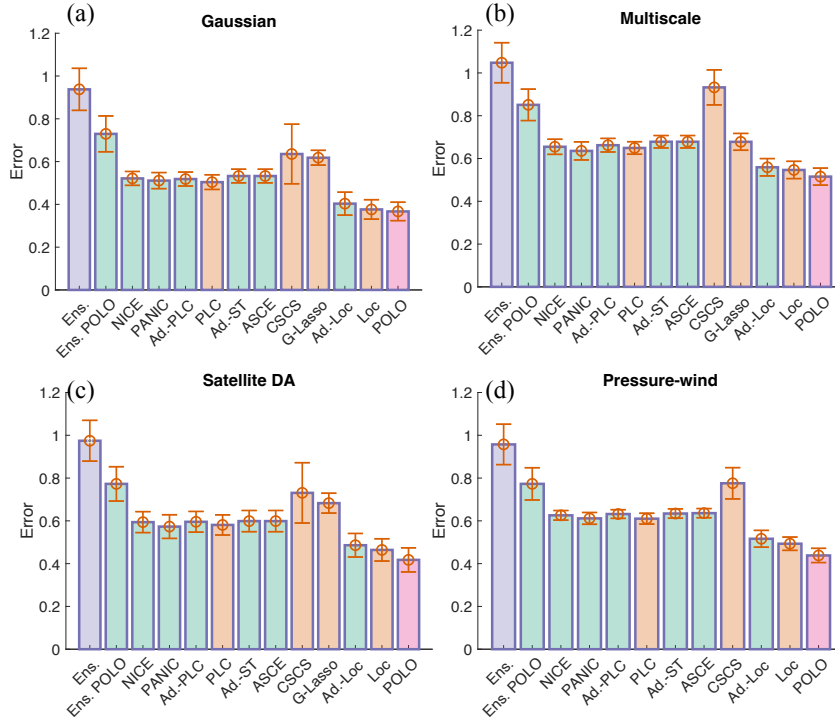
3. *Satellite data assimilation covariance.* Bishop et al. (2017) consider the covariance matrix

$$[\mathbf{P}]_{ij} = \sqrt{\frac{ij}{n^2}} \exp\left(-\frac{1}{2}\left(\frac{i-j}{\ell_1}\right)^2\right) + \sqrt{\left(1 - \frac{i}{n}\right)\left(1 - \frac{j}{n}\right)} \exp\left(-\frac{1}{2}\left(\frac{i-j}{\ell_2}\right)^2\right).$$

724 as a toy problem for satellite data assimilation. Following Bishop et al. (2017), we  
 725 chose the length scales to be  $\ell_1 = 1$  and  $\ell_2 = 8$ . Note that this covariance ma-  
 726 trix is “nonstationary” covariance because the elements of  $\mathbf{P}$  depend on  $i$  and  $j$ ,  
 727 not just of  $i - j$ .

4. *Pressure-wind covariance.* We consider two spatially extended fields  $u$  (pressure) and  $w$  (wind), related by a derivative such that  $w = du/dx$ . We assume that the pressure variable has a Gaussian covariance kernel with length scale  $\ell = 5$  and

730



**Figure 4.** Error (mean and one standard deviation error bars) in covariance matrix estimates for various covariance types with dimension  $n_x = 100$ . The ensemble size is  $n_e = 20$ . (a) Gaussian covariance kernel. (b) Multi-scale covariance kernel. (c) Satellite data assimilation covariance matrix. (d) Pressure-wind covariance matrix. The bar chart is color coded so that the vanilla method (Ens.) appears in blue, tuning-free/adaptive methods (Ens.-POLO, NICE, PANIC, Ad.-PLC, Ad.-ST, ASCE, Ad.-Loc) appear in green, tuned methods (PLC, CSCS, G-Lasso, Loc) appear in orange, and the infeasible method (POLO) appears in pink (rightmost bar in each panel).

731 we construct the covariance of  $w$ , as well as the cross covariances between  $u$  and  
 732  $w$ , using a finite difference operator (see Morzfeld and Hodyss (2023) for more details). We note that if both  $u$  and  $w$  have 100 components, the overall dimension  
 733 of the problem is  $n_x = 200$ .  
 734

735 We apply all 13 covariance estimation techniques listed above for all but the pressure-  
 736 wind covariance, for which we do not apply G-Lasso because the code runs very slowly  
 737 on this 200-dimensional problem. Note that all four covariance matrices we consider here  
 738 have exponentially small correlations in the far field (away from the diagonal), so that  
 739 the use of a localization and PANIC are appropriate. Results are summarized in Figure 4,  
 740 which shows the average error ( $10^3$  trials) for each method and covariance type along  
 741 with one standard deviation error bars. The numerical experiments support the follow-  
 742 ing conclusions.

- 743
- 744 1. For all four covariance types, *all* covariance estimation techniques are more ac-  
 745 curate than the sample covariance matrix, which always has the largest error.
  - 746 2. POLO with ensemble correlations (Ens.-POLO) improves the covariance estimates  
 in all four cases, but not to the extent of the other methods we tried.

- 747 3. NICE, Ad.-PLC, Ad.-ST and ASCE lead to similar errors which are in turn com-  
748 parable to the errors of a finely tuned PLC. The fact that all four adaptive meth-  
749 ods perform as well as a related finely tuned method suggests that the discrep-  
750 ancy principle and the pre-computed noise level are robustly applicable to adap-  
751 tive covariance estimation.
- 752 4. The adaptive localization (Ad.-Loc) leads to errors almost as small as the errors  
753 obtained by a finely tuned localization (Loc). This reiterates our previous point,  
754 i.e., that adapting localization/covariance estimation parameters via a discrepancy  
755 principle is a robust idea.
- 756 5. The errors of PANIC are slightly smaller than the errors of NICE, which suggests  
757 that reducing the (non-adaptive) far-field correlations has a positive effect.
- 758 6. Localization (Loc) comes close to the optimal errors obtained by POLO and Loc  
759 and POLO lead to the smallest errors in all four examples.
- 760 7. G-Lasso and CSCS lead to smaller errors than Ens.-POLO, but the errors are larger  
761 than for the new adaptive methods. G-Lasso and CSCS also require significantly  
762 more computations than the competing methods, and we conclude that G-Lasso  
763 and CSCS are not competitive in these examples. Recall, however that G-Lasso  
764 and CSCS are designed to estimate the precision matrix (not the covariance ma-  
765 trix as we do here). CSCS further targets applications with a natural ordering of  
766 the data.

767 During the trials of our experiments we monitored if a covariance matrix estimate  
768 was PSD or not. When the exponent in PLC was chosen adaptively (Ad.-PLC) or via  
769 tuning, we encountered *no* negative eigenvalues, while POLO, Ens. POLO, and Ad.-ST  
770 often produced non-PSD estimates. This is an interesting result because POLO is the  
771 estimator with the lowest errors and yet it is not always PSD. Our error metric here, how-  
772 ever, does not account for this deficiency, violating the PSD property may cause insta-  
773 bility within EnKFs or EKI (see Section 2.3).

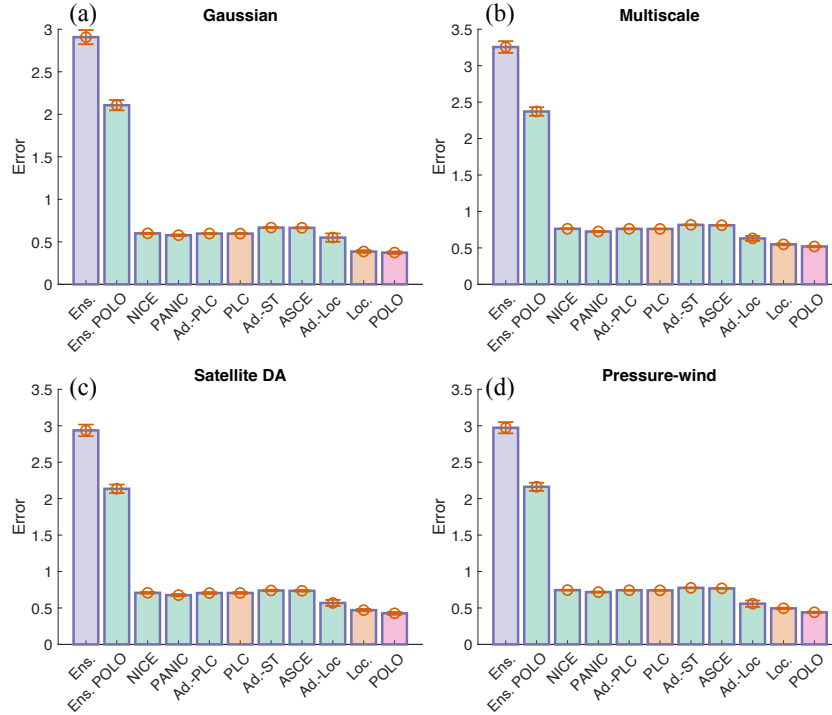
774 When we increase the dimension of the problem, the decrease in errors is more dra-  
775 matic (Hodyss & Morzfeld, 2023). Figure 5 summarizes results obtained for problems  
776 of dimension  $n = 1000$ . We note qualitatively the same results as in the  $n = 100$  di-  
777 mensional example: NICE, Ad.-PLC, Ad.-ST and ASCE are comparable and, even though  
778 these methods do *not* require tuning, they are as good as a tuned PLC. These four meth-  
779 ods, however, do not lead to errors as small as those obtained by localization (tuned or  
780 adaptive) or an optimal correction (POLO).

781 Finally, note that the correlations decay with distance in all above examples, which  
782 is exploited by classical (or adaptive) localization, but this correlation structure is *dis-*  
783 *covered* by the adaptive methods (NICE, Ad.-PLC, Ad.-ST and ASCE). Our first set  
784 of simple tests thus suggests that NICE, Ad.-PLC, Ad.-ST and ASCE can be viable op-  
785 tions in problems where assumptions about the underlying correlation structure are un-  
786 available or in problems where one wishes to reduce the tuning costs.

## 787 4.2 Cycling Data Assimilation Experiments with the Lorenz '96 model

788 We perform cycling data assimilation (DA) experiments with the Lorenz'96 model  
789 (L'96, Lorenz (1996)) and an ensemble Kalman filter (stochastic EnKF implementation,  
790 Burgers et al. (1998); Evensen (2009, 1994)). Specifically, we apply, within the EnKF,  
791 the covariance estimation methods NICE, PANIC, Ad.-PLC, ASCE, PLC, Ad.-Loc, lo-  
792 calization, and a version of POLO that indicates a best-case scenario at the expense of  
793 requiring a very large ensemble (hence being infeasible in practice). As is common in DA,  
794 we apply the covariance estimation (NICE, etc.) in conjunction with a covariance *infla-*  
795 *tion*. For the inflation, we simply set

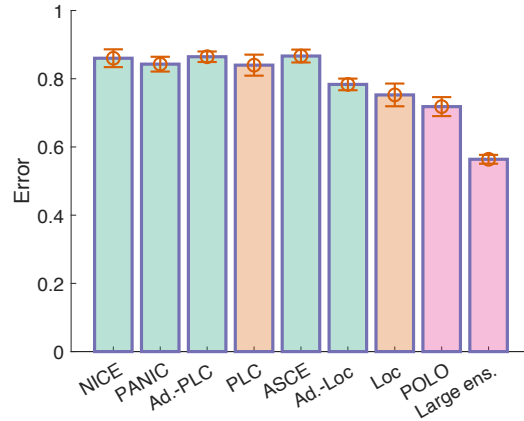
$$796 \mathbf{P} \leftarrow (1 + \kappa)\mathbf{P}, \quad (49)$$



**Figure 5.** Same as Figure 4, but with  $n = 1000$ . Error (mean and one standard deviation error bars) in covariance matrix estimates for various covariance types with dimension  $n = 1000$ . (a) Gaussian covariance kernel. (b) Multi-scale covariance kernel. (c) Satellite data assimilation covariance matrix. (d) Pressure-wind covariance matrix. The bar chart is color coded so that the vanilla method (Ens.) appears in blue, tuning-free methods (Ens.-POLO, NICE, PANIC, Ad.-PLC, Ad.-ST, ASCE, Ad.-Loc) appear in green, tuned methods (PLC, Loc) appear in orange, and the infeasible method (POLO) appears in pink (rightmost bars of each panel).

797 where  $\kappa > 0$  is an inflation parameter (tuned, see below).

798 The tuning of covariance estimation and/or the inflation is as follows. For the adap-  
 799 tive methods (NICE, Ad.-PLC, ASCE and Ad.-Loc), we only need to tune the inflation  
 800 parameter  $\kappa$ . For PANIC, we also only tune the inflation and set the length scale for the  
 801 spatial localization to  $\ell = 10$ , which is wide enough to expect that correlations beyond  
 802 that length scale are unreasonable. For localization, we tune the length scale jointly with  
 803 the inflation parameter  $\kappa$ . Similarly, for PLC we tune the exponent  $\beta$  jointly with the  
 804 inflation parameter  $\kappa$ . In all cases, the tuning is done by running 2,000 DA cycles, dis-  
 805 regarding the first 200 cycles as “spin-up,” and recording the associated, time-averaged  
 806 root mean square error (RMSE) for each inflation parameter and, if needed, additional  
 807 covariance estimation parameters. The parameters that lead to the smallest time-averaged  
 808 RMSE in the training experiment are subsequently used in another, independent exper-  
 809 iment in which we perform 1,000 DA cycles, disregard the first 100 cycles as spin-up, and  
 810 average RMSE after the spin-up period. Throughout the experiments, we hold the en-  
 811 semble size constant at  $n_e = 20$ . The state dimension is  $n = 40$  and we observe every  
 812 other variable, i.e., the number of observations is equal to 20. All observation error vari-  
 813 ances are equal to one. Observations are collected every 0.4 (dimensionless) time units  
 814 and the time step of the numerical integrator (a 4th order forward Runge-Kutta method)  
 815 is set to  $\Delta t = 0.05$  (this is the same setup as in Hodyss and Morzfeld (2023)).



**Figure 6.** Time-averaged analysis RMSE after spin up along with one standard deviation error bars in EnKF with various covariance estimation techniques. The bar chart is color coded so that tuning-free methods (NICE, PANIC, Ad.-PLC, ASCE, Ad.-Loc) appear in green, tuned methods (PLC, Loc) appear in orange, and infeasible methods (POLO and the large ensemble EnKF) appear in pink (two bars on the right).

816 To establish a best-case scenario, we use an EnKF *without* inflation or localization/  
 817 covariance estimation but with a large ensemble size  $n_e = 500$ . We further apply POLO  
 818 to an EnKF with  $n_e = 20$  (with tuned inflation), but run, in parallel, the large ensemble  
 819 size EnKF ( $n_e = 500$ ) to obtain the correlation information. These latter experi-  
 820 ments can indicate what a near-optimal localization may achieve (assuming the large en-  
 821 semble size EnKF reveals the main features of the “true” correlation).

822 The results of our numerical experiments are summarized in Figure 6, which shows  
 823 the time average of the analysis RMSE of EnKFs with various covariance estimation/  
 824 localization techniques. The results of the cycling DA experiments follow a similar pat-  
 825 tern as the simpler tests with “static” covariances of the previous section.

- 826 1. NICE, Ad.-PLC, ASCE and the tuned PLC lead to nearly identical errors, and  
 827 the adaptive localization (Ad.-Loc) comes fairly close to the tuned localization (Loc),  
 828 reiterating that the discrepancy principle is robustly applicable to adaptive covari-  
 829 ance estimation.
- 830 2. PANIC reduces the error as compared to NICE, because the assumption of zero  
 831 (or near-zero) correlations in the far-field is valid for L’96. The additional error  
 832 reduction that PANIC achieves over NICE, however, is minor (as in the previous,  
 833 non-cycling examples).
- 834 3. Localization leads to smaller errors than NICE, PANIC, Ad.-PLC or PLC, but  
 835 the errors are still larger than what can be achieved with a large ensemble size or  
 836 a nearly optimal localization (POLO).
- 837 4. POLO based on correlations extracted from a large ensemble leads to smaller er-  
 838 rors than all other techniques, but still cannot reach the low error achieved by a  
 839 large ensemble size. This could be due to the Gaussian assumption underpinning  
 840 POLO, which is not satisfied in cycling DA experiments with L’96, or it could in-  
 841 dicate more general limitations of correlation-based covariance corrections.
- 842 5. We encountered no negative eigenvalues during the cycling DA experiments with  
 843 PLC or Ad.-PLC.



844 We further note that all covariance estimation methods (NICE, PANIC, Ad.-PLC, PLC,  
 845 ASCE, Ad.-Loc, Loc, POLO) lead to much smaller errors than a vanilla EnKF without  
 846 inflation or localization/covariance estimation. The vanilla EnKF diverges and, there-  
 847 fore, leads to macroscopic error.

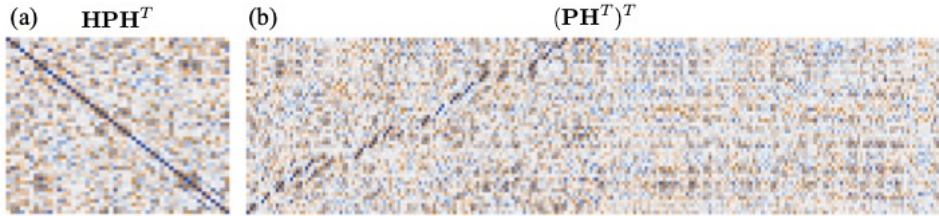
848 Our test with the L’96 model re-iterates our conclusions based on the simpler ex-  
 849 periments of the previous section: NICE, Ad.-PLC and ASCE reduce the error in co-  
 850 variance estimates and, therefore, in a cycling EnKF *without* making any assumptions  
 851 about the underlying correlation structure and *without* tuning (only inflation is tuned  
 852 for these methods). The fact that localization leads to smaller errors than NICE, Ad.-  
 853 PLC or ASCE stems from the heavy tuning and, perhaps more importantly, from the  
 854 fact that the underlying correlation structure here is consistent with the assumptions of  
 855 classical localization.

### 856 4.3 Cycling Data Assimilation Experiments with a Geomagnetic Proxy 857 Model

858 We consider cycling DA experiments with an EnKF on a proxy model for geomag-  
 859 netic data assimilation, described in detail by Gwirtz et al. (2021). The model consists  
 860 of a (chaotic) Kuramoto-Sivashinsky (KS) equation coupled to an induction equation,  
 861 and describes the spatial and temporal variations of a velocity field coupled, via induc-  
 862 tion, to a magnetic field. We consider the model in a 2D configuration on a square and  
 863 discretize the partial differential equations (PDE) by a spectral method (Fourier series),  
 864 which leads to a state dimension of  $n = 1920$  Fourier coefficients. Following Gwirtz et  
 865 al. (2021), we collect observations of Fourier modes of the magnetic field with wavenum-  
 866 bers in the  $x$ - and  $y$ -directions that are less than or equal to three (for a sum total of  
 867 48 Fourier coefficients). The time interval between two consecutive observations is about  
 868 7% of the model’s  $e$ -folding time. Note that the velocity field is entirely unobserved. This  
 869 setup is somewhat indicative of what to expect in a larger numerical dynamo model for  
 870 decadal-scale forecasts of the geomagnetic field (Gwirtz et al., 2021).

871 We assimilate the spectral observations using a stochastic EnKF with ensemble size  
 872  $n_e = 100$ , essentially repeating the DA experiment reported in Section 4.2 of Gwirtz  
 873 et al. (2021). Since we observe Fourier coefficients, we have no natural notion of a “spa-  
 874 tial” distance, and we therefore resort to NICE and Ad.-PLC to correct the covariances  
 875 within the EnKF. We have tried hard, but failed to find a localization based on a spa-  
 876 tial decay of correlation that reduces errors, see also Gwirtz et al. (2021). Note that the  
 877 state dimension is large ( $n_x = 1920$ ), but the number of observations is small ( $n_y =$   
 878 48), so that it is natural to estimate the matrices  $\mathbf{H}\hat{\mathbf{P}}\mathbf{H}^T$  ( $48 \times 48$ ) and  $\hat{\mathbf{P}}\mathbf{H}^T$  ( $1920 \times$   
 879 48), rather than the ensemble covariance  $\hat{\mathbf{P}}$  ( $1920 \times 1920$ ). The results reported below,  
 880 however, do not change much if we estimate the ensemble covariance  $\hat{\mathbf{P}}$  using the same  
 881 methods. A more detailed discussion of the differences between these two approaches in  
 882 the context of localization can be found in Campbell et al. (2010).

883 The apparent absence of correlation structure in covariance matrices within an EnKF  
 884 is described in detail in Gwirtz et al. (2021) (see, e.g., Figures 5a, 5b and 10 of Gwirtz  
 885 et al. (2021)), and is also illustrated in Figure 7, where we plot correlation matrices as-  
 886 sociated with  $\mathbf{P}\mathbf{H}^T$  and  $\mathbf{H}\mathbf{P}\mathbf{H}^T$  during one cycle of an EnKF with a large ensemble size  
 887  $n_e = 1000$ . It is clear from the figure that correlations are strong throughout the sys-  
 888 tem, but also that there is no coherent pattern. The lack of discernible correlation pat-  
 889 terns makes estimating covariances from a small ensemble difficult, but, as we will see,  
 890 NICE and Ad.-PLC handle this problem well. Moreover, the correlations change from  
 891 one DA cycle to the next (see Fig. 10 in Gwirtz et al. (2021)), but since NICE and Ad.-  
 892 PLC are adaptive, these methods can capture the time-varying correlation structure within  
 893 this cycling EnKF.

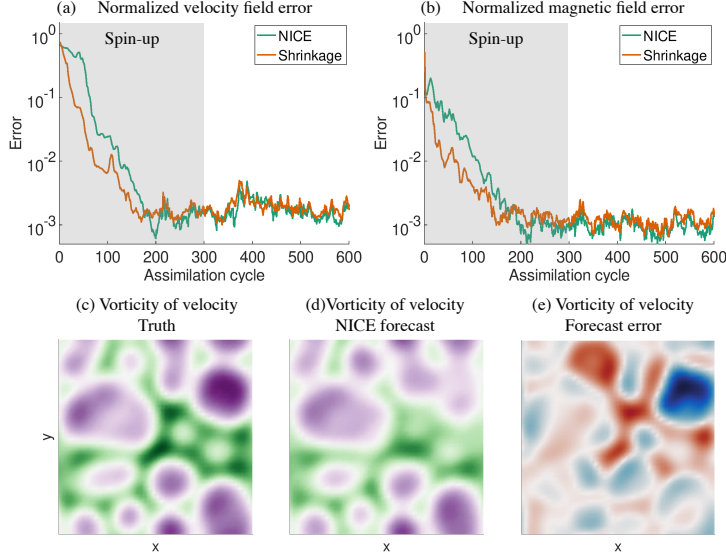


**Figure 7.** Correlations in a cycling EnKF for the geomagnetic model (large ensemble size). (a) The  $48 \times 48$  correlation matrix associated with  $\mathbf{HPH}^T$  during one cycle (post spin-up) of an EnKF with ensemble size  $n_e = 1,000$ . (b) The correlation matrix associated with  $\mathbf{PH}^T$ , during one cycle (post spin-up) of an EnKF with ensemble size  $n_e = 1,000$ . The correlation matrix associated with  $\mathbf{PH}^T$  is truncated at wave number five so that its size is  $240 \times 48$ . To make the figure easier to read, we transpose the correlation matrix associated with  $\mathbf{PH}^T$ . What is important to note from this figure is that (i) there are strong correlations across various variables represented in  $\mathbf{HPH}^T$  and  $\mathbf{PH}^T$ ; and (ii) the correlations follow no discernible pattern and, what’s worse, large and small correlations switch places across various assimilation cycles (not shown, but see Gwirtz et al. (2021)). Color indicates the matrix elements with blue corresponding to one, white to zero, and red to minus one.

894 For our numerical tests, we set the ensemble size to  $n_e = 100$  and use NICE and  
 895 Ad.-PLC, along with a 6% covariance inflation of both  $\mathbf{HPH}^T$  and  $\mathbf{PH}^T$  (Gwirtz et al.,  
 896 2021). For each EnKF, we perform 600 DA cycles, with the first 300 cycles being dis-  
 897 carded as “spin-up.” Figure 8 illustrates the results of our numerical experiments for NICE  
 898 (results for Ad.-PLC are similar). Panels (a) and (b) show errors (truth minus a one-  
 899 cycle forecast) as a function of the DA cycle for the velocity field and magnetic field for  
 900 EnKFs with NICE (green). We note the spin-up period and the subsequent stable DA  
 901 phase. The errors in the figure are normalized by the macroscopic error, which is the er-  
 902 ror one would expect without any data assimilation. Panels (c) - (e) illustrate a forecast  
 903 based on an EnKF using NICE for covariance estimation. Shown is the vorticity of the  
 904 velocity field approximately 4.7  $e$ -folding times *after* the last assimilation cycle (panel (c)),  
 905 along with the NICE-EnKF forecast (panel (d)) and the difference of the two (panel (e)).  
 906 It is notable that the EnKF with a NICE covariance estimation can be used to create  
 907 forecasts that are accurate on practically relevant time scales.

908 We compare the performance of an EnKF with a small ensemble size ( $n_e = 100$ )  
 909 using NICE and Ad.-PLC, to an EnKF with a large ensemble ( $n_e = 1,000$ ) but *with-*  
 910 *out* covariance corrections. In this context, it is important to note that Gwirtz et al. (2021)  
 911 showed that an EnKF without covariance corrections stabilized on this problem with an  
 912 ensemble size of  $n_e = 800$ . We further consider an EnKF with  $n_e = 100$ , and with co-  
 913 variance estimation based on a shrinkage estimator, which decreases the magnitude of  
 914 all off-diagonal elements of a covariance matrix. The shrinkage estimator is taken from  
 915 Gwirtz et al. (2021), where it was heavily tuned, and was found to be “the best” covari-  
 916 ance estimation method for this problem.

917 The results of our numerical experiments and relevant results reported in Gwirtz  
 918 et al. (2021) are summarized in Table 2, which lists errors in magnetic (observed) and  
 919 velocity (unobserved) fields. We note a similar pattern as in our earlier experiments: NICE  
 920 and Ad.-PLC are as good or better than a finely tuned estimator (Shrinkage) and the  
 921 adaptive covariance estimation methods indeed come quite close to the performance of  
 922 an EnKF with a much larger ensemble size. Moreover, both methods succeed in prop-  
 923 agating information from the observed magnetic field to the unobserved velocity field,



**Figure 8.** Illustration of covariance estimation within a cycling EnKF for a geomagnetic proxy model. (a) Normalized error in the unobserved velocity field as a function of assimilation cycle for two covariance estimation methods (NICE and shrinkage). (b) Normalized error in the partially observed magnetic field as a function of assimilation cycle for two covariance estimation methods (NICE and shrinkage). (c) Vorticity of the velocity field. (d) Forecast of the vorticity of the velocity field based on data assimilation with NICE. (e) Forecast error (difference between panels (c) and (d)).

	Error in mag. field	Error in vel. field
Shrinkage (tuned)	1.2	2.0
Ad.-PLC	1.2	2.5
NICE	1.0	1.8
Large ens.	0.7	1.1

**Table 2.** Normalized errors scaled by the respective macroscopic errors and multiplied by 10<sup>3</sup> for three covariance estimation methods and for an EnKF with a large ensemble applied to a geomagnetic proxy model.

924 as indicated by the small errors in the unobserved velocity field. In this example, NICE  
 925 leads to smaller errors than Ad.-PLC (in both fields). Nonetheless, the fact that both  
 926 adaptive methods succeed with essentially no tuning on a problem that is much harder  
 927 and much more high-dimensional than the previous test problems is reassuring and speaks  
 928 to the robustness of the proposed techniques. Moreover, NICE leads to smaller errors  
 929 than the best method thus far reported in the literature (the heavily tuned shrinkage es-  
 930 timator of Gwirtz et al. (2021)).

931 Finally, we report (again) that even though Ad.-PLC does not guarantee positive  
 932 semi-definite covariance estimates, all covariances ( $\mathbf{HPH}^T$ ) that were estimated with this  
 933 method in this example turned out to be positive semi-definite.

#### 934 4.4 Inversion of Electromagnetic Data

935 We now apply ensemble Kalman inversion (EKI) to a marine electromagnetic (EM)  
 936 inverse problem. The goal of the inversion is to compute resistivity as a function of depth  
 937 from measurements of apparent resistivity and phase, both as a function of period. The  
 938 seafloor magnetotelluric (MT) data (ten apparent resistivities along with ten phases, see  
 939 Figure 10(d)) are collected off-shore of New Jersey (Gustafson et al., 2019; Blatter et al.,  
 940 2019). The data are equipped with error estimates in the form of standard deviations.  
 941 The MT model uses a standard recursion relationship (Ward and Hohmann (2012), see  
 942 also Blatter et al. (2022b, 2022a)), and is discretized with 60 layers, each 20m thick,

943 As is common in geophysical inversion, we use a quadratic regularization, i.e., we  
 944 minimize the cost function

$$945 F(\mathbf{x}) = \left\| \mathbf{R}_d^{-\frac{1}{2}} (\mathbf{d} - \mathcal{M}(\mathbf{x})) \right\|_2^2 + \mu \left\| \mathbf{B}^{-\frac{1}{2}} \mathbf{x} \right\|_2^2, \quad (50)$$

946 where  $\mathbf{d}$  are the data,  $\mathbf{x}$  are the unknown resistivities,  $\mathcal{M}$  is the MT model,  $\mathbf{R}_d$  is a di-  
 947 agonal matrix that contains the variances associated with the data on its diagonal, and  
 948 where  $\mathbf{B}$  is a regularization matrix, which we chose to be a covariance matrix with a Gaus-  
 949 sian kernel and length scale  $\ell = 200\text{m}$ . The regularization parameter  $\mu$  was obtained  
 950 via an Occam inversion (Constable et al., 1987). We note that discovering an appropri-  
 951 ate regularization strength  $\mu$  in EM inversions is an interesting subject in itself, but for  
 952 the purposes of this numerical demonstration, it is sufficient to think of  $\mu$  as being given.  
 953 A similar EM inverse problem was considered by Tong and Morzfeld (2023), also in the  
 954 context of localizing EKI.

955 To apply EKI to this regularized problem, we recast the cost function as

$$956 f(x) = \left\| \mathbf{R}^{-\frac{1}{2}} (\mathbf{y} - \mathcal{G}(\mathbf{x})) \right\|_2^2, \quad (51)$$

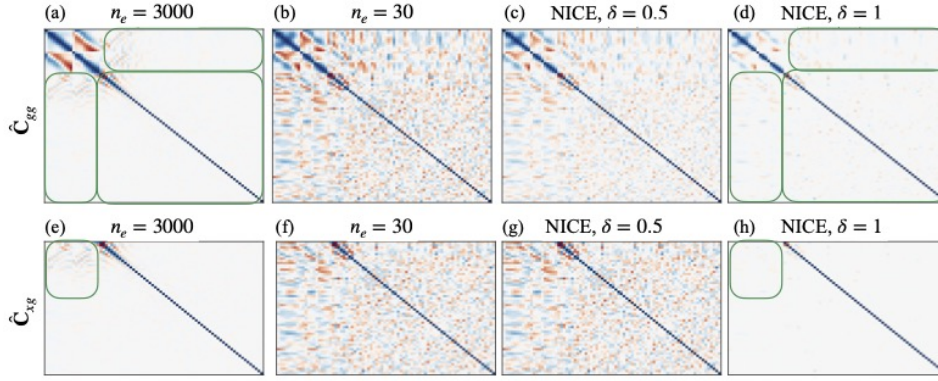
957 where

$$958 \mathbf{R}^{-\frac{1}{2}} = \begin{pmatrix} \mathbf{R}_d^{-\frac{1}{2}} & \mathbf{0} \\ \mathbf{0} & \sqrt{\mu} \mathbf{B}^{-\frac{1}{2}} \end{pmatrix}, \quad \mathbf{y} = \begin{pmatrix} \mathbf{d} \\ \mathbf{0} \end{pmatrix}, \quad \mathcal{G}(\mathbf{x}) = \begin{pmatrix} \mathcal{M}(\mathbf{x}) \\ \mathbf{x} \end{pmatrix}. \quad (52)$$

959 This “trick” is explained in detail in Chada et al. (2020) where the resulting method is  
 960 called Tikhonov regularized ensemble Kalman inversion (TEKI). Note that the EKI frame-  
 961 work (see Section 2.2) can now be directly applied, but at the expense that the “data-  
 962 data” correlations in  $\hat{\mathbf{C}}_{gg}$  are stacks of an ensemble of model outputs and the ensemble  
 963 itself.

964 Recall that the EKI iteration requires that we repeatedly estimate the covariances  
 965  $\hat{\mathbf{C}}_{gg}$  and  $\hat{\mathbf{C}}_{xg}$  from the ensemble. We correct these covariances using NICE, Ad-PLC and  
 966 ASCE. For all three methods, our numerical experiments indicate that the tunable pa-  
 967 rameter  $\delta$  in the discrepancy principle needs to be decreased when we correct the data-  
 968 to-unknown covariances  $\hat{\mathbf{C}}_{xg}$ . A factor of  $\delta = 0.5$  leads to good results, whereas  $\delta =$   
 969  $1$  leads to TEKI iterations that do not reduce the error as low as with  $\delta = 0.5$ . The rea-  
 970 son for reducing  $\delta$  is that a smaller  $\delta$  leads to a softer correction, which is needed because  
 971 several of the “true” data-to-unknown covariances are small, and it is advantageous to  
 972 keep them, rather than to remove them, in order to propagate information from the data  
 973 to the unknown variables. This effect is illustrated in Figure 9: NICE with a “strong”  
 974 correction ( $\delta = 1$ ) is adequate for the data-data correlations (top row), but inadequate  
 975 for the cross correlations (bottom row).

976 Implementing a spatial localization is neither intuitive nor easy in this example,  
 977 but we tried it nonetheless. First, we apply a localization to  $\hat{\mathbf{C}}_{gg}$ , although this has lit-  
 978 tle physical motivation. We chose a localization matrix with a Gaussian kernel and a length  
 979 scale  $\ell = 200\text{m}$  after some initial tries (no careful tuning). Performing a localization  
 980 on  $\hat{\mathbf{C}}_{xg}$  ( $60 \times 80$ ) is more tricky. We apply *no* localization to the first 20 columns of  $\hat{\mathbf{C}}_{xg}$ ,



**Figure 9.** Correlation matrices corresponding to  $\mathbf{C}_{gg}$  (top row) and  $\mathbf{C}_{xg}$  (bottom row) during one step of a TEKI. Panels (a) and (e) show estimates of the correlations for a large ensemble size. Panels (b) and (f) show estimates of the correlations for a small ensemble size. Panels (c) and (g) show the NICE estimator with  $\delta = 0.5$ . Panels (d) and (h) show the NICE estimator with  $\delta = 1$ . In panel (a), green squares highlight areas in which correlations are weak, which NICE with  $\delta = 1$  (panel (d)) dampens, but NICE with  $\delta = 0.5$  keeps. In panel (e), a green square highlights an area in which correlations are present, but which are dampened too strongly by NICE with  $\delta = 1$ , as in panel (h), whereas NICE with  $\delta = 0.5$  “keeps” these correlations. Color indicates the matrix elements with blue corresponding to one, white to zero, and red to minus one.

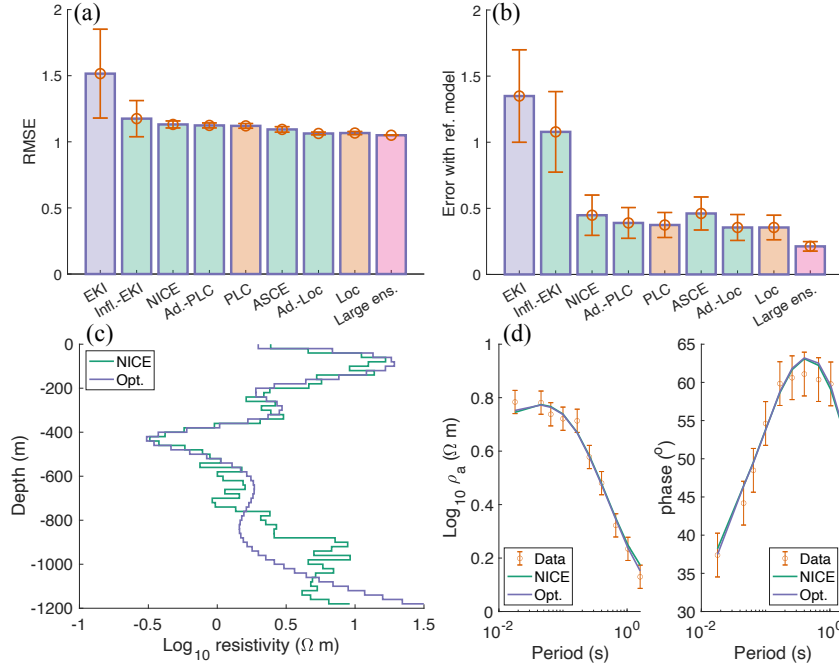
981 which corresponds to the covariances computed from the ensemble of model outputs. We  
 982 apply a Gaussian localization with length scale  $\ell = 200\text{m}$  to the remaining 60 columns.  
 983 With this same setup, we can also apply an adaptive localization (Ad.-Loc).

984 We now run TEKIs with various covariance estimation schemes and covariance in-  
 985 flation (see equation (49)). The inflation depends on the root mean square error (RMSE),  
 986 defined by

$$987 \quad \text{RMSE} = \sqrt{\frac{1}{n_d} \sum_{i=1}^{n_d} \left( \frac{\mathbf{d}_i - \hat{\mathbf{d}}_i}{\sigma_i} \right)^2}, \quad (53)$$

988 where  $\mathbf{d}_i$ ,  $i = 1, \dots, n_d$ , are the  $n_d$  data points,  $\sigma_i$  are the corresponding observation  
 989 error standard deviations (given as part of the MT data set as the diagonal elements of  
 990  $\mathbf{R}_d^{1/2}$ );  $\hat{\mathbf{d}} = \mathcal{M}(\hat{\mathbf{x}})$  are model predictions based on the mean of the TEKI ensemble,  $\hat{\mathbf{x}}$ .  
 991 The inflation is  $\kappa = 15\%$  when  $\text{RMSE} > 1.2$ ,  $\kappa = 10\%$  when  $1.1 \leq \text{RMSE} \leq 1.2$ ,  
 992 and we turn the inflation off ( $\kappa = 0$ ) when  $\text{RMSE} < 1.1$ . We did not tune the infla-  
 993 tion systematically.

994 We use TEKIs with ensemble size  $n_e = 30$  and 200 iterations. For each TEKI,  
 995 we perform 100 independent experiments, each with a different random initial ensem-  
 996 ble and then average the results. Our findings are summarized in Figure 10. Panel (a)  
 997 shows the averaged RMSE for each TEKI. Note that an RMSE of approximately one is  
 998 good because then the TEKI estimate fits the data to within the assumed error level (stan-  
 999 dard deviation of the data). First, we note as before that all covariance estimation meth-  
 1000 ods (NICE, Ad.-PLC, PLC, ASCE, Ad.-Loc, Loc) lead to TEKIs which can achieve an  
 1001 acceptably low RMSE and that the adaptive methods are nearly as good as the tuned  
 1002 methods or a TEKI with a larger ensemble ( $n_e = 200$ ). Second, we note that the in-  
 1003 flation already has a large effect on the RMSE: An inflated TEKI reaches an RMSE that  
 1004 is lower than a “vanilla” TEKI without any covariance estimation or inflation.



**Figure 10.** Summary of results for the electromagnetic inversion. (a) RMSE (see (53)) of various TEKI implementations. (b) Error with respect to a reference model, obtained via gradient-based optimization (see (54)) of various TEKI implementations. In panels (a) and (b), the bars are averages over results obtained by randomizing the initial TEKI ensemble and the error bars denote one standard deviation. The bars are color coded so that blue labels a “vanilla” TEKI, green labels a TEKI with an adaptive covariance estimation, orange labels a TEKI with a tuned covariance estimation, and pink (furthest to the right) labels a large ensemble result. Panel (c) shows (log) resistivity as a function of depth obtained via Gauss-Newton optimization (blue) and TEKI with NICE (green). Panel (d) shows the EM data and the model output resulting from TEKI with NICE (green) and Gauss-Newton method (purple). Averages and standard deviations are computed from 100 independent numerical experiments.

1005  
1006  
1007

We further assess the “quality” of our TEKI inversions by comparing the TEKI results to a gradient-based optimization (Gauss-Newton). We measure the difference between the TEKI result and the Gauss-Newton result by the error

1008

$$\text{Ref. Error} = \frac{\|\text{res.GN} - \text{res.teki}\|_2^2}{\|\text{res.GN}\|_2^2} \quad (54)$$

1009  
1010  
1011  
1012  
1013  
1014  
1015  
1016  
1017  
1018  
1019

where res. refers to the (log) resistivity and the subscript GN refers to the Gauss-Newton method and subscript teki refers to a TEKI result. The error with respect to a reference model is shown in Figure 10(b). We see that the reference error behaves very similarly to the RMSE (not surprisingly): The covariance estimation methods all lead to a small reference error and all methods perform similarly. The inflated TEKI (no additional covariance estimation) leads to a significantly larger reference error than the other TEKIs, although the RMSE is comparable. The two errors can be different here because many different models fit the MT data similarly well. The large reference error indicates that the model obtained with an inflated TEKI is quite different from the reference model. Thus, the covariance estimation is helpful here to “smooth” the models so that they are similar to the reference model, obtained by Gauss-Newton (see also Figure 10(c)).

1020 Finally, Figures 10(c) and (d) show a typical result obtained with TEKI and NICE.  
 1021 Panel (c) shows the (log) resistivity as a function of depth and panel (d) shows the as-  
 1022 sociated fit to the data. For comparison, we also show the resistivity and data fit we ob-  
 1023 tain via Gauss-Newton. The TEKI approximation with NICE is very similar to the Gauss-  
 1024 Newton result for depths up to about 600m, where the data are most informative (small  
 1025 error with respect to the reference model in Figure 10(b)) and the fit to the data for TEKI  
 1026 and Gauss-Newton is nearly identical (small RMSE in Figure 10(a)).

#### 1027 4.5 Training Feed-Forward Neural Networks with Time-Averaged Data

1028 Our last example is a simplification of a climate sciences problem in which sub-grid  
 1029 parameterizations of a climate model are represented by neural networks (NN). The train-  
 1030 ing strategy for the neural network is to define a loss function in terms of time-averaged  
 1031 data of the climate model and to adjust the weights and biases of the NN to minimize  
 1032 the loss function. The usual back propagation (gradient descent) cannot be used in this  
 1033 context because the “map” from the NN weights and biases to the time-averaged data  
 1034 of the climate model may not be differentiable, or derivatives may be difficult to obtain  
 1035 (Schneider et al., 2024).

1036 As a “cartoon” for this difficult problem, we consider a modified Lorenz model (mL’96)  
 1037 as a stand-in for a climate model and we parameterize the forcing of mL’96 by a sim-  
 1038 ple feed-forward neural network. Specifically, the mL’96 model is

$$1039 \frac{dx_i}{dt} = (x_{i+1} - x_{i-2})x_{i-1} - x_i + F_i, \quad (55)$$

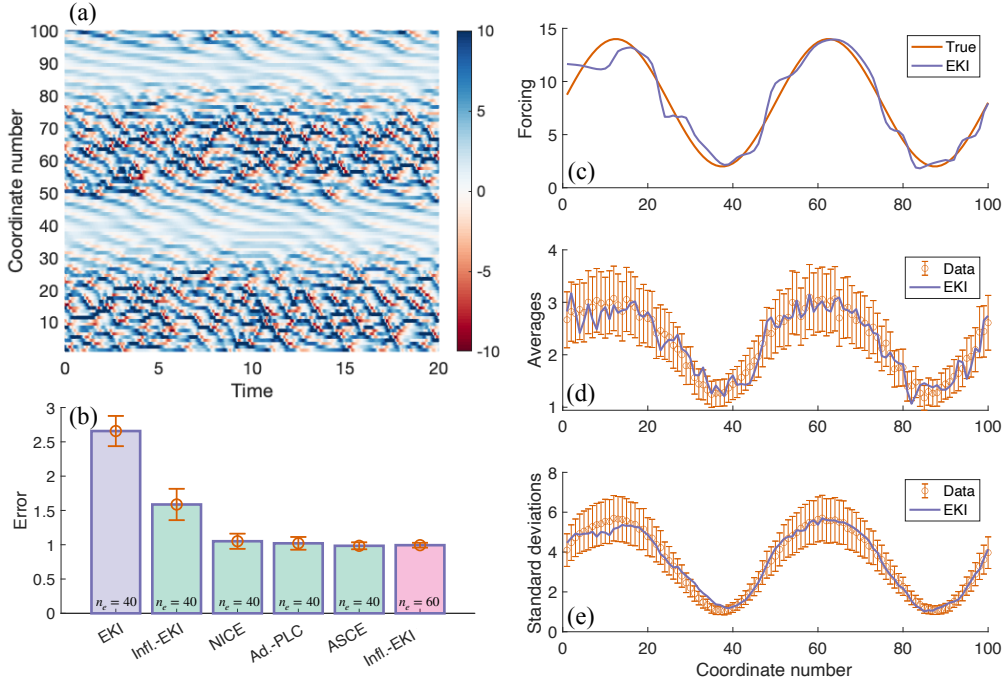
1040 where  $x_{-1} = x_{n_x-1}$ ,  $x_0 = x_{n_x}$ ,  $x_{n_x+1} = x_1$  (periodicity) and where

$$1041 F_i = 8 + 6 \sin\left(\frac{4\pi}{n_x} i\right), \quad (56)$$

1042 is a coordinate-dependent forcing. Note that the forcing is the only modification we make,  
 1043 and our modification is inspired by the storm-track model of Bishop et al. (2017). We  
 1044 choose the state dimension to be  $n_x = 100$ . Figure 11(a) shows a Hovmöller diagram  
 1045 of the mL’96 model and illustrates the time evolution of all  $n_x = 100$  coordinates as  
 1046 a function of time. Due to the sinusoidal forcing, we can identify regions of chaotic dy-  
 1047 namics (larger  $F_i$ ) and regions with more predictable characteristics (smaller  $F_i$ ).

1048 Our goal is to recover the forcing  $F_i$ ,  $i = 1, \dots, n_x$  from time-averaged data, which  
 1049 are the means and standard deviations of all  $n_x$  coordinates over a period of  $T = 500$   
 1050 time units ( $2n_x = 200$  data points). The noise in the data are independent mean-zero  
 1051 Gaussians with standard deviations equal to 10% of that of the data points. The neu-  
 1052 ral network that parameterizes the forcing is a feed-forward neural net with one input  
 1053 layer, one hidden layer and one output layer (Goodfellow et al., 2016). The total num-  
 1054 ber of weights and biases in the network is 91, largely due to the size of the hidden layer,  
 1055 which we adjusted so that the neural network is expressive enough to capture the sinu-  
 1056 soidal forcing.

1057 EKI requires an initial ensemble which we generate using ideas from transfer learn-  
 1058 ing. We draw  $n_e$  realizations of a smooth Gaussian process (Gaussian kernel, length scale  
 1059 is  $\ell = 5$ , (Rasmussen & Williams, 2005)) and then train a NN on each random func-  
 1060 tion draw. Here, we use back-propagation, as is standard in simple function approxima-  
 1061 tion tasks, because the NN is differentiable – the time-averaged data are not. The weights  
 1062 and biases of the NNs we obtain from training on random smooth functions represent  
 1063 the initial ensemble for our EKIs. This simple strategy works well for small ensembles  
 1064 (up to  $n_e = 60$ ), but it leads to instabilities with EKIs with larger ensemble sizes. More  
 1065 sophisticated initialization may make it possible to run EKI with large  $n_e$  on this prob-  
 1066 lem, but since our focus is on EKI and small ensemble sizes, we do not pursue initial-  
 1067 ization of NNs in EKI further.



**Figure 11.** (a) Hovmöller diagram of the mL'96 model, showing the time evolution of all  $n_x$  coordinates as a function of time. (b) Average RMSE (bars) and standard deviations (error bars) of several EKI variants, computed over ten independent experiments, each using a different set of perturbations within the various EKIs. (c)-(e) Results of a typical EKI inversion with NICE covariance estimation. (c) Recovered forcing, parameterized by an NN, trained with EKI (purple) and true forcing (orange). (d) Averages of the  $n_x = 100$  mL'96 coordinates (error bars) and EKI-NN reconstructions (purple). (e) Standard deviations of the  $n_x = 100$  mL'96 coordinates (error bars) and EKI-NN reconstructions (purple).

1068  
1069  
1070  
1071  
1072

A typical result we obtain with EKI and NICE is illustrated in Figure 11(c)-(e), which shows the recovered forcing (panel (c)) and data fits (panels (d) and (e)). The EKI can train the NN so that the mL'96 model with the NN parameterization fits the data to within the assumed errors. Moreover, the recovered NN captures the sinusoidal variation of the forcing.

1073  
1074  
1075  
1076  
1077  
1078

We now follow our usual procedure and compare EKIs of various flavors: (i) EKI with NICE; (ii) EKI with ASCE; and (iii) EKI with Ad.-PLC. All EKIs apply covariance estimation to  $\hat{\mathbf{C}}_{gg}$  and  $\hat{\mathbf{C}}_{xg}$ , and we again adjust the tuning factor  $\delta$  to be equal to 0.5 when estimating  $\mathbf{C}_{xg}$ . The EKIs with NICE, ASCE or Ad.-PLC further inflate the covariance matrices with the same strategy as described in Section 4.4. We compare the above EKIs to a vanilla EKI, as well as to an EKI with inflation.

1079  
1080  
1081  
1082  
1083  
1084  
1085  
1086

The results of our comparison are illustrated in Figure 11(b), which shows the average RMSE of the various EKIs, computed over ten independent experiments, each using different random perturbations during 30 iterations. The EKIs with NICE, ASCE or Ad.-PLC use an ensemble size  $n_e = 40$  and we compare their performance to EKIs with or without inflation and ensemble sizes 40 and 60. The results we obtain in this example are in line with our earlier findings: NICE, ASCE and Ad.-PLC perform very similarly, reduce the error compared to inflated or vanilla EKIs and lead to a good fit to the data. Moreover, NICE, ASCE or Ad.-PLC result in similar errors as an inflated EKI with



1087 a larger ensemble size ( $n_e = 60$ ). In summary, we can apply EKI to train a neural net-  
 1088 work that parameterizes a chaotic dynamical system, and covariance estimation meth-  
 1089 ods such as NICE, ASCE or Ad.-PLC help with the computational efficiency of the in-  
 1090 version because they enable us to run the EKI with a small ensemble size.

## 1091 5 Summary and Conclusions

1092 We consider the problem of estimating a covariance matrix from a small number  
 1093 of samples in the context of Earth science applications. Our focus is on problems in which  
 1094 the correlation structure is *unknown*, because the problem of high-dimensional covari-  
 1095 ance estimation with a priori assumptions about the correlation structure is essentially  
 1096 solved (i.e. covariance localization in numerical weather prediction).

1097 A new method for covariance estimation, called NICE (noise-informed covariance  
 1098 estimation), is built on a single fundamental fact we know about estimating correlations:  
 1099 Small correlations are notoriously hard to compute, while it is relatively easy to com-  
 1100 pute large correlations. We translate this simple idea into an efficient and adaptive co-  
 1101 variance estimation method that guarantees a symmetric positive semi-definite covari-  
 1102 ance estimate.

1103 Adaptivity of NICE is achieved by (i) estimating a noise level for the correlation  
 1104 matrix; and (ii) adjusting the correlation corrections so that the resulting correlation es-  
 1105 timate is compatible with the noise level. We also used these ideas to design a few other  
 1106 adaptive covariance estimation methods: adaptive power law corrections (Ad.-PLC), adap-  
 1107 tive localization (Ad.-Loc), adaptive soft-thresholding (Ad.-ST), and adaptive sparse co-  
 1108 variance estimation (ASCE).

1109 We compared our new covariance estimation methods to several other methods on  
 1110 a large set of numerical experiments with correlation structures that are not easy to an-  
 1111 ticipate or decipher. Our tests include cycling data assimilation with a geomagnetic proxy  
 1112 model, geophysical inversion of field data, and the training of a feed-forward neural net-  
 1113 work with time-averaged data from a chaotic dynamical system. *All* new covariance es-  
 1114 timation methods we created perform well on this diverse set of numerical tests and are  
 1115 similar in accuracy to related tuned methods, which speaks for the robustness of our ap-  
 1116 proach to adaptive covariance estimation. NICE, however, has the advantage of guar-  
 1117 anteeing a positive semi-definite covariance estimator at a low computational cost.

## 1118 Data Availability Statement

1119 The code and data used in this manuscript are available at (Vishny et al., 2024).

## 1120 Acknowledgments

1121 We thank Jeff Anderson of NCAR and an anonymous reviewer for their careful reviews  
 1122 that helped us improve the manuscript.

1123 DV was supported by the NSF through grant AGS-2202991. MM is supported by  
 1124 the US Office of Naval Research (ONR) grant N00014-21-1-2309. KG is supported by  
 1125 an appointment to the NASA Postdoctoral Program at Goddard Space Flight Center,  
 1126 administered by Oak Ridge Associated Universities under contract with NASA. EB is  
 1127 supported by the Foster and Coco Stanback Postdoctoral Fellowship. ORAD was sup-  
 1128 ported by the generosity of Eric and Wendy Schmidt by recommendation of the Schmidt  
 1129 Futures program, National Science Foundation Grant AGS-1835860, the Defense Advanced  
 1130 Research Projects Agency (Agreement No. HR00112290030), the Heising-Simons Foun-  
 1131 dation, Audi Environmental Foundation, and the Cisco Foundation. DH is supported  
 1132 by the US Office of Naval Research (ONR) grant N0001422WX00451.

We thank Dr. Dani Blatter of Las Positas College for help and advice with the mag-  
neto telluric (MT) data, the MT forward modeling code, and the geophysical inversion.

## References

- Agapiou, S., Papaspiliopoulos, O., Sanz-Alonso, D., & Stuart, A. M. (2017). Importance sampling: Intrinsic dimension and computational cost. *Statist. Sci.*, *32*(3), 405–431.
- Al Ghattas, O., & Sanz-Alonso, D. (2022). Non-asymptotic analysis of ensemble Kalman updates: Effective dimension and localization. *arXiv:2208.03246*.
- Anderson, J. L. (2001). An ensemble adjustment Kalman filter for data assimilation. *Monthly weather review*, *129*(12), 2884–2903.
- Anderson, J. L. (2012). Localization and sampling error correction in ensemble Kalman filter data assimilation. *Monthly Weather Review*, *140*(7), 2359–2371.
- Anderson, J. L., & Lei, L. (2013). Empirical localization of observation impact in ensemble Kalman filters. *Mon. Wea. Rev.*, *141*, 4140–4153.
- Anzengruber, S. W., & Ramlau, R. (2009, dec). Morozov’s discrepancy principle for Tikhonov-type functionals with nonlinear operators. *Inverse Problems*, *26*(2), 025001. doi: 10.1088/0266-5611/26/2/025001
- Bannister, R. N. (2017). A review of operational methods of variational and ensemble-variational data assimilation. *Quarterly Journal of the Royal Meteorological Society*, *143*(703), 607–633. doi: <https://doi.org/10.1002/qj.2982>
- Bickel, P., & Levina, E. (2008). Covariance regularization by thresholding. *The Annals of Statistics*, *36*(6), 2577 – 2604. doi: 10.1214/08-AOS600
- Bickel, P., & Lindner, M. (2012). Approximating the inverse of banded matrices by banded matrices with applications to probability and statistics. *Theory of Probability & Its Applications*, *56*(1), 1–20. doi: 10.1137/S0040585X97985224
- Bieli, M., Dunbar, O. R. A., de Jong, E. K., Jaruga, A., Schneider, T., & Bischoff, T. (2022). An efficient Bayesian approach to learning droplet collision kernels: Proof of concept using “Cloudy,” a new n-moment bulk microphysics scheme. *Journal of Advances in Modeling Earth Systems*, *14*(8), e2022MS002994.
- Bishop, C. H., & Hodyss, D. (2007). Flow-adaptive moderation of spurious ensemble correlations and its use in ensemble-based data assimilation. *Quarterly Journal of the Royal Meteorological Society*, *133*(629), 2029–2044.
- Bishop, C. H., & Hodyss, D. (2009a). Ensemble covariances adaptively localized with ECO-RAP. Part 1: Tests on simple error models. *Tellus A: Dynamic Meteorology and Oceanography*, *61*(1), 84–96. doi: 10.1111/j.1600-0870.2007.00371.x
- Bishop, C. H., & Hodyss, D. (2009b). Ensemble covariances adaptively localized with ECO-RAP. Part 2: A strategy for the atmosphere. *Tellus A: Dynamic Meteorology and Oceanography*, *61*(1), 97–111. doi: 10.1111/j.1600-0870.2007.00372.x
- Bishop, C. H., & Hodyss, D. (2011). Adaptive ensemble covariance localization in ensemble 4D-VAR state estimation. *Monthly Weather Review*, *139*(4), 1241 – 1255. doi: 10.1175/2010MWR3403.1
- Bishop, C. H., Whitaker, J. S., & Lei, L. (2017). Gain form of the ensemble transform Kalman filter and its relevance to satellite data assimilation with model space ensemble covariance localization. *Monthly Weather Review*, *145*(11), 4575 – 4592.
- Blatter, D., Key, K., Ray, A., Gustafson, C., & Evans, R. (2019). Bayesian joint inversion of controlled source electromagnetic and magnetotelluric data to image freshwater aquifer offshore New Jersey. *Geophysical Journal International*, *218*(3), 1822–1837. doi: 10.1093/gji/ggz253
- Blatter, D., Morzfeld, M., Key, K., & Constable, S. (2022a, 06). Uncertainty quantification for regularized inversion of electromagnetic geophysical data – Part

- 1186 II: Application in 1-D and 2-D problems. *Geophysical Journal International*,  
 1187 *231*(2), 1075-1095. doi: 10.1093/gji/ggac242
- 1188 Blatter, D., Morzfeld, M., Key, K., & Constable, S. (2022b, 06). Uncertainty quan-  
 1189 tification for regularized inversion of electromagnetic geophysical data-Part I:  
 1190 Motivation and theory. *Geophysical Journal International*, *231*(2), 1057-1074.  
 1191 doi: 10.1093/gji/ggac241
- 1192 Bocquet, M. (2016). Localization and the iterative ensemble Kalman smoother.  
 1193 *Quarterly Journal of the Royal Meteorological Society*, *142*(695), 1075-1089.
- 1194 Bocquet, M., & Sakov, P. (2014). An iterative ensemble Kalman smoother. *Quar-*  
 1195 *terly Journal of the Royal Meteorological Society*, *140*(682), 1521-1535.
- 1196 Buehner, M. (2005). Ensemble-derived stationary and flow-dependent background-  
 1197 error covariances: Evaluation in a quasi-operational NWP setting. *Quarterly*  
 1198 *Journal of the Royal Meteorological Society*, *131*(607), 1013-1043. doi: [https://](https://doi.org/10.1256/qj.04.15)  
 1199 [doi.org/10.1256/qj.04.15](https://doi.org/10.1256/qj.04.15)
- 1200 Buehner, M. (2012). Evaluation of a spatial/spectral covariance localization ap-  
 1201 proach for atmospheric data assimilation. *Monthly Weather Review*, *140*(2),  
 1202 617–636. doi: 10.1175/MWR-D-10-05052.1
- 1203 Buehner, M., Mourneau, J., & Charette, C. (2013). Four-dimensional ensemble-  
 1204 variational data assimilation for global deterministic weather prediction. *Non-*  
 1205 *lin. Processes Geophys.*, *20*, 669-682.
- 1206 Buehner, M., & Shlyaeva, A. (2015). Scale-dependent background-error covariance  
 1207 localisation. *Tellus A: Dynamic Meteorology and Oceanography*, *67*(1). doi: 10  
 1208 .3402/tellusa.v67.28027
- 1209 Burgers, G., Leeuwen, P. V., & Evensen, G. (1998). Analysis scheme in the ensem-  
 1210 ble Kalman filter. *Monthly weather review*, *126*(6), 1719–1724.
- 1211 Campbell, W., Bishop, C., & Hodyss, D. (2010). Vertical covariance localization for  
 1212 satellite radiances in ensemble kalman filters. *Monthly Weather Review*, *138*,  
 1213 282-290.
- 1214 Chada, N. K., Chen, Y., & Sanz-Alonso, D. (2021). Iterative ensemble Kalman  
 1215 methods: A unified perspective with some new variants. *Foundations of Data*  
 1216 *Science*, *3*(3), 331-369. doi: 10.3934/fods.2021011
- 1217 Chada, N. K., Stuart, A. M., & Tong, X. T. (2020). Tikhonov regularization within  
 1218 ensemble Kalman inversion. *SIAM Journal on Numerical Analysis*, *58*(2),  
 1219 1263–1294.
- 1220 Chada, N. K., & Tong, X. T. (2022). Convergence acceleration of ensemble Kalman  
 1221 inversion in nonlinear settings. *Mathematics of Computation*, *91*(335), 1247–  
 1222 1280.
- 1223 Chen, Y., & Oliver, D. (2010). Cross-covariances and localization for EnKF in multi-  
 1224 phase flow data assimilation. *Computational Geosciences*, *14*(4), 579–601.
- 1225 Chen, Y., & Oliver, D. (2013). Levenberg-Marquardt forms of the iterative ensemble  
 1226 smoother for efficient history matching and uncertainty quantification. *Compu-*  
 1227 *tational Geosciences*, *17*, 689–703.
- 1228 Chen, Y., & Oliver, D. (2017). Localization and regularization for iterative ensemble  
 1229 smoothers. *Computational Geosciences*, *21*(1), 13–30.
- 1230 Chevrotière, M. D. L., & Harlim, J. (2017). A data-driven method for improving the  
 1231 correlation estimation in serial ensemble Kalman filters. *Monthly Weather Re-*  
 1232 *view*, *145*(3), 985-1001. doi: 10.1175/MWR-D-16-0109.1
- 1233 Chorin, A. J., & Morzfeld, M. (2013). Conditions for successful data assimilation.  
 1234 *Journal of Geophysical Research: Atmospheres*, *118*(20), 11,522-11,533. doi: 10  
 1235 .1002/2013JD019838
- 1236 Cleary, E., Garbuno-Inigo, A., Lan, S., Schneider, T., & Stuart, A. M. (2021). Cali-  
 1237 brate, emulate, sample. *Journal of Computational Physics*, *424*, 109716.
- 1238 Constable, S., Parker, R., & Constable, C. (1987). Occam’s inversion: A practical  
 1239 algorithm for generating smooth models from electromagnetic sounding data.  
 1240 *Geophysics*, *3*(52), 289–300.

- 1241 Dunbar, O. R. A., Lopez-Gomez, I., Garbuno-Iñigo, A., Huang, D. Z., Bach, E., &  
 1242 Wu, J.-L. (2022). EnsembleKalmanProcesses.jl: Derivative-free ensemble-  
 1243 based model calibration. *Journal of Open Source Software*, 7(80), 4869. doi:  
 1244 10.21105/joss.04869
- 1245 Emerick, A. A., & Reynolds, A. (2011). Combining sensitivities and prior infor-  
 1246 mation for covariance localization in the ensemble Kalman filter for petroleum  
 1247 reservoir applications. *Computational Geosciences*, 15(2), 251–269.
- 1248 Emerick, A. A., & Reynolds, A. C. (2013). Ensemble smoother with multiple data  
 1249 assimilation. *Computers & Geosciences*, 55, 3-15.
- 1250 Evensen, G. (1994). Sequential data assimilation with a nonlinear quasi-geostrophic  
 1251 model using Monte Carlo methods to forecast error statistics. *Journal of Geo-  
 1252 physical Research: Oceans*, 99(C5), 10143-10162. doi: 10.1029/94JC00572
- 1253 Evensen, G. (2009). *Data assimilation: the ensemble Kalman filter* (2nd ed.).  
 1254 Springer.
- 1255 Flowerdew, J. (2015). Towards a theory of optimal localisation. *Tellus A: Dynamic  
 1256 Meteorology and Oceanography*, 67(1), 25257. doi: 10.3402/tellusa.v67.25257
- 1257 Friedman, J., Hastie, T., & Tibshirani, R. (2007, 12). Sparse inverse covariance es-  
 1258 timation with the graphical lasso. *Biostatistics*, 9(3), 432-441. doi: 10.1093/  
 1259 biostatistics/kxm045
- 1260 Furrer, R., & Bengtsson, T. (2007). Estimation of high-dimensional prior and pos-  
 1261 terior covariance matrices in Kalman filter variants. *Journal of Multivariate  
 1262 Analysis*, 98(2), 227-255.
- 1263 Gaspari, G., & Cohn, S. (1999). Construction of correlation functions in two and  
 1264 three dimensions. *Quarterly Journal of the Royal Meteorological Society*, 125,  
 1265 723–757.
- 1266 Gharamti, M. E., Reader, K., & Anderson, J. L. (2019). Comparing adaptive prior  
 1267 and posterior inflation for ensemble filters using an atmospheric general circula-  
 1268 tion model. *Mon. Wea. Rev.*, 147, 2535-2553.
- 1269 Gilpin, S., Matsuo, T., & Cohn, S. E. (2023). A generalized, compactly supported  
 1270 correlation function for data assimilation applications. *Quarterly Journal of the  
 1271 Royal Meteorological Society*, 149(754), 1953-1989.
- 1272 Goodfellow, I., Bengio, Y., & Courville, A. (2016). *Deep learning*. MIT Press.
- 1273 Guillot, D., & Rajaratnam, B. (2015). Functions preserving positive definiteness for  
 1274 sparse matrices. *Transactions of the American Mathematical Society*, 367(1),  
 1275 627 - 649.
- 1276 Gustafson, C., Key, K., & Evans, R. (2019). Aquifer systems extending far offshore  
 1277 on the U.S. Atlantic margin. *Scientific Reports*, 9, 8709.
- 1278 Gwartz, K., Morzfeld, M., Kuang, W., & Tangborn, A. (2021). A testbed for ge-  
 1279 omagnetic data assimilation. *Geophysical Journal International*, 227, 2180-  
 1280 2203.
- 1281 Hamill, T. M., & Snyder, C. (2000). A hybrid ensemble Kalman filter–3D variational  
 1282 analysis scheme. *Monthly Weather Review*, 128(8), 2905 - 2919.
- 1283 Hamill, T. M., Whitaker, J. S., Anderson, J. L., & Snyder, C. (2009). Comments on  
 1284 “Sigma-Point Kalman Filter Data Assimilation Methods for Strongly Nonlin-  
 1285 ear Systems”. *Journal of the Atmospheric Sciences*, 66(11), 3498 - 3500.
- 1286 Harty, T., Morzfeld, M., & Snyder, C. (2021). Eigenvector-spatial localisa-  
 1287 tion. *Tellus A: Dynamic Meteorology and Oceanography*, 73(1), 1-18. doi:  
 1288 10.1080/16000870.2021.1903692
- 1289 Hodyss, D., Bishop, C. H., & Morzfeld, M. (2016). To what extent is your data  
 1290 assimilation scheme designed to find the posterior mean, the posterior mode or  
 1291 something else? *Tellus A*, 68, 1-17.
- 1292 Hodyss, D., Campbell, W. F., & Whitaker, J. S. (2016). Observation-dependent pos-  
 1293 terior inflation for the ensemble Kalman filter. *Monthly Weather Review*, 144,  
 1294 2667-2684.
- 1295 Hodyss, D., & Morzfeld, M. (2023). How sampling errors in covariance estimates

- 1296 cause bias in the Kalman gain and impact ensemble data assimilation. *Monthly*  
1297 *Weather Review*, 151(9), 2413 - 2426.
- 1298 Horn, R. A., & Johnson, C. R. (1991). *Topics in matrix analysis*. Cambridge Univer-  
1299 sity Press. doi: 10.1017/CBO9780511840371
- 1300 Houtekamer, P. L., & Mitchell, H. L. (1998). Data assimilation using an ensemble  
1301 Kalman filter technique. *Mon. Wea. Rev.*, 126, 796–811.
- 1302 Houtekamer, P. L., & Mitchell, H. L. (2001). A sequential ensemble Kalman filter  
1303 for atmospheric data assimilation. *Monthly Weather Review*, 129(1), 123–137.
- 1304 Huang, D., Huang, J., Reich, S., & Stuart, A. (2022). Efficient derivative-free  
1305 Bayesian inference for large-scale inverse problems. *Inverse Problems*, 38(12),  
1306 125006.
- 1307 Iglesias, M. A., Law, K. J., & Stuart, A. M. (2013). Ensemble Kalman methods for  
1308 inverse problems. *Inverse Problems*, 29(4), 045001.
- 1309 Khare, K., Oh, S.-Y., Rahman, S., & Rajaratnam, B. (2019). A scalable sparse  
1310 Cholesky based approach for learning high-dimensional covariance matrices in  
1311 ordered data. *Machine Learning*, 108, 2061 - 2086.
- 1312 Kovachki, N., & Stuart, A. (2019). Ensemble Kalman inversion: A derivative-free  
1313 technique for machine learning tasks. *Inverse Problems*, 35(9), 095005.
- 1314 Kuhl, D. D., Rosmond, T. E., Bishop, C. H., McLay, J., & Baker, N. L. (2013).  
1315 Comparison of hybrid ensemble/4DVar and 4DVar within the NAVDAS-AR  
1316 data assimilation framework. *Monthly Weather Review*, 141, 2740-2758.
- 1317 Lee, Y. (2021a).  $l_p$  regularization for ensemble Kalman inversion. *SIAM Journal on*  
1318 *Scientific Computing*, 43(5), A3417-A3437. doi: 10.1137/20M1365168
- 1319 Lee, Y. (2021b). *Sampling error correction in ensemble Kalman inversion*. arXiv.  
1320 doi: 10.48550/ARXIV.2105.11341
- 1321 Lorenc, A. (2003). The potential of the ensemble Kalman filter for NWP – a com-  
1322 parison with 4D-Var. *Quarterly Journal of the Royal Meteorological Society*,  
1323 129(595), 3183-3203. doi: 10.1256/qj.02.132
- 1324 Lorenc, A. (2017). Improving ensemble covariances in hybrid variational data assim-  
1325 ilation without increasing ensemble size. *Quarterly Journal of the Royal Meteoro-*  
1326 *logical Society*, 143(703), 1062-1072. doi: 10.1002/qj.2990
- 1327 Lorenz, E. (1996). Predictability: A problem partly solved. *In: Proceedings of the*  
1328 *ECMWF Seminar on predictability*, 1, 1-18.
- 1329 Luk, E., Bach, E., Baptista, R., & Stuart, A. (2024, June). *Learning Optimal Filters*  
1330 *Using Variational Inference* (No. arXiv:2406.18066). arXiv.
- 1331 Luo, X., Bhakta, T., & Nævdal, G. (2018). Correlation-based adaptive localization  
1332 with applications to ensemble-based 4D-seismic history matching. *SPE Jour-*  
1333 *nal*, 23(02), 396–427.
- 1334 Ménétrier, B., Montmerle, T., Michel, Y., & Berre, L. (2015). Linear filtering of  
1335 sample covariances for ensemble-based data assimilation. Part I: Optimal-  
1336 ity criteria and application to variance filtering and covariance localization.  
1337 *Monthly Weather Review*, 143, 1622-1643. doi: 10.1175/MWR-D-14-00157.1
- 1338 Miyoshi, T., & Kondo, K. (2013). A multi-scale localization approach to an en-  
1339 semble Kalman filter. *SOLA (Scientific Online Letters on the Atmosphere)*, 9,  
1340 170-173. doi: 10.2151/sola.2013-038
- 1341 Morozov, V. (1984). *Methods for solving incorrectly posed problems*. Springer.
- 1342 Morzfeld, M., & Hodyss, D. (2023). A theory for why even simple covariance local-  
1343 ization is so useful in ensemble data assimilation. *Monthly Weather Review*.
- 1344 Morzfeld, M., Tong, X. T., & Marzouk, Y. M. (2019). Localization for MCMC: Sam-  
1345 pling high-dimensional posterior distributions with local structure. *Journal of*  
1346 *Computational Physics*, 380, 1–28. doi: 10.1016/j.jcp.2018.12.008
- 1347 Ott, E., Hunt, B., Szunyogh, I., Zimin, A., Kostelich, E., Corazza, M., . . . Yorke,  
1348 J. (2004). A local ensemble Kalman filter for atmospheric data assimilation.  
1349 *Tellus A*, 56, 415-428. doi: 10.1111/j.1600-0870.2004.00076.x
- 1350 Poterjoy, J., & Zhang, F. (2015). Systematic comparison of four-dimensional data

- 1351 assimilation methods with and without the tangent linear model using hybrid  
 1352 background error covariance: E4DVar versus 4DEnVar. *Monthly Weather*  
 1353 *Review*, *143*(5), 1601–1621.
- 1354 Pourahmadi, M. (2011). Covariance Estimation: The GLM and Regularization Per-  
 1355 spectives. *Statistical Science*, *26*(3), 369 – 387. doi: 10.1214/11-STS358
- 1356 Rasmussen, C. E., & Williams, C. (2005). *Gaussian processes for machine learning*  
 1357 (*adaptive computation and machine learning*). The MIT Press.
- 1358 Schillings, C., & Stuart, A. M. (2017). Analysis of the ensemble Kalman filter for in-  
 1359 verse problems. *SIAM Journal on Numerical Analysis*, *55*(3), 1264–1290.
- 1360 Schillings, C., & Stuart, A. M. (2018). Convergence analysis of ensemble Kalman in-  
 1361 version: the linear, noisy case. *Applicable Analysis*, *97*(1), 107–123.
- 1362 Schneider, T., Leung, L. R., & Wills, R. C. J. (2024). Opinion: Optimizing climate  
 1363 models with process-knowledge, resolution, and AI. *EGUsphere*, *2024*, 1–26.  
 1364 doi: 10.5194/egusphere-2024-20
- 1365 Schneider, T., Stuart, A. M., & Wu, J.-L. (2021, 12). Learning stochastic closures  
 1366 using ensemble Kalman inversion. *Transactions of Mathematics and Its Appli-*  
 1367 *cations*, *5*(1), ttab003. doi: 10.1093/imatrm/tnab003
- 1368 Schoenberg, I. J. (1942). Positive definite functions on spheres. *Duke Mathematical*  
 1369 *Journal*, *9*(1), 96 – 108. doi: 10.1215/S0012-7094-42-00908-6
- 1370 Schur, J. (1911). Bemerkungen zur Theorie der beschränkten Bilinearformen mit  
 1371 unendlich vielen Veränderlichen. *Journal für die reine und angewandte Mathe-*  
 1372 *matik*, *140*, 1-28.
- 1373 Talagrand, O., & Courtier, P. (1987). Variational assimilation of meteorological ob-  
 1374 servations with the adjoint vorticity equation. I: Theory. *Quarterly Journal of*  
 1375 *the Royal Meteorological Society*, *113*(478), 1311–1328.
- 1376 Tippett, M., Anderson, J., Bishop, C., Hamill, T., & Whitaker, J. (2003). Ense-  
 1377 mble square root filters. *Monthly Weather Review*, *131*(7), 1485 - 1490. doi: 10  
 1378 .1175/1520-0493(2003)131<1485:ESRF>2.0.CO;2
- 1379 Tong, X., & Morzfeld, M. (2023). Localized ensemble Kalman inversion. *Inverse*  
 1380 *Problems*, *39*(6), 064002.
- 1381 Vishny, D., Morzfeld, M., & Gwirtz, K. (2024, July). *Matlab code & data for “High-*  
 1382 *dimensional covariance estimation from a small number of samples”*. Zenodo.  
 1383 Retrieved from <https://doi.org/10.5281/zenodo.12701351> doi: 10.5281/  
 1384 zenodo.12701351
- 1385 Wainwright, M. J. (2019). *High-dimensional statistics: A non-asymptotic viewpoint*.  
 1386 Cambridge University Press. doi: 10.1017/9781108627771
- 1387 Ward, S. H., & Hohmann, G. W. (2012). Electromagnetic theory for geophysi-  
 1388 cal applications. In *Electromagnetic methods in applied geophysics: Volume 1,*  
 1389 *theory* (p. 130-311). Society of Exploration Geophysicists. doi: 10.1190/1  
 1390 .9781560802631.ch4
- 1391 Whitaker, J., & Hamill, T. (2012). Evaluating methods to account for system errors  
 1392 in ensemble data assimilation. *Mon. Wea. Rev.*, *140*, 3078-3089.
- 1393 Xue, L., Ma, S., & Zou, H. (2012). Positive-definite  $\ell_1$ -penalized estimation of large  
 1394 covariance matrices. *Journal of the American Statistical Association*, *107*(500),  
 1395 1480-1491. doi: 10.1080/01621459.2012.725386
- 1396 Zhang, F., Zhang, M., & Hansen, J. A. (2009). Coupling ensemble Kalman filter  
 1397 with four-dimensional variational data assimilation. *Advances in Atmospheric*  
 1398 *Sciences*, *26*, 1-8.
- 1399 Zhen, Y., & Zhang, F. (2014). A probabilistic approach to adaptive covariance local-  
 1400 ization for serial ensemble square root filters. *Monthly Weather Review*, *142*,  
 1401 4499-4518. doi: 10.1175/MWR-D-13-00390.1