

# Parallel Direct Solution of the Covariance-Localized Ensemble Square Root Kalman Filter Equations with Matrix Functions

JEFFREY L. STEWARD

*University of California, Los Angeles, Los Angeles, California*

JOSE E. ROMAN AND ALEJANDRO LAMAS DAVIÑA

*Universitat Politècnica de València, València, Spain*

ALTUĞ AKSOY

*Cooperative Institute for Marine and Atmospheric Studies, University of Miami, and Hurricane Research Division, NOAA/AOML, Miami, Florida*

(Manuscript received 20 January 2018, in final form 20 March 2018)

## ABSTRACT

Recently, the serial approach to solving the square root ensemble Kalman filter (ESRF) equations in the presence of covariance localization was found to depend on the order of observations. As shown previously, correctly updating the localized posterior covariance in serial requires additional effort and computational expense. A recent work by Steward et al. details an all-at-once direct method to solve the ESRF equations in parallel. This method uses the eigenvectors and eigenvalues of the forward observation covariance matrix to solve the difficult portion of the ESRF equations. The remaining assimilation is easily parallelized, and the analysis does not depend on the order of observations. While this allows for long localization lengths that would render local analysis methods inefficient, in theory, an eigenpair-based method scales as the cube number of observations, making it infeasible for large numbers of observations. In this work, we extend this method to use the theory of matrix functions to avoid eigenpair computations. The Arnoldi process is used to evaluate the covariance-localized ESRF equations on the reduced-order Krylov subspace basis. This method is shown to converge quickly and apparently regains a linear scaling with the number of observations. The method scales similarly to the widely used serial approach of Anderson and Collins in wall time but not in memory usage. To improve the memory usage issue, this method potentially can be used without an explicit matrix. In addition, hybrid ensemble and climatological covariances can be incorporated.

## 1. Introduction

Data assimilation of increasingly plentiful satellite and radar observations requires efficient and accurate algorithms. A single overpass of a polar orbiting satellite over a regional numerical weather prediction (NWP) domain can produce tens of thousands of potentially usable observations, especially when all-sky observations are considered. The Japanese K computer assimilates radar observations every 30 s with a 100-m grid spacing (Miyoshi et al. 2016), and with the next-generation GOES-16 (Schmit et al. 2017) and Himawari-8 (Bessho et al. 2016) geostationary observing platforms providing

observations with approximately kilometer resolution approximately every 5 minutes, data assimilation algorithms need to handle increasingly large data volumes to keep pace. In this paper, we describe a new, efficient, and parallel technique for solving the covariance-localized square root ensemble Kalman filter equations that overcomes several issues in previously described implementations.

The ensemble Kalman filter, first introduced by Evensen (1994), is one of the most widely used methods for data assimilation. Using an ensemble with a relatively small number of members to estimate the flow-dependent background error covariance from the Kalman filter as originally formulated (Kalman 1960) made it feasible to run statistical data assimilation problems even on very large domains. However, two main issues became

---

*Corresponding author:* Jeffrey L. Steward, [jsteward@jifresse.ucla.edu](mailto:jsteward@jifresse.ucla.edu)

apparent in the implementation of the ensemble Kalman filter. The first is that using the same observations to update the mean and ensemble perturbations leads to a systematic underestimation of covariance. Second, the unlocalized estimated covariances contain sample error due to the low number of ensemble members used, leading to spurious relationships.

The issue of systematic covariance underestimation was first solved by perturbing observations with independently sampled noise for each ensemble member (Houtekamer and Mitchell 1998; Burgers et al. 1998). While this solves the underestimation of covariance, adding additional noise increases sampling error, causing the filter to be suboptimal, especially when the ensemble size is small (Whitaker and Hamill 2002). Subsequently, the ensemble square root filter (ESRF) was introduced; it corrects for the underrepresentation of error covariance by adding a square root term to the Kalman update for the ensemble. Various flavors of ESRF have been developed (Bishop et al. 2001; Anderson 2001; Whitaker and Hamill 2002), which Tippett et al. (2003) showed are all equivalent in the sense that they perform analysis in the same vector space and find the same covariance. These methods as originally formulated assume the rank of the covariance matrices is the number of ensemble members.

Independently from covariance underestimation, the issue of spurious correlations due to small ensemble size has been addressed in two main ways: covariance localization and local analysis. Sakov and Bertino (2011) demonstrated that these two approaches are approximately equal, and the choice of approach is therefore dependent upon other factors. Critically, the localization radius used in local methods will determine their efficiency, and large localization radii will require repetitive solution of large problems for each grid point. In this work, we investigate covariance localization, which uses a Schur product (component-wise multiplication) to zero out correlations farther than a specified distance (Gaspari and Cohn 1999; Houtekamer and Mitchell 2001; Hamill et al. 2001). This causes the rank of the forward observation covariance matrix used in the inverse of the Kalman gain to increase beyond the number of ensemble members. As shown in Steward et al. (2017, hereafter S17), a relatively short localization radius will lead to a full-rank forward observation covariance matrix, while a long localization radius will lead to a rank deficient one.

The combination of these factors leads to several different possibilities for scalable parallel implementations of the ensemble Kalman filter equations. Local methods with perturbed observations and covariance localization include Keppenne and Rienecker (2002),

Houtekamer et al. (2014), Bishop et al. (2015), and Niño-Ruiz et al. (2015), while local analysis methods based on the ESRF equations include Ott et al. (2004), Anderson (2003), Zhang et al. (2005), Hunt et al. (2007), Wang et al. (2013), and Niño-Ruiz et al. (2018). Note that the widely used local ensemble Kalman transform filter (LETKF) of Hunt et al. (2007) applies a localization strategy based on the observation error covariance matrix  $\mathbf{R}$  rather than on the sample covariance matrices estimated by the ensemble. The widely used and highly efficient method of Anderson and Collins (2007, hereafter AC07) is a “global” analysis (i.e., nonlocal) parallel implementation based on the serial assimilation of the ESRF equations with covariance localization. This method also treats the observations as part of an augmented state in order to update the observations in parallel without requiring excessive communication. Houtekamer and Mitchell (2001) describe a global analysis method with perturbed observations and covariance localization.

Because of the difficulties in solving the global ESRF equations directly, in implementations such as Anderson (2001), Whitaker and Hamill (2002), AC07, and Aksoy (2013), a serial approach is utilized where a single observation is assimilated at a time. This approach is provably identical to the global analysis without covariance localization and linear observation operators. However, with covariance-based localization, the ordering of observations affects the analysis, as shown in Nerger (2015) and Bishop et al. (2015), due to the nonlinear nature of covariance localization. In other words, in the presence of ensemble sample covariance localization, serially assimilating observation  $A$  before observation  $B$  may give different results than assimilating observation  $B$  before  $A$ . The magnitude of this issue has not yet been fully explored.

As shown in Bishop et al. (2015), the issue of observation-ordering-dependent analysis in serial covariance-localized methods stems from the inconsistent application of the high-rank localized covariance matrices. In particular, when covariance localization is used, the matrix to be inverted in the Kalman gain becomes full rank or nearly full rank (as shown in, e.g., S17). Without covariance localization (or in a local analysis method that does not increase the rank of the matrix using a Schur product), as shown in Tippett et al. (2003), the Sherman–Woodbury update is sufficient for an unlocalized matrix, as the rank of the matrix is at most the number of ensemble members (Godinez and Moulton 2012). However, the fundamental shift to high-rank matrices requires additional effort to correct.

Several strategies have been proposed to handle this observation-ordering dependence within a serial filter.

Bishop et al. (2015) propose the consistent hybrid ensemble filter (CHEF) with local analysis and perturbed observations that will ensure the analysis is consistent and does not depend on the order of assimilation. Kotsuki et al. (2017) present a study of observation ordering with a Lorenz-96 model and investigate rules for observation assimilation ordering to minimize analysis forecast error. The method of correcting sample correlation described in Anderson (2012) has also been used to reduce the dependence of observation ordering in a serial filter (J. Anderson 2017, personal communication).

Extending upon these works, as an alternative to attempting to apply and update the high-rank localized matrices serially in a consistent way, we propose assimilating all observations within the assimilation window in a single pass as a potential alternative. In other words, we do not utilize the single observation processing strategy normally employed for serial filter solutions and instead solve the ESRF equations directly. This is done by dividing the necessary matrix operators across the set of processing elements in a “top down” fashion, as opposed to the “bottom up” approach of local analysis. This method was utilized in S17 to provide a global, “all at once,” parallel, direct solution of the covariance-localized ESRF equations. Note that “all at once” here is used to refer to assimilating all observations that the serial filter would assimilate one by one but not all observations within all assimilation windows at once; that is, the method in S17 as well as the one presented below are both sequential filters in that batches of observations can also be assimilated. The benefit of this approach is that the analysis consistently applies the high-rank covariance-localized matrices and, as a result, does not depend on the order of observations. It provides a solution to the ESRF equations with a proven error bounds that can be used as a benchmark against other methodologies.

The cost of this approach is that a product with the entire full-rank matrix inverse (which also requires a square root term) of the forward observation error covariance is required. S17 solve for eigenvalues and eigenvectors of the observation covariance matrix and use the ESRF matrix function “scalarized” on the eigenvalues to find the required matrix inverses and products. As eigenpairs are extremely convenient for mathematical analysis, the approach in S17 also includes an error bounds related to the smallest eigenvalue used. The final analysis is also shown not to depend on the ordering of observations. This error-bounded method, which directly solves the ESRF equations, is therefore a highly accurate solution to the ESRF equations known to be the minimum variance solution to the data assimilation problem.

However, as predicted by theory and shown in this work, while the method described in S17 is accurate to within a configurable tolerance, it is impractical for large numbers of observations due to the nature of the eigenproblem, where, for general matrices, finding a large number of eigenpairs scales as  $O(n^3)$  for a matrix of size  $n \times n$  (Golub and Van Loan 1996). The quantity  $n$  is the number of quality-controlled observations in this case. This paper extends S17 to take advantage of recent improvements in the theory and computation of matrix functions to transform the problem of solving the difficult inverse and square root portion of the ESRF equations into computing matrix-vector products that are used to build up a Krylov subspace and, through a library call, applying the matrix function directly to a small dense matrix. This small dense matrix represents the compression of the larger localized forward observation covariance matrix onto the reduced-order Krylov subspace basis.

As we show below, this matrix function method gives results that are practically identical to the error-bounded methodology of S17 but is much more computationally efficient. As only a matrix-vector product with the observation covariance matrix is required, this matrix function approach is well suited for a matrix-free implementation where the covariance matrix is not explicitly formed. This method is also amenable to hybrid covariance models using both ensemble and climatological covariances.

We implement the matrix function method and compare the performance results with both S17 as well as the parallel augmented-state method of AC07. As a proof-of-concept application, we test this method on the difficult, highly nonlinear case of first-cycle tropical cyclone (TC) data assimilation. In this case, the background ensemble can contain position errors of features, and the posterior analysis increment can be large (e.g., Chang et al. 2014). As we show, the order-dependence issue of a serial filter is nontrivial in this case. To demonstrate the unique properties of our new method, we investigate TC assimilation with a long covariance length scale that would be impractical for local analysis methods. As we show, the matrix function method is roughly comparable in terms of wall-time performance to AC07 and far superior to S17. The analysis results do not depend on observation ordering, like S17 but contrary to AC07. However, our results demonstrate the memory scaling of the matrix function method is inferior to AC07 and suggest that matrix-free methods would be required to scale this method to the order of millions of observations at once.

This paper is organized as follows. Section 2 summarizes S17 in order to build upon it. In section 3, the

eigenpair computation of S17 is replaced with a much more efficient matrix function (MFN)-based approach that uses a basis for the Krylov space to compress the forward observation covariance matrix and apply the covariance-localized ESRF matrix functions to this reduced-order matrix. Section 4 summarizes AC07. Section 5 presents numerical results of the matrix function approach and a performance comparison to S17 and AC07. Finally, section 6 presents conclusions and a discussion.

## 2. Eigenvalue–eigenvector solution of S17

In this section, we briefly review S17 in order to introduce the new matrix function method that extends it. Given an ensemble  $\mathbf{X}_f$  of a previous forecast, the updated analysis to the ensemble mean  $\bar{\mathbf{x}}_f$  of size  $N_{\text{state}} \times 1$ , and ensemble perturbations  $\mathbf{X}'_f$  of size  $N_{\text{state}} \times N_{\text{ens}}$  the square root ensemble Kalman filter without perturbed observations (Whitaker and Hamill 2002) is

$$\begin{aligned}\bar{\mathbf{x}}_a &= \bar{\mathbf{x}}_f + \mathbf{K} \left[ \mathbf{y} - \overline{H(\mathbf{X}_f)} \right], \\ \mathbf{X}'_a &= \mathbf{X}'_f + \tilde{\mathbf{K}}(\mathbf{0} - \mathbf{H}\mathbf{X}),\end{aligned}\quad (1)$$

where  $\mathbf{y}$  ( $N_{\text{obs}} \times 1$ ) are the observations,  $\overline{H(\mathbf{X}_f)}$  ( $N_{\text{obs}} \times 1$ ) is the mean of the forward-calculated observation operators, and  $\mathbf{H}\mathbf{X}_{ij} = h_i[\mathbf{X}_f^{(j)}] - \bar{h}_i(\mathbf{X}_f)$  is the mean-subtracted  $i$ th observation operator acting on the  $j$ th ensemble member  $\mathbf{X}_f^{(j)}$ . The  $\mathbf{H}\mathbf{X}$  matrix is  $N_{\text{obs}} \times N_{\text{ens}}$  (as is  $\mathbf{0}$ , a matrix filled with zeros). The traditional Kalman gain  $\mathbf{K}$  ( $N_{\text{state}} \times N_{\text{obs}}$ ) is

$$\mathbf{K} = \mathbf{C}_{\mathbf{x},H\mathbf{x}} \mathbf{D}^{-1}, \quad (2)$$

where  $\mathbf{C}_{\mathbf{x},H\mathbf{x}} = \text{cov}(\mathbf{x}_f, H(\mathbf{x}_f))$  is the localized covariance between  $\mathbf{x}_f$  (an  $N_{\text{state}} \times 1$  random variable representing the previous forecast) and  $H(\mathbf{x}_f)$  (the observation operator acting on this random variable). The matrix  $\mathbf{D} = \mathbf{C}_{H\mathbf{x},H\mathbf{x}} + \mathbf{R}$  for  $\mathbf{C}_{H\mathbf{x},H\mathbf{x}} = \text{cov}(H(\mathbf{x}_f), H(\mathbf{x}_f))$ , the localized forward observation covariance, and  $\mathbf{R}$  is the observation error covariance  $\text{cov}(\mathbf{y}_t, H(\mathbf{x}_f))$  for a random variable  $\mathbf{y}_t$  representing the true observations without observation noise.

The  $\tilde{\mathbf{K}}$  matrix ( $N_{\text{state}} \times N_{\text{obs}}$ ), the correction from using nonperturbed observations, is

$$\tilde{\mathbf{K}} = \mathbf{C}_{\mathbf{x},H\mathbf{x}} \mathbf{D}^{-1/2} \left( \sqrt{\mathbf{D}} + \sqrt{\mathbf{R}} \right)^{-1}. \quad (3)$$

As detailed in S17, the covariance matrices we consider can include localized ensemble-based correlations

in observation-space and/or variational-style model-space localization. For observation-space localization, a component-wise multiplication  $\circ$  between two matrices is used as

$$\mathbf{C}_{H\mathbf{x},H\mathbf{x}}^{\text{obs}} = \boldsymbol{\rho}_{\mathbf{y},\mathbf{y}} \circ \mathbf{Q}_{H\mathbf{x},H\mathbf{x}}, \quad (4)$$

where  $\boldsymbol{\rho}_{\mathbf{y},\mathbf{y}}$  is the localization matrix arising from a localization function (Gaspari and Cohn 1999)  $\ell$  such that

$$\left( \boldsymbol{\rho}_{\mathbf{y},\mathbf{y}} \right)_{ij} = \ell \left( d_{ij} \middle| L_{ij} \right), \quad (5)$$

where  $d_{ij}$  is the distance between the location of the  $i$ th and  $j$ th observations, and  $L_{ij}$  is the characteristic length scale for the localization function  $\ell$ . The  $\mathbf{Q}_{H\mathbf{x},H\mathbf{x}}$  is the sample covariance matrix

$$\mathbf{Q}_{H\mathbf{x},H\mathbf{x}} = \frac{\mathbf{H}\mathbf{X}(\mathbf{H}\mathbf{X})^T}{N_{\text{ens}} - 1}. \quad (6)$$

Likewise, the observation-space localized model and observation cross-covariance is given by

$$\mathbf{C}_{\mathbf{x},H\mathbf{x}}^{\text{obs}} = \boldsymbol{\rho}_{\mathbf{y},\mathbf{y}} \circ \mathbf{Q}_{\mathbf{x},H\mathbf{x}} \quad (7)$$

for

$$\left( \boldsymbol{\rho}_{\mathbf{y},\mathbf{y}} \right)_{ij} = \ell \left( d_{ij} \middle| L_{ij} \right), \quad (8)$$

where  $d_{ij}$  is the distance between the location of the model state  $i$  and observation  $j$  with the same localization function as Eq. (5), and

$$\mathbf{Q}_{\mathbf{x},H\mathbf{x}} = \frac{\mathbf{X}'_f(\mathbf{H}\mathbf{X})^T}{N_{\text{ens}} - 1}. \quad (9)$$

As noted in Campbell et al. (2010), integrated observations such as satellite scans do not have a particular vertical location to ascribe. In these cases, model-space localization is more applicable. For model-space localization, the observation operator tangent-linear  $\mathbf{H}$  and adjoint  $\mathbf{H}^T$  are applied to the localized model covariance as

$$\mathbf{C}_{H\mathbf{x},H\mathbf{x}}^{\text{model}} = \mathbf{H} \left( \boldsymbol{\rho}_{\mathbf{x},\mathbf{x}} \circ \mathbf{Q}_{\mathbf{x},\mathbf{x}} \right) \mathbf{H}^T, \quad (10)$$

where

$$\mathbf{Q}_{\mathbf{x},\mathbf{x}} = \frac{\mathbf{X}'_f(\mathbf{X}'_f)^T}{N_{\text{ens}} - 1} \quad (11)$$

for the ensemble perturbations  $\mathbf{X}'_f$ , and

$$\left(\boldsymbol{\rho}_{\mathbf{x}\mathbf{x}}\right)_{ij} = \ell\left(d_{ij}\middle|L_{ij}\right), \quad (12)$$

where  $d_{ij}$  is the distance between the location of two model states  $i$  and  $j$  with the same localization function as Eq. (5). Equation (7) is changed analogously as

$$\mathbf{C}_{\mathbf{x},H\mathbf{x}}^{\text{model}} = \left(\boldsymbol{\rho}_{\mathbf{x}\mathbf{x}} \circ \mathbf{Q}_{\mathbf{x}\mathbf{x}}\right)\mathbf{H}^T. \quad (13)$$

Note that all of these localized matrices are sparse, and zero elements (i.e., the correlations farther than the specified localization distance) are not stored in memory or computed. Thus, for example, only those elements of  $\mathbf{Q}_{H\mathbf{x},H\mathbf{x}}$  that will be nonzero after localization are calculated. Furthermore, the full model-space matrix  $\mathbf{Q}_{\mathbf{x}\mathbf{x}}$  will never be explicitly formed due to its prohibitively large size. See S17 for more details.

As we will allow for full-rank matrices, our method is compatible with either of these localization methods, a linear combination of the two, or any other “reasonable” modeled covariance between  $H\mathbf{x}$  and  $H\mathbf{x}$  and  $\mathbf{x}$  and  $H\mathbf{x}$ , which we denote in general  $\mathbf{C}_{H\mathbf{x},H\mathbf{x}}$  and  $\mathbf{C}_{\mathbf{x},H\mathbf{x}}$ . Note that in this work, however, we only present results for the observation-based localization of Eqs. (4) and (7).

We now return to solving Eq. (1). Both S17 and the matrix function approach utilize a preprocessing step of a transformation first introduced in Bishop et al. (2001) to whiten the observations as  $\mathbf{y} = \mathbf{R}_{\text{old}}^{-1/2}\mathbf{y}_{\text{old}}$ , where the “old” subscript represents the untransformed observations. The observation operator is also scaled as  $H(\mathbf{x}) = \mathbf{R}_{\text{old}}^{-1/2}H_{\text{old}}(\mathbf{x})$ . As a result of this preprocessing transformation, the  $\mathbf{R}$  matrix is now identity, which makes Eq. (3) much easier to solve. For the diagonal observation error matrix  $\mathbf{R}_{\text{old}}$  typically used in data assimilation (which assumes uncorrelated observation errors), multiplying by  $\mathbf{R}_{\text{old}}^{-1/2}$  is equivalent to dividing each observation by the standard deviation of the observation error, and for nondiagonal  $\mathbf{R}_{\text{old}}$ , this transformation removes that off-diagonal correlation using principal components.

As  $\mathbf{D}_{\text{new}} = \mathbf{C}_{H\mathbf{x},H\mathbf{x}} + \mathbf{I}$  by this transformation (note that we drop the “new” subscript in what follows, as it could be applied to virtually all matrices; i.e., we write  $\mathbf{D}_{\text{new}}$  as  $\mathbf{D}$  in a slight abuse of notation), this leads to

$$\tilde{\mathbf{K}} = \mathbf{C}_{\mathbf{x},H\mathbf{x}} \mathbf{M}^{-1} \quad (14)$$

for  $\mathbf{M} = \mathbf{D} + \sqrt{\mathbf{D}}$ . Let  $\lambda_i, \mathbf{v}_i$  denote the  $i$ th eigenpair of  $\mathbf{C}_{H\mathbf{x},H\mathbf{x}}$ . Then,

$$\mathbf{M}\mathbf{v}_i = \mathbf{C}_{H\mathbf{x},H\mathbf{x}}\mathbf{v}_i + \mathbf{v}_i + \left(\mathbf{C}_{H\mathbf{x},H\mathbf{x}} + \mathbf{I}\right)^{1/2}\mathbf{v}_i. \quad (15)$$

As shown in S17, we have

$$\mathbf{M}\mathbf{v}_i = \left[\lambda_i + 1 + (\lambda_i + 1)^{1/2}\right]\mathbf{v}_i. \quad (16)$$

Therefore,

$$\mathbf{M}^{-1}\mathbf{v}_i = \lambda'_i\mathbf{v}_i \quad (17)$$

for

$$\lambda'_i = \frac{1}{\lambda_i + 1 + (\lambda_i + 1)^{1/2}}. \quad (18)$$

We find the largest  $r$  eigenvalues and corresponding eigenvectors of  $\mathbf{C}_{H\mathbf{x},H\mathbf{x}}$ , where  $r$  is chosen such that  $\lambda_{r+1} \leq \varepsilon_\lambda$  for some small constant  $\varepsilon_\lambda$ , and we can therefore solve

$$\mathbf{M}^{-1}(\mathbf{0} - H\mathbf{X})_j \approx \sum_{i=1}^r \lambda'_i \alpha_{ij} \mathbf{v}_i \quad (19)$$

for  $\alpha_{ij} = -\mathbf{v}_i^T H\mathbf{X}_j$ . An error bound on this approximation related to  $\varepsilon_\lambda$  is proved in S17.

Similarly, for the mean update,

$$\mathbf{D}^{-1}\left[\mathbf{y} - \overline{H(\mathbf{X}_f)}\right] \approx \sum_{i=1}^r \frac{\beta_i}{\lambda_i + 1} \mathbf{v}_i, \quad (20)$$

where  $\beta_i = \mathbf{v}_i^T [\mathbf{y} - \overline{H(\mathbf{X}_f)}]$ .

The  $N_{\text{obs}} - N_{\text{ens}}$  matrix ( $\mathbf{E}|\mathbf{g}$ ), where

$$\mathbf{E}_j = \sum_{i=1}^r \lambda'_i \alpha_{ij} \mathbf{v}_i, \quad (21)$$

and

$$\mathbf{g} = \sum_{i=1}^r \frac{\beta_i}{\lambda_i + 1} \mathbf{v}_i, \quad (22)$$

is then distributed to all processing elements. The remaining Kalman gain from Eq. (1) only requires multiplication with  $\mathbf{C}_{\mathbf{x},H\mathbf{x}}$  which can proceed in an embarrassingly parallel fashion. This makes an efficient parallel method that only requires the eigenpairs of the  $N_{\text{obs}} \times N_{\text{obs}}$  sparse, to be a positive semidefinite symmetric matrix  $\mathbf{C}_{H\mathbf{x},H\mathbf{x}}$ . The Scalable Library for Eigenvalue Problem Computations (SLEPc; Hernandez et al. 2005), which is built upon the Portable Extensible Toolkit for Scientific Computing (PETSc; Balay et al. 1997, 2016, 2017), is used to solve this eigenproblem using sparse matrices in a manner that scales well as a function of the number of processors, as shown in S17.

### 3. New matrix function approach

We first note that while S17 evaluates the largest  $r$  eigenpairs of  $\mathbf{C}_{H\mathbf{x},H\mathbf{x}}$  in order to solve Eq. (1), only those



eigenvectors  $i$  such that  $\alpha_{i,j}$  for all  $j$  and  $\beta_i \neq 0$  are required. This suggests a more efficient solution that does not require all eigenpairs. In this section, we develop such a solution that requires only the matrix-vector product  $\mathbf{C}_{H_x, H_x} \mathbf{b}$  for some vector  $\mathbf{b}$  to compute a reduced-order, accurate basis for representation of the ESRF matrix functions.

In addition to solving the eigenproblem, SLEPc can also evaluate the action of a matrix function on a vector  $\mathbf{z} = f(\mathbf{A})\mathbf{b}$ , where  $\mathbf{z}$  and  $\mathbf{b}$  are vectors,  $\mathbf{A}$  is a matrix, and  $f$  is a matrix function in the sense given in Higham (2008). In the case of the mean  $\mathbf{K}$  in Eq. (2) given above,

$$f_1(\mathbf{D}) = \mathbf{D}^{-1}, \tag{23}$$

while for  $\tilde{\mathbf{K}}$  in Eq. (14),

$$f_2(\mathbf{D}) = (\mathbf{D} + \sqrt{\mathbf{D}})^{-1}. \tag{24}$$

Recall that  $\mathbf{D} = \mathbf{C}_{H_x, H_x} + \mathbf{I}$ . Also note that  $f_1$  involves the standard linear system of equations  $\mathbf{D}\mathbf{x} = \mathbf{b}$  solving for  $\mathbf{x}$ , which is normally handled by other methods; in this work, we test using the matrix function approach for both the mean and the perturbations.

The matrix function solvers in SLEPc are based on Krylov subspace methods (Higham 2008, chapter 13). Earlier works using Krylov subspace methods to approximate matrix functions include Van Der Vorst (1987), Saad (1992), and Hochbruck and Lubich (1997). These methods are appropriate for the case of our large, high-rank matrix  $\mathbf{D}$ , as they compute the result  $\mathbf{z}$  without explicitly building the matrix  $f(\mathbf{D})$ . The calculation of  $f(\mathbf{D})\mathbf{b}$  proceeds in a manner similar to the Arnoldi method (Arnoldi 1951) for finding eigenpairs. At the first step,  $\mathbf{V}_1 = \mathbf{b}/\|\mathbf{b}\|_2$  and at step  $m$ , given an  $N_{\text{obs}} \times (m-1)$  orthonormal basis  $\mathbf{V}_{m-1}$  of the Krylov subspace  $\mathcal{K}_{m-1}(\mathbf{D}, \mathbf{b}) = \text{span}\{\mathbf{b}, \mathbf{D}\mathbf{b}, \mathbf{D}^2\mathbf{b}, \dots, \mathbf{D}^{m-2}\mathbf{b}\}$  we seek the orthonormal basis  $\mathbf{V}_m$  that spans  $\mathcal{K}_m(\mathbf{D}\mathbf{b})$ . This is done by the Arnoldi relation  $\mathbf{D}\mathbf{V}_{m-1} = \mathbf{V}_{m-1}\mathbf{H}_{m-1} + h_{m,m-1}\mathbf{v}_m\mathbf{e}_{m-1}^T$ , where  $\mathbf{H}_{m-1}$  is an  $(m-1) \times (m-1)$  upper-Hessenberg matrix that contains the values of the projections of  $\mathbf{D}$  onto the basis  $\mathbf{V}_{m-1}$ ,  $\mathbf{v}_m$  is the  $m$ th column to be added to  $\mathbf{V}_m$  this iteration, and  $h_{m,m-1}$  is the  $(m, m-1)$  entry in the  $\mathbf{H}_m$  matrix. Term  $\mathbf{e}_{m-1}$  is the  $m-1$  unit coordinate vector, so  $h_{m,m-1}\mathbf{v}_m\mathbf{e}_{m-1}^T$  is the  $N_{\text{obs}} \times (m-1)$  zero matrix, except column  $m-1$ , which is  $h_{m,m-1}\mathbf{v}_m$ . Once  $\mathbf{V}_m$  is found, the approximation of  $\mathbf{z}$  can be computed as

$$\tilde{\mathbf{z}}_m = \beta \mathbf{V}_m f(\mathbf{H}_m) \mathbf{e}_1, \tag{25}$$

where  $\beta = \|\mathbf{b}\|_2$ . The  $\mathbf{e}_1$  is the first coordinate vector, so right multiplying by it gives the first column of  $\beta \mathbf{V}_m f(\mathbf{H}_m)$  in Eq. (25). Note that  $\mathbf{b} = \beta \mathbf{V}_m \mathbf{e}_1$ . In addition, note that

$\mathbf{H}_m$  represents the compression of  $\mathbf{D}$  onto  $\mathcal{K}_m(\mathbf{D}, \mathbf{b})$  with respect to the basis  $\mathbf{V}_m$ . Hence, the problem of computing the function of a large matrix  $\mathbf{D}$  of order  $N_{\text{obs}}$  is reduced to computing the function of a small matrix  $\mathbf{H}_m$  of order  $m$  with  $m \ll N_{\text{obs}}$ . For the latter task, we can employ algorithms for dense matrices as discussed below.

Note that in the above description, the Arnoldi process requires a numerically stabilized Gram–Schmidt process to orthonormalize the basis vectors in a way that the final result is not overly affected by numerical noise. Furthermore, the parallelization of this stabilized process requires careful implementation to avoid negatively impacting performance by creating bottlenecks. Thus, the relatively straightforward (conceptually) Gram–Schmidt process becomes rather complex when implemented in a parallel setting as discussed in Björck (1994) and Frayssé et al. (1998). SLEPc utilizes an efficient parallel version of the iterated classical Gram–Schmidt (ICGS) in the Arnoldi process that does not require global communication but maintains numerical stability. As in the high-level description given above, the resulting projections in the ICGS process onto the previous basis vectors are stored in the  $\mathbf{H}_m$  matrix. For more details on the orthogonalization process in SLEPc, see Hernandez et al. (2007).

The  $m$  parameter is of paramount importance for this method. If  $m$  is too small, the Krylov subspace will not contain enough information to build an accurate approximation. On the other hand, if  $m$  is too large, the memory requirements for storing  $\mathbf{V}_m$  (as well as the computational cost) will be prohibitive. For this reason, SLEPc implements a restarted variant of the method, where  $m$  is prescribed to a fixed value; here, we use  $m = 150$  which, as shown below, is based on testing for our particular application. When the subspace reaches this size, a restart is carried out by keeping part of the data computed so far and discarding unnecessary information. Investigation into restarting matrix function iterations is still an area of active research (Afanasjew et al. 2008; Eiermann et al. 2011; Frommer et al. 2017). SLEPc implements the Eiermann–Ernst restart (Eiermann and Ernst 2006), in which only the last basis vector  $\mathbf{v}_{m+1}$  is kept (in order to continue the Arnoldi recurrence) along with the matrix  $\mathbf{H}_m$  that is “glued” together with the previous ones. After  $k$  restarts, the matrix used in the approximation (25) has the form

$$\mathbf{H}_{k \times m} = \begin{bmatrix} \mathbf{H}_{(k-1)m} & \mathbf{0}_m \\ h_{m+1,m}^{(k-1)} \mathbf{e}_1 \mathbf{e}_{(k-1)m}^T & \mathbf{H}_m^{(k)} \end{bmatrix}, \tag{26}$$

where  $\mathbf{H}_m^{(k)}$  is the matrix computed by the Arnoldi method in the  $k$ th restart. Note that in the Eiermann–Ernst restart, the glued matrix [(26)] is not used directly

in (25) because  $\mathbf{H}_{k \times m}$  has size  $km \times km$ , but  $\mathbf{V}_m$  has only  $m$  columns. Therefore, only the last  $m$  components of the vector  $f(\mathbf{H}_{k \times m})\mathbf{e}_1$  are used in (25) to give a correction to be added to the approximation available in the previous restart. This correction is given by  $\tilde{\mathbf{z}}^{(k)} = \tilde{\mathbf{z}}^{(k-1)} + \mathbf{c}^{(k)}$ , where

$$\mathbf{c}^{(k)} = \beta \mathbf{V}_m^{(k)} [\mathbf{0}, \mathbf{I}_m] f(\mathbf{H}_{k \times m})\mathbf{e}_1, \quad (27)$$

and  $\mathbf{V}_m^{(k)}$  is the basis computed in the last restart. Equations (25)–(27) are implemented in a numerically efficient way in SLEPc.

SLEPc bases the stopping criterion on the norm of the correction; that is, restarting continues until  $\|\mathbf{c}^{(k)}\|_2 < \beta \times \varepsilon_{\text{tol}}$  for some user-defined  $\varepsilon_{\text{tol}}$  ( $10^{-8}$  by default for 8-byte floating point precision). As noted in Eiermann and Ernst (2006), the Arnoldi method converges rapidly with superlinear behavior for smooth functions. The convergence behavior when including restarting is presented in Afanasjew et al. (2008) for a related method.

In this work, we are interested in solving  $\mathbf{g} = f_1(\mathbf{D})[\mathbf{y} - H(\mathbf{X}_f)]$  to replace Eq. (22) and  $\mathbf{E}_j = f_2(\mathbf{D})(\mathbf{0} - H\mathbf{X})_j$  to replace Eq. (21) for  $j = 1, \dots, N_{\text{ens}}$ . Applying the method described above leads to the evaluation of  $f_1(\mathbf{H}_1)$  and  $f_2(\mathbf{H}_2)$  explicitly for small dense matrices  $\mathbf{H}_1$  and  $\mathbf{H}_2$  of the form in Eq. (26). Note that these matrices are not symmetric even though  $\mathbf{D}$  is symmetric, and also note that the matrices grow at each restart of the Krylov method.

SLEPc allows flexibility in the definition of functions by combining two simpler functions. In our case, we define  $f_1(\cdot)$  as the reciprocal of the identity function and  $f_2(\cdot)$  as the reciprocal of another function, which in turn is defined as the sum of two functions (identity and the square root). All these subfunctions can be evaluated easily, except the matrix square root. For this, SLEPc implements a reduction to (real) Schur form followed by a block version of a Schur algorithm (Higham 1987; Deadman et al. 2012).

Note that only the matrix action  $\mathbf{D}\mathbf{b}$  is required in this algorithm, allowing for matrix-free implementations. This could be potentially useful for defining matrix-vector products using the “modulation product” defined in Bishop and Hodyss (2009) or for variational-style covariances that use fast Fourier transforms (FFTs) to define the action of a circulant covariance matrix. Hybrid methods are also possible; as long as the action of the covariance  $\mathbf{C}_{H\mathbf{x}, H\mathbf{x}}$  as well as  $\mathbf{C}_{\mathbf{x}, H\mathbf{x}}$  can be applied, any such modeled covariance can be imposed on the analysis through the ESRF equations through this approach.

#### 4. Serial augmented-state filter of AC07

To compare the performance of our new matrix function approach to an existing method, we briefly

summarize the method of AC07 here. AC07 details a highly scalable approach to solving the ESRF equations in serial that is provably identical to the global solution with linear observation operators and without covariance localization. With covariance localization, however, the results will depend upon the ordering of observations as discussed above, although to what extent this difference will impact ensemble NWP forecasts has not yet been explored.

AC07 describe an algorithm that loops over each observation in serial. Each observation is owned by a particular processing element. For each observation  $n$ , the owner of that observation broadcasts the observation details (including the observation location, ensemble forward-calculated values  $h_n(\mathbf{x}_j)$  for  $j = 1, \dots, N_{\text{ens}}$  and QC status) to the other processing elements, which then each process the observation in parallel. An important innovation of AC07 is the treatment of observations themselves as part of the augmented state vector. In other words, just as water vapor, temperature, and other geophysical variables are updated by the Kalman filter equations, the observations (which are assumed to have a particular location in space) are also updated during the assimilation process. Thus, the  $n$ th observation that is broadcast by the owner-processing element will have been potentially updated by observations 1 through  $n - 1$ . This saves the computational expense of having to communicate in order to recompute the observation operators.

A scalar form of the ESRF [Eq. (1)] is used to efficiently update all of the covariance-localized state points and observations. The mean of each state  $i$  is updated as indicated with an overbar

$$\bar{\mathbf{x}}_i = \bar{\mathbf{x}}_i + k_{i,n} \left[ \mathbf{y}_n - \overline{H(\mathbf{X}_f)_n} \right] \quad (28)$$

for the Kalman gain  $k_{i,n}$  from Eq. (2), scalarized for point  $i$  for observation  $n$  as

$$k_{i,n} = \frac{\rho_{i,n}}{d_n} \frac{1}{N_{\text{ens}} - 1} \sum_{j=1}^{N_{\text{ens}}} (\mathbf{X}'_f)_{ij} (H\mathbf{X})_{nj}. \quad (29)$$

Here,

$$d_n = \frac{1}{N_{\text{ens}} - 1} \sum_{j=1}^{N_{\text{ens}}} (H\mathbf{X})_{nj}^2 + \mathbf{R}_{n,n}, \quad (30)$$

where  $\mathbf{R}_{n,n}$  ( $\mathbf{R}$  is assumed diagonal) is the observation error variance of the  $n$ th observation, and  $\rho_{i,n}$  is the localization factor between the state point  $i$  and observation  $n$ ; that is, it corresponds to the  $(i, n)$  component of

the  $\boldsymbol{\rho}_{\mathbf{x}\mathbf{x}}$  matrix in Eq. (8), although this matrix is not formed in this implementation.

Similarly, given the scalar square root correction

$$\beta_n = \frac{1}{1 + \sqrt{r_d}}, \quad (31)$$

where

$$r_d = \frac{\mathbf{R}_{n,n}}{d_n}, \quad (32)$$

the  $j$ th ensemble perturbation at state point  $i$  is updated as

$$\mathbf{X}'_{i,j} = \mathbf{X}'_{i,j} + \beta_n k_{i,n} [\mathbf{0} - (H\mathbf{X})_{n,j}]. \quad (33)$$

Note that the analogous equations are used to update the approximations of the forward observation mean  $[H(\mathbf{X}_f)]_k$  and perturbations  $H(\mathbf{X}_f)_{k,j}$  for  $k = n + 1$  to  $N_{\text{obs}}$ ; that is, the remaining unassimilated forward observations are treated as part of the augmented state vector.

## 5. Numerical results

The implementation described in section 3 was used to replace the computation of  $(\mathbf{E}|\mathbf{g})$  from Eqs. (21) and (22) from S17, retaining the remaining components. For comparison, the serial method of AC07 was implemented and tested as well. To ensure consistent comparisons, an object-oriented approach was incorporated in the Hurricane Ensemble Data Assimilation System (HEDAS; Aksoy et al. 2012, 2013; Aksoy 2013; Vukicevic et al. 2013; Aberson et al. 2015) to maintain consistency in observation processing, quality control (QC), and disk input/output among all three implementations. Only the filter aspect differs.

All timings were tested on the NOAA jet supercomputing system xjet installed in 2015/16, where each node has 24 cores with a 2.3-GHz Intel Haswell CPU and 2.66 GB RAM connected via FDR Infiniband. As a proof of concept for this method, we ran two experiments, each with 30 Hurricane WRF (HWRF; Gopalakrishnan et al. 2011) ensemble members, using the Hurricane Edouard (2014) study described in Christophersen et al. (2017). Both of these experiments use quality-controlled observations from sources including satellite retrievals and the NASA AV6 Global Hawk 20140916GH Storm Survey mission (Zawislak et al. 2016; Rogers et al. 2016; Christophersen et al. 2017).

To illustrate the performance on a relevant single cycle as in Christophersen et al. (2017), the first experiment uses HWRF to spin up 30 GFS ensemble

members initialized at 1200 UTC 16 September 2014 for 4 h, then assimilates 15 200 quality-controlled observations from this set at 1600 UTC 16 September 2014  $\pm$  30 min using HEDAS. The localization length scale was set to  $L = 240$  as  $c = L/2$  from Eq. (4.10) of Gaspari and Cohn (1999), as described in S17. Figure 1 shows the analyzed water vapor field at level 20 (out of 60) for the eigenproblem-based solution (EPS), MFN, and serial implementation of AC07. Ten different random observation orderings were assimilated. The mean and standard deviation of the 10 different AC07 analyses are shown in Figs. 1a and 1b. As shown, the standard deviation of these different orderings can reach up to approximately  $1.5 \text{ g kg}^{-1}$ . The same 10 random orderings were assimilated with the MFN solution as shown in Figs. 1c and 1d. Each time, the MFN analysis was identical to within  $10^{-7}$ ; the standard deviation is less than  $10^{-7}$  (“zero”) as well. For comparison, the absolute difference between the average serial analysis and the EPS analysis is shown in Fig. 1e, which as shown is greater than  $2 \text{ g kg}^{-1}$  in places. The absolute difference between the MFN and EPS solution is shown in Fig. 1f, which is also “zero.”

To emphasize the order-independence issue, Fig. 2 shows the assimilation of the first two random observation orderings assimilated in Fig. 1 (order 1 and order 2). No effort was made to maximize this difference for AC07—the first two random orderings were chosen—but likewise, no attempt was made to minimize forecast impact in AC07 by optimizing the ordering as in Kotsuki et al. (2017). The differences at this level reach up to  $3.5 \text{ g kg}^{-1}$ . The root-mean-squared difference of the entire domain at this level was approximately  $0.5 \text{ g kg}^{-1}$ . However, the MFN analyzed solutions with different orderings were found to be identical to within  $10^{-8}$ . A similar tolerance was found by comparing the MFN and EPS solutions.

Figure 3a shows the level 20 water vapor standard deviation (across the ensemble) of the prior ensemble perturbations  $\mathbf{X}'_f$ , while the standard deviation of the MFN posterior perturbations  $\mathbf{X}'_a$  with orderings 1 and 2 (which are numerically equivalent up to single precision) is shown in Fig. 3b. Figures 3c and 3d show the standard deviation of  $\mathbf{X}'_a$  at this level for orderings 1 and 2, respectively, with the AC07 filter. Figure 3e shows the two standard deviations differ by up to  $0.1 \text{ g kg}^{-1}$ , while the difference between the AC07 order 1  $\mathbf{X}'_a$  and the EPS solution is up to  $0.35 \text{ g kg}^{-1}$ . As in the mean, the MFN perturbations and the EPS perturbations are identical to within  $10^{-7}$ .

As shown in Figs. 1–3, the differences in the  $\bar{\mathbf{x}}_a$  analysis with random orderings using the AC07 filter are large enough that they are comparable to the posterior



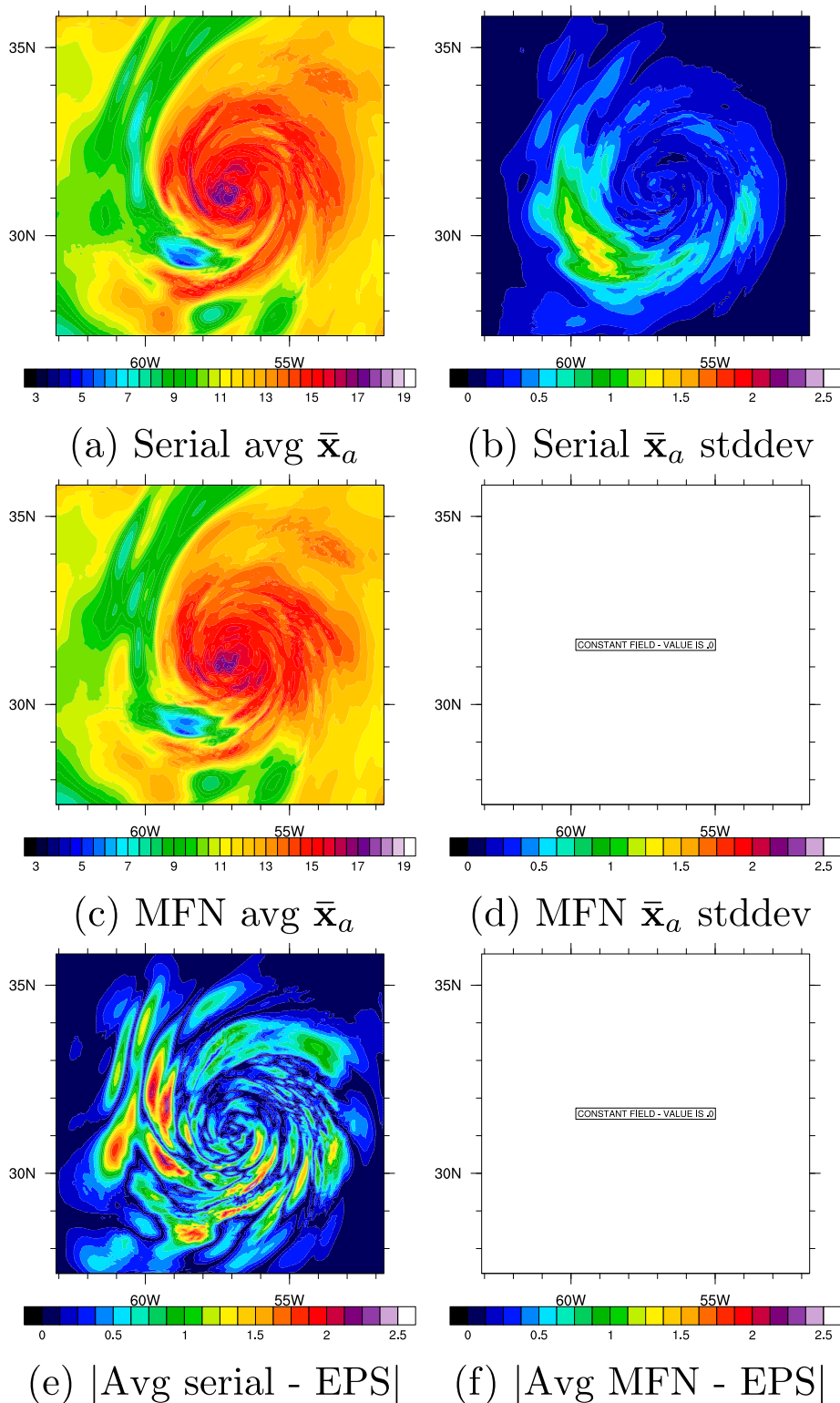


FIG. 1. Comparison between water vapor ( $\text{g kg}^{-1}$ ) at level 20 (of 60 total, corresponding to a height of approximately 2.5 km) of the Hurricane Edouard single cycle case of [Christophersen et al. \(2017\)](#) with 15 200 observations and  $L = 240$  as described in [S17](#). Ten different random orderings of observations were used. (a) The average of the 10 AC07  $\bar{x}_a$  analyses. (b) The standard deviation of

covariance in certain locations. This is likely due to the highly nonlinear nature of the first-cycle tropical cyclone data assimilation problem. In this application, flights are used as observing platforms to narrow the inner-core uncertainty as shown in Fig. 3. The first-cycle background contains ensemble members with simulated tropical cyclones with features centered at different locations, leading to large analysis updates. The main area of uncertainty in the AC07 analyses is actually outside of the inner core in the southwest quadrant near an area of dry air inflow. As shown, over the different serial courses of assimilation, the order-dependent error standard deviation of this region can grow to be roughly equivalent in magnitude to the posterior covariance. The matrix function approach, however, is order independent and therefore removes this source of error and is thus more numerically consistent with the eigenpair-based solution to the ESRF equations.

Having established that in this case the matrix function solution is numerically similar to the EPS method, which has a proven error bounds, we now turn our attention to the computational performance of the new method. For this purpose, we use a second experimental setup that combines the observations at all times that fall within the same domain as the first experiment. This leads to up to 35 420 quality-controlled observations that can be used for performance testing.

Keeping 1/2 of the total number of observations from all cycles fixed at 17 700, the scaling as a function of number of cores is shown in Fig. 4. The matrix function method scales nearly linearly as a function of the number of processing elements as in S17, but overall, the wall time remains bound by I/O time.

As a function of the number of observations, the MFN implementation scales much better than the EPS as shown in Fig. 5, where the number of processing elements is fixed at 386,  $L = 240$  for the correlation length scale, and the number of observations varies. As discussed in S17,  $L = 240$  leads to points across more than half of the domain being correlated, which in turn leads to a relatively dense, nearly full-rank matrix. As predicted by theory, the EPS solution appears to scale as the cube of the number of observations. However, the MFN approach apparently scales linearly. Times for the EPS solution longer than 45 min are not shown. With 17 700 observations on 386 processing elements, the EPS

solution took 41:28 to complete from start to finish (including expensive disk reading and writing), while the MFN solution took only 16:42. The MFN solution continues to scale well even at 35 400 observations, completing in 30:45, which is still more than 10 min faster than the EPS solution with half as many observations. Therefore, as shown, the MFN approach scales much better as function of the number of observations than the EPS solution.

The MFN solution is also roughly comparable to the AC07 solution in terms of wall time. While the MFN approach is actually slightly faster for small numbers of observations, for the largest number of observations tested (35 400 observations), the serial filter is faster with a wall time of 28:24, as opposed to 30:45. However, the wall-time differences are small enough that the observation order independence of MFN apparently makes it competitive with AC07 for these numbers of observations. This is somewhat surprising, as the only communication used by the AC07 filter is to broadcast observations, while distributed matrix multiplications are required by the MFN approach. However, the MFN approach has the potential benefit that it does not serially iterate over the observations, but instead can process all observations in parallel.

The number of matrix multiplications, and hence the overall timing of the matrix function solution, is directly related to the number of restarts and  $m$ , the maximum basis size before restarting. Increasing  $m$  leads to fewer restarts but requires additional memory and dense matrix processing time. The number of Eiermann–Ernst restarts necessary for convergence with  $m = 150$ , as used in our study, ranged from 1 for the smallest number of observations (2760) to 2 for the largest number of observations (35 420). The SLEPc error estimate at the end of each restart iteration for the smallest number of observations was on the order of  $10^{-2}$  for  $k = 0$  and  $10^{-15}$  for  $k = 1$ , while for the largest, the error was on the order of  $10^{-2}$  for  $k = 0$ ,  $10^{-8}$  for  $k = 1$ , and  $10^{-13}$  for  $k = 2$ . It appears the number of restarts grows very weakly with  $N_{\text{obs}}$ .

Table 1 shows the time necessary to solve the matrix function portion of the ESRF equations per ensemble member with  $L = 240$  for the 17 700-observation case as a function of varying the  $m$  parameter. As shown,  $m$  less than 100 requires an excessive number of restarts

---

←

these 10 AC07  $\bar{x}_a$  analyses. (c) The average of the 10 MFN  $\bar{x}_a$  analyses. (d) The standard deviation of the MFN analyses, which is less than  $10^{-7}$  at all points. (e) The absolute difference between (a) and the EPS solution. (f) The absolute difference between (c) and the EPS solution (also less than  $10^{-7}$  for all points).

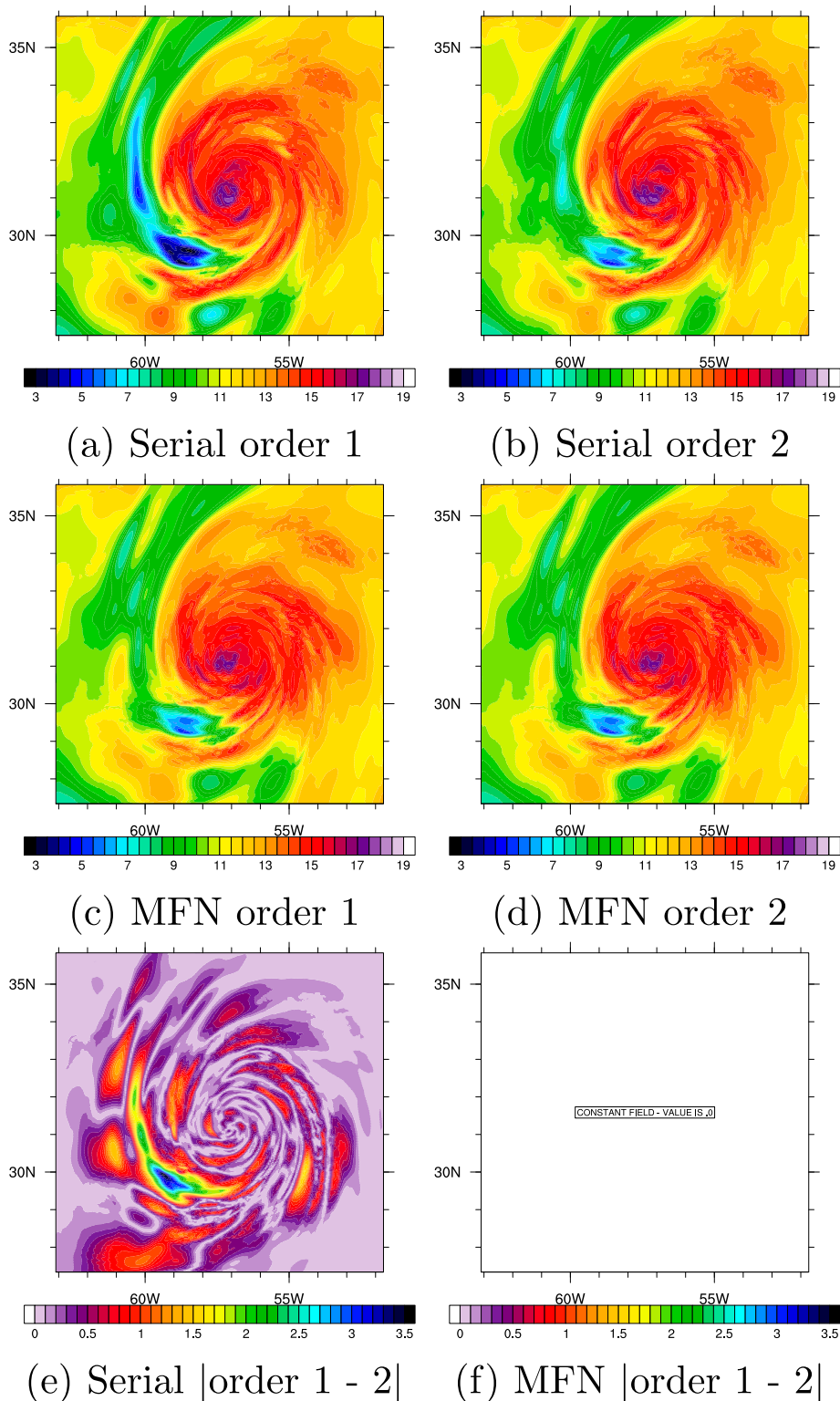


FIG. 2. Comparison between analyzed  $\bar{x}_q$  level 20 water vapor, as in Fig. 1, for two of the 10 different random orderings of observations. The serial filter of AC07 with (a) ordering 1 and (b) ordering 2. (c) MFN analyzed  $\bar{x}_q$ , ordering 1. (d) MFN  $\bar{x}_q$ , ordering 2. (e) The absolute value difference between (a) and (b). (f) The difference between the two MFN orderings in (c),(d), which is less than  $10^{-7}$ . The difference between the MFN and EPS analysis for this case is also less than  $10^{-7}$  at all levels.

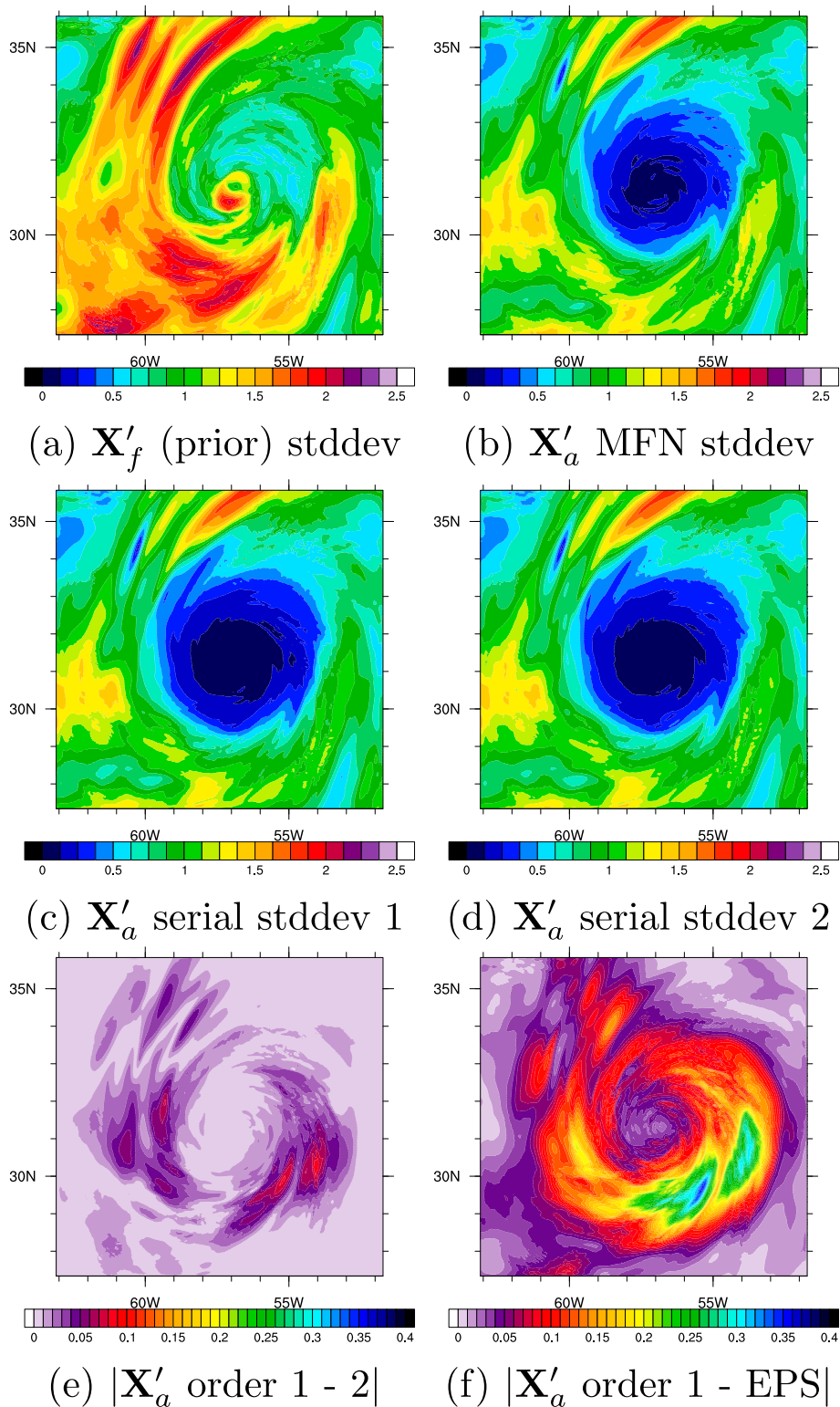


FIG. 3. Ensemble spread (i.e., standard deviations over the ensemble) of water vapor ( $\text{g kg}^{-1}$ ) at level 20, as in Figs. 1 and 2. Here, the first two random orderings of observations were used as in Fig. 2. (a) The standard deviation of the prior distribution  $\mathbf{X}'_f$  at this level. (b) The standard deviation of the MFN posterior distribution  $\mathbf{X}'_a$  (orderings 1, 2, and the EPS solution are the same to within  $10^{-7}$ ).

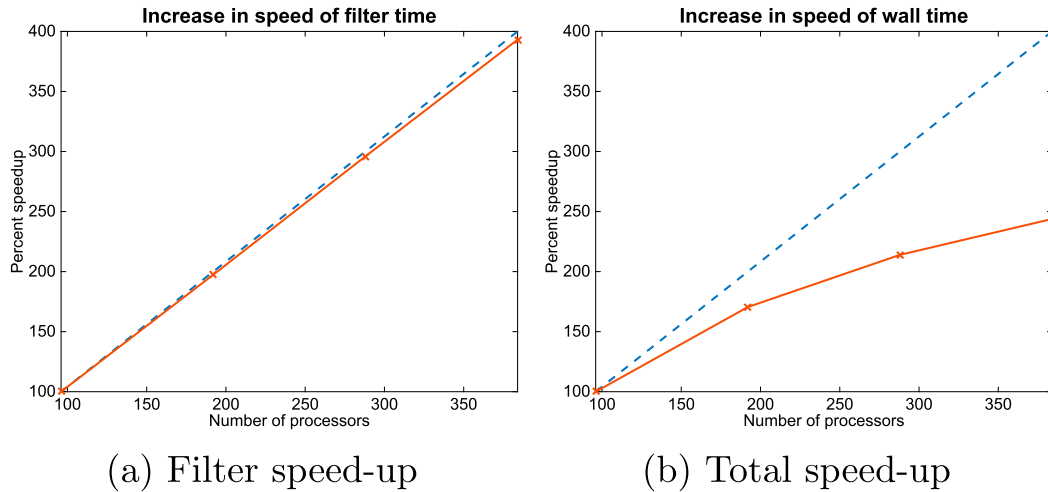


FIG. 4. (a) Speed increase of applying the MFN filter, including the time to calculate  $(E|g)$  using the MFN approach and solve Eq. (1), as a function of number of processing elements with the number of observations fixed at 17 000. The speedup is nearly linear and is dominated by the time applying  $C_{H_x, H_x}$ . This should be compared with Fig. 6d from S17, which likewise shows a nearly linear speed increase as a function of number of processors during filter time. (b) Total speed increase of wall time including disk reads and writes. As the process is I/O bound, the total speed increase is sublinear. Compare with Fig. 6f from S17, which likewise shows a sublinear increase (and even an eventual decrease) as a function of total wall time due to degradation in parallel I/O performance.

and total matrix product evaluations; for  $m$  greater than 100, the overall performance is dependent upon the exact number of matrix product evaluations required to reach the numerical accuracy of  $\epsilon_{\text{tol}} = 10^{-8}$ . For this case,  $m = 125$  requires the fewest matrix-product evaluations, which is highly correlated with the total MFN solve time. Table 2 shows the same results with the localization length scale  $L = 60$ . In this case,  $m = 150$  gives the optimal results. The best particular value of  $m$  therefore depends upon the factorization of the total number of evaluations required. An  $m$  larger than 100 is recommended to avoid excessive restarting, and an  $m$  less than 200 is recommended due to the expense of dense matrix evaluations. We choose  $m = 150$  to split the difference.

The scaling of memory usage on 386 xjet processors as a function of number of observations is shown in Fig. 6. As shown, and as expected by theory, the EPS solution memory usage scales cubically as a function of the number of observations. The serial filter of AC07 apparently scales linearly, as it only processes a single observation at once. The MFN solution, which currently stores the entire sparse  $C_{H_x, H_x}$  matrix in memory, scales

better than S17 but apparently worse than linearly. This is because with  $L = 240$ , the  $C_{H_x, H_x}$  matrix is relatively dense. For a dense matrix, the memory requirements would be quadratic, while for a sparse matrix, the memory requirements would be closer to linear. The memory scaling here is consistent with a factor somewhere in between quadratic and linear. Note, however, that the expense here is related to the representation of  $C_{H_x, H_x}$  and not directly to the MFN approach.

Indeed, the computational performance of the MFN method comes down to computing the matrix product. As mentioned, as only  $D_b$  is required in this method, it is not necessary to explicitly store the matrix  $D$  in memory. This so-called “matrix-free method” was implemented and tested successfully. As a first test, we used a simple implementation that brute-force recalculated the elements of  $C_{H_x, H_x}$  when required and avoided storing these elements in memory. While the memory usage decreased as expected, the time necessary to recompute the covariances made the method uncompetitive with the stored-in-memory matrix approach. The matrix-free implementation took 29:19 on 386 processors for 4500 observations versus just 5:40 with a stored matrix.

←

Standard deviation of  $X'_q$  for AC07 (c) ordering 1 and (d) ordering 2. (e) The absolute difference between the serial analysis with orderings 1 and 2 from (c),(d). (f) The absolute difference between the EPS solution and ordering 1 from (c).



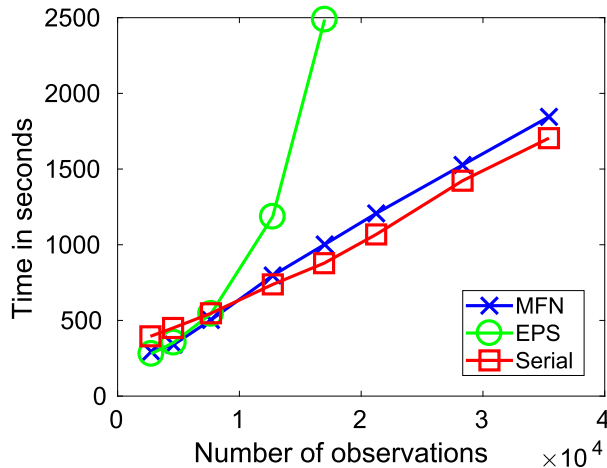


FIG. 5. Scaling as a function of number of observations with 386 processors. The MFN approach described in this paper appears to scale approximately linearly ( $y = 4.86 \times 10^{-2}x + 152$ ) while the EPS scales consistent with a cubic fit ( $y = 4.43 \times 10^{-10}x^3 + 302$ ). The serial filter of AC07 likewise scales linearly ( $y = 4.04 \times 10^{-2}x + 224.72$ ). Times longer than 2500 s for the EPS solution are not shown.

A more suitable matrix-free implementation, such as one based on FFT, would make this feature of the matrix function algorithm more attractive. Additional research is required in this area.

As an additional note, the MFN approach for solving the mean  $\mathbf{x} = f_1(\mathbf{D})[\mathbf{y} - \overline{H(\mathbf{X}_f)}]$  was compared with the more traditional method of solving for  $\mathbf{D}\mathbf{x} = \mathbf{y} - \overline{H(\mathbf{X}_f)}$  using GMRES. In this particular case, the MFN was found to be competitive with GMRES. This may be because  $\mathbf{D}$  is relatively dense, and an efficient preconditioner for use with GMRES was not found. Regardless, the novel contribution here is computing the more difficult  $f_2(\mathbf{D})(\mathbf{0} - H\mathbf{X})$  using MFN.

TABLE 1. Time to complete the solution, number of restarts (per control vector), and total number of matrix product evaluations as a function of  $m$ , the size of the Krylov subspace before restarting, required to solve the perturbation update matrix function  $f_2$  in Eq. (24) with  $L = 240$  [in Eqs. (5) and (8)] and 17 700 observations, as described in section 5. The timings are with a single MPI process on an Intel Core i7 server. Note these times are for a single ensemble member.

$m$	Time (s)	Restarts	Total evals
25	$2.5743 \times 10^4$	74	1875
50	$4.8219 \times 10^3$	12.2	660
75	$3.2298 \times 10^3$	5	450
100	$2.8686 \times 10^3$	3	400
125	$2.6897 \times 10^3$	2	375
150	$3.2460 \times 10^3$	2	450
175	$2.5257 \times 10^3$	1	350
200	$2.8950 \times 10^3$	1	400

TABLE 2. As in Table 1, but with  $L = 60$ . The reduction in time vs  $L = 240$  is due to the increased sparsity of the localization matrices  $\rho_{y,y}$  and  $\rho_{x,y}$ .

$m$	Time (s)	Restarts	Total evals
25	$1.6822 \times 10^4$	66.2	1705
50	$3.8763 \times 10^3$	15.0333	802
75	$2.7369 \times 10^3$	6.9	593
100	$2.2974 \times 10^3$	4	500
125	$2.2872 \times 10^3$	3	500
150	$2.0749 \times 10^3$	2	450
175	$2.4319 \times 10^3$	2	525
200	$2.2491 \times 10^3$	1.43 333	487

### 6. Discussion and conclusions

In this work, we describe the utilization of matrix functions, a powerful linear algebra tool, to derive numerically accurate and efficient solutions of the ESRF equations. With this method, high-rank localized covariance matrices can be applied consistently in such a way that the final analysis does not depend upon the ordering of observations. For the number of observations investigated, this method is roughly competitive in terms of wall time with the highly efficient serial filter of AC07.

The matrix function approach is built on the Arnoldi iteration, which provides a basis for the Krylov subspace spanned by the covariance matrix of the forward-computed observations  $\mathbf{C}_{H\mathbf{x},H\mathbf{x}}$  and a vector  $\mathbf{b}$ . This basis allows for evaluation of the ESRF matrix functions over a much smaller, upper-Hessenberg matrix. The Scalable Library for Eigenvalue Problem Computations

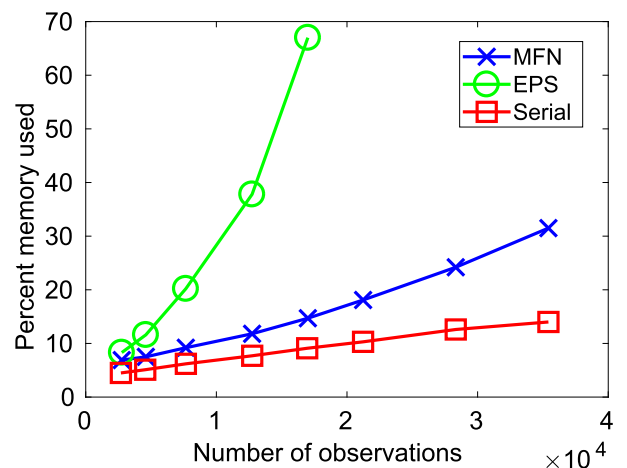


FIG. 6. Memory usage as a function of the number of observations with 386 processors. The EPS scales cubically as predicted by theory, while the serial filter of AC07 scales linearly in memory usage. The MFN approach apparently scales worse than linearly. A matrix-free implementation of MFN improves memory scaling.

(SLEPc; Hernandez et al. 2005) includes an efficient implementation of the matrix function method, along with the Eiermann–Ernst restart (Eiermann and Ernst 2006). Only the matrix-vector product is required, which can be used to provide matrix-free implementations; however, for performance reasons, storing the entire sparse  $\mathbf{C}_{H_x, H_x}$  matrix across processing elements may be preferable, as shown in our case. The ability to consistently incorporate high-rank covariance models with a known error bounds provides a platform to investigate hybrid ensemble–climatological covariances, as well as observation- versus model-space covariance issues.

Additional effort will be needed to fully understand the computational performance of this method in comparison to other existing parallel EnKF techniques, but a few basic conclusions can be drawn. First, in comparison to the eigenpair solution method of S17, the matrix function approach scales much better as a function of the number of observations assimilated and uses less memory while maintaining independence of observation ordering and achieving nearly identical numerical results. Second, while this method and the consistent hybrid ensemble filter (CHEF) of Bishop et al. (2015) are similarly independent of the order of observations for high-rank covariance models, as the matrix function approach applies the high-rank covariance matrices globally, it may be more computationally efficient than CHEF (which applies the matrices locally), especially for long localization lengths. This approach also solves the ESRF equations rather than using perturbed observations. Finally, the matrix function method is competitive with the serial implementation of AC07 in terms of wall time for the cases tested here. While it uses more memory, the matrix function approach is shown to be more faithful to the eigenpair-based solution of the ESRF equations than AC07. It is unknown if this additional precision will have a positive impact on forecasts. The recent work of Emanuel and Zhang (2017) demonstrates the crucial impact of inner-core moisture on TC predictability, and the two serial AC07 analyses shown in Fig. 2 with merely different observation orderings differ on the extent of dry air near the inner core. As shown, the two water vapor analyses for this difficult first-cycle TC case can differ by up to  $3 \text{ g kg}^{-1}$ , and therefore it is reasonable to expect the two serial analyses shown in Fig. 2 may produce qualitatively different medium-term forecasts. A method that can increase fidelity to the ESRF equations, known to be the minimum variance solution (e.g., Bishop et al. 2015), for tropical cyclone cases may be worth the additional computational expense. Because of the efficiency and ease of implementation of the serial filter, continued research into minimizing observation ordering impact is also likely to be beneficial.

Comparison of this method to other local analysis methods remains more unclear. The performance of local analysis methods is most critically related to the radius of influence. For large radii as considered here, this would likely make local analysis methods inefficient, as the problem for each local grid point becomes nearly as large as the entire domain. However, in such cases, when sample-based covariance localization is utilized with the ESRF approach, the matrix function approach could also potentially be used to improve performance versus  $O(n^3)$  algorithms, such as finding eigenpairs or the Cholesky decomposition. This may be unnecessary, however, if the number of local observations does not exceed  $\approx 10^2$ .

At the moment, a major weakness of the nonlocal matrix function approach in comparison to the AC07 serial approach is the memory usage scaling. Extrapolating the results presented in Fig. 6 on 386 processors and keeping the number of processors constant, with approximately 80 000 (assuming quadratic growth) to 115 000 observations (assuming linear growth), the matrix function approach would run out of memory. By comparison, the serial filter would run out of memory (assuming linear growth) at approximately 3.2 million observations. A matrix-free implementation would address this issue. Since in the matrix function approach, computational performance comes down efficient methods of applying the matrix product, we aim to investigate application of the modulation product of Bishop and Hodyss (2009) to apply correlations in order improve the memory scaling issue. In the meantime, batch processing of large numbers of observations is one potential workaround.

The algorithm described in this paper requires a distributed sparse matrix implementation, such as that available in the Portable Extensible Toolkit for Scientific Computing (PETSc; Balay et al. 1997, 2016, 2017), which SLEPc is built upon. In addition, the restarted Arnoldi process (including a numerically stable parallel Gram–Schmidt orthogonalization process) must be implemented to estimate the required reduced-order matrix function products. When using the SLEPc library that provides this functionality, this approach is not more difficult than the eigenpair implementation of S17. However, either implementation is certainly more complex than the serial approximation.

Finally, while the order-dependency issue shown here is nontrivial, the TC first-cycle case is likely to be a “worst case” scenario due to the highly nonlinear nature of feature misalignment. While Neger (2015) hypothesized that the effect of the observation-order dependency in the serial implementation is small when the analysis is not far from the prior, the filter described

here may be useful to test the practical effect of this hypothesis in a variety of large-scale cases and to develop mitigation solutions for the serial approach when necessary.

*Acknowledgments.* This research was partially funded by the NOAA Hurricane Forecast Improvement Project Award NA14NWS4680022. This work was partially supported by Agencia Estatal de Investigación (AEI) under Grant TIN2016-75985-P, which includes European Commission ERDF funds. Alejandro Lamas Daviña was supported by the Spanish Ministry of Education, Culture and Sport through a grant with reference FPU13-06655. The fourth author's work was in part carried out under the auspices of CIMAS, a joint institute of the University of Miami and NOAA, Cooperative Agreement NA15OAR4320064. The authors acknowledge the NOAA Research and Development High Performance Computing Program for providing computing and storage resources that have contributed to the research results reported within this paper (<http://rdhpcs.noaa.gov>). We thank Jeff Anderson, Shu-Chih Yang, and three anonymous reviewers for their helpful comments and contributions. We also thank Hui Christophersen for providing technical assistance.

#### REFERENCES

- Aberson, S. D., A. Aksoy, K. J. Sellwood, T. Vukicevic, and X. Zhang, 2015: Assimilation of high-resolution tropical cyclone observations with an ensemble Kalman filter using HEDAS: Evaluation of 2008–11 HWRf forecasts. *Mon. Wea. Rev.*, **143**, 511–523, <https://doi.org/10.1175/MWR-D-14-00138.1>.
- Afanasjew, M., M. Eiermann, O. G. Ernst, and S. Güttel, 2008: Implementation of a restarted Krylov subspace method for the evaluation of matrix functions. *Linear Algebra Appl.*, **429**, 2293–2314, <https://doi.org/10.1016/j.laa.2008.06.029>.
- Aksoy, A., 2013: Storm-relative observations in tropical cyclone data assimilation with an ensemble Kalman filter. *Mon. Wea. Rev.*, **141**, 506–522, <https://doi.org/10.1175/MWR-D-12-00094.1>.
- , S. Lorsolo, T. Vukicevic, K. J. Sellwood, S. D. Aberson, and F. Zhang, 2012: The HWRf Hurricane Ensemble Data Assimilation System (HEDAS) for high-resolution data: The impact of airborne Doppler radar observations in an OSSE. *Mon. Wea. Rev.*, **140**, 1843–1862, <https://doi.org/10.1175/MWR-D-11-00212.1>.
- , S. D. Aberson, T. Vukicevic, K. J. Sellwood, S. Lorsolo, and X. Zhang, 2013: Assimilation of high-resolution tropical cyclone observations with an ensemble Kalman filter using NOAA/AOML/HRD's HEDAS: Evaluation of the 2008–11 vortex-scale analyses. *Mon. Wea. Rev.*, **141**, 1842–1865, <https://doi.org/10.1175/MWR-D-12-00194.1>.
- Anderson, J. L., 2001: An ensemble adjustment Kalman filter for data assimilation. *Mon. Wea. Rev.*, **129**, 2884–2903, [https://doi.org/10.1175/1520-0493\(2001\)129<2884:AEAKFF>2.0.CO;2](https://doi.org/10.1175/1520-0493(2001)129<2884:AEAKFF>2.0.CO;2).
- , 2003: A local least squares framework for ensemble filtering. *Mon. Wea. Rev.*, **131**, 634–642, [https://doi.org/10.1175/1520-0493\(2003\)131<0634:ALLSFF>2.0.CO;2](https://doi.org/10.1175/1520-0493(2003)131<0634:ALLSFF>2.0.CO;2).
- , 2012: Localization and sampling error correction in ensemble Kalman filter data assimilation. *Mon. Wea. Rev.*, **140**, 2359–2371, <https://doi.org/10.1175/MWR-D-11-00013.1>.
- , and N. Collins, 2007: Scalable implementations of ensemble filter algorithms for data assimilation. *J. Atmos. Oceanic Technol.*, **24**, 1452–1463, <https://doi.org/10.1175/JTECH2049.1>.
- Arnoldi, W. E., 1951: The principle of minimized iterations in the solution of the matrix eigenvalue problem. *Q. Appl. Math.*, **9**, 17–29, <https://doi.org/10.1090/qam/42792>.
- Balay, S., W. D. Gropp, L. C. McInnes, and B. F. Smith, 1997: Efficient management of parallelism in object-oriented numerical software libraries. *Modern Software Tools for Scientific Computing*, E. Arge, A. M. Bruaset, and H. P. Langtangen, Eds., Springer, 163–202, [https://doi.org/10.1007/978-1-4612-1986-6\\_8](https://doi.org/10.1007/978-1-4612-1986-6_8).
- , and Coauthors, 2016: PETSc Users Manual. Argonne National Laboratory Tech. Rep. ANL-95/11, Revision 3.7, 272 pp., <http://www.mcs.anl.gov/petsc/petsc-current/docs/manual.pdf>.
- , and Coauthors, 2017: PETSc: Portable, Extensible Toolkit for Scientific Computation. PETSc/Tao, <https://www.mcs.anl.gov/petsc/>.
- Bessho, K., and Coauthors, 2016: An introduction to Himawari-8/9—Japan's new-generation geostationary meteorological satellites. *J. Meteor. Soc. Japan*, **94**, 151–183, <https://doi.org/10.2151/jmsj.2016-009>.
- Bishop, C. H., and D. Hodyss, 2009: Ensemble covariances adaptively localized with ECO-RAP. Part 2: A strategy for the atmosphere. *Tellus*, **61A**, 97–111, <https://doi.org/10.1111/j.1600-0870.2008.00372.x>.
- , B. J. Etherton, and S. J. Majumdar, 2001: Adaptive sampling with the ensemble transform Kalman filter. Part I: Theoretical aspects. *Mon. Wea. Rev.*, **129**, 420–436, [https://doi.org/10.1175/1520-0493\(2001\)129<0420:ASWTET>2.0.CO;2](https://doi.org/10.1175/1520-0493(2001)129<0420:ASWTET>2.0.CO;2).
- , B. Huang, and X. Wang, 2015: A nonvariational consistent hybrid ensemble filter. *Mon. Wea. Rev.*, **143**, 5073–5090, <https://doi.org/10.1175/MWR-D-14-00391.1>.
- Björck, Å., 1994: Numerics of Gram-Schmidt orthogonalization. *Linear Algebra Appl.*, **197–198**, 297–316, [https://doi.org/10.1016/0024-3795\(94\)90493-6](https://doi.org/10.1016/0024-3795(94)90493-6).
- Burgers, G., P. J. van Leeuwen, and G. Evensen, 1998: Analysis scheme in the ensemble Kalman filter. *Mon. Wea. Rev.*, **126**, 1719–1724, [https://doi.org/10.1175/1520-0493\(1998\)126<1719:ASITEK>2.0.CO;2](https://doi.org/10.1175/1520-0493(1998)126<1719:ASITEK>2.0.CO;2).
- Campbell, W. F., C. H. Bishop, and D. Hodyss, 2010: Vertical covariance localization for satellite radiances in ensemble Kalman filters. *Mon. Wea. Rev.*, **138**, 282–290, <https://doi.org/10.1175/2009MWR3017.1>.
- Chang, C.-C., S.-C. Yang, and C. Keppenne, 2014: Applications of the mean recentering scheme to improve typhoon track prediction: A case study of Typhoon Nanmadol (2011). *J. Meteor. Soc. Japan*, **92**, 559–584, <https://doi.org/10.2151/jmsj.2014-604>.
- Christophersen, H., A. Aksoy, J. Dunion, and K. Sellwood, 2017: The impact of NASA Global Hawk unmanned aircraft dropwindsonde observations on tropical cyclone track, intensity, and structure: Case studies. *Mon. Wea. Rev.*, **145**, 1817–1830, <https://doi.org/10.1175/MWR-D-16-0332.1>.
- Deadman, E., N. J. Higham, and R. Ralha, 2012: Blocked Schur algorithms for computing the matrix square root. *Applied Parallel and Scientific Computing*, P. Manninen and P. Öster, Eds., Lecture Notes in Computer Science Series, Vol. 7782, Springer, 171–182, [https://doi.org/10.1007/978-3-642-36803-5\\_12](https://doi.org/10.1007/978-3-642-36803-5_12).

- Eiermann, M., and O. Ernst, 2006: A restarted Krylov subspace method for the evaluation of matrix functions. *SIAM J. Numer. Anal.*, **44**, 2481–2504, <https://doi.org/10.1137/050633846>.
- , —, and S. Güttel, 2011: Deflated restarting for matrix functions. *SIAM J. Matrix Anal. Appl.*, **32**, 621–641, <https://doi.org/10.1137/090774665>.
- Emanuel, K., and F. Zhang, 2017: The role of inner-core moisture in tropical cyclone predictability and practical forecast skill. *J. Atmos. Sci.*, **74**, 2315–2324, <https://doi.org/10.1175/JAS-D-17-0008.1>.
- Evensen, G., 1994: Sequential data assimilation with a nonlinear quasi-geostrophic model using Monte Carlo methods to forecast error statistics. *J. Geophys. Res.*, **99**, 10143–10162, <https://doi.org/10.1029/94JC00572>.
- Frayssé, V., L. Giraud, and H. Kharraz-Aroussi, 1998: On the influence of the orthogonalization scheme on the parallel performance of GMRES. *Euro-Par'98 Parallel Processing*, D. Pritchard and J. Reeve, Eds., Lecture Notes in Computer Science Series, Vol. 1470, Springer, 751–762, <https://doi.org/10.1007/BFb0057927>.
- Frommer, A., K. Lund, M. Schweitzer, and D. Szyld, 2017: The Radau–Lanczos method for matrix functions. *SIAM J. Matrix Anal. Appl.*, **38**, 710–732, <https://doi.org/10.1137/16M1072565>.
- Gaspari, G., and S. E. Cohn, 1999: Construction of correlation functions in two and three dimensions. *Quart. J. Roy. Meteor. Soc.*, **125**, 723–757, <https://doi.org/10.1002/qj.49712555417>.
- Godinez, H. C., and J. D. Moulton, 2012: An efficient matrix-free algorithm for the ensemble Kalman filter. *Comput. Geosci.*, **16**, 565–575, <https://doi.org/10.1007/s10596-011-9268-9>.
- Golub, G. H., and C. F. Van Loan, 1996: *Matrix Computations*. Vol. 3. Johns Hopkins University Press, 694 pp.
- Gopalakrishnan, S. G., F. Marks, X. Zhang, J.-W. Bao, K.-S. Yeh, and R. Atlas, 2011: The experimental HWRF system: A study on the influence of horizontal resolution on the structure and intensity changes in tropical cyclones using an idealized framework. *Mon. Wea. Rev.*, **139**, 1762–1784, <https://doi.org/10.1175/2010MWR3535.1>.
- Hamill, T. M., J. S. Whitaker, and C. Snyder, 2001: Distance-dependent filtering of background error covariance estimates in an ensemble Kalman filter. *Mon. Wea. Rev.*, **129**, 2776–2790, [https://doi.org/10.1175/1520-0493\(2001\)129<2776:DDFOBE>2.0.CO;2](https://doi.org/10.1175/1520-0493(2001)129<2776:DDFOBE>2.0.CO;2).
- Hernandez, V., J. E. Roman, and V. Vidal, 2005: SLEPc: A scalable and flexible toolkit for the solution of eigenvalue problems. *ACM Trans. Math. Software*, **31**, 351–362, <https://doi.org/10.1145/1089014.1089019>.
- , —, and A. Tomas, 2007: Parallel Arnoldi eigensolvers with enhanced scalability via global communications rearrangement. *Parallel Comput.*, **33**, 521–540, <https://doi.org/10.1016/j.parco.2007.04.004>.
- Higham, N. J., 1987: Computing real square roots of a real matrix. *Linear Algebra Appl.*, **88–89**, 405–430, [https://doi.org/10.1016/0024-3795\(87\)90118-2](https://doi.org/10.1016/0024-3795(87)90118-2).
- Higham, N., 2008: *Functions of Matrices: Theory and Computation*. Other Titles in Applied Mathematics Series, Vol. 104, Society for Industrial and Applied Mathematics, 425 pp.
- Hochbruck, M., and C. Lubich, 1997: On Krylov subspace approximations to the matrix exponential operator. *SIAM J. Numer. Anal.*, **34**, 1911–1925, <https://doi.org/10.1137/S0036142995280572>.
- Houtekamer, P. L., and H. L. Mitchell, 1998: Data assimilation using an ensemble Kalman filter technique. *Mon. Wea. Rev.*, **126**, 796–811, [https://doi.org/10.1175/1520-0493\(1998\)126<0796:DAUAEK>2.0.CO;2](https://doi.org/10.1175/1520-0493(1998)126<0796:DAUAEK>2.0.CO;2).
- , and —, 2001: A sequential ensemble Kalman filter for atmospheric data assimilation. *Mon. Wea. Rev.*, **129**, 123–137, [https://doi.org/10.1175/1520-0493\(2001\)129<0123:ASEKFF>2.0.CO;2](https://doi.org/10.1175/1520-0493(2001)129<0123:ASEKFF>2.0.CO;2).
- , B. He, and H. L. Mitchell, 2014: Parallel implementation of an ensemble Kalman filter. *Mon. Wea. Rev.*, **142**, 1163–1182, <https://doi.org/10.1175/MWR-D-13-00011.1>.
- Hunt, B. R., E. J. Kostelich, and I. Szunyogh, 2007: Efficient data assimilation for spatiotemporal chaos: A local ensemble transform Kalman filter. *Physica D*, **230**, 112–126, <https://doi.org/10.1016/j.physd.2006.11.008>.
- Kalman, R. E., 1960: A new approach to linear filtering and prediction problems. *J. Basic Eng.*, **82**, 35–45, <https://doi.org/10.1115/1.3662552>.
- Keppenne, C. L., and M. M. Rienecker, 2002: Initial testing of a massively parallel ensemble Kalman filter with the Poseidon isopycnal ocean general circulation model. *Mon. Wea. Rev.*, **130**, 2951–2965, [https://doi.org/10.1175/1520-0493\(2002\)130<2951:ITOAMP>2.0.CO;2](https://doi.org/10.1175/1520-0493(2002)130<2951:ITOAMP>2.0.CO;2).
- Kotsuki, S., S. J. Greybush, and T. Miyoshi, 2017: Can we optimize the assimilation order in the serial ensemble Kalman filter? A study with the Lorenz-96 model. *Mon. Wea. Rev.*, **145**, 4977–4995, <https://doi.org/10.1175/MWR-D-17-0094.1>.
- Miyoshi, T., and Coauthors, 2016: “Big data assimilation” revolutionizing severe weather prediction. *Bull. Amer. Meteor. Soc.*, **97**, 1347–1354, <https://doi.org/10.1175/BAMS-D-15-00144.1>.
- Nerger, L., 2015: On serial observation processing in localized ensemble Kalman filters. *Mon. Wea. Rev.*, **143**, 1554–1567, <https://doi.org/10.1175/MWR-D-14-00182.1>.
- Niño-Ruiz, E. D., A. Sandu, and J. Anderson, 2015: An efficient implementation of the ensemble Kalman filter based on an iterative Sherman–Morrison formula. *Stat. Comput.*, **25**, 561–577, <https://doi.org/10.1007/s11222-014-9454-4>.
- , —, and X. Deng, 2018: A parallel implementation of the ensemble Kalman filter based on modified Cholesky decomposition. *J. Comput. Sci.*, <https://doi.org/10.1016/j.jocs.2017.04.005>, in press.
- Ott, E., and Coauthors, 2004: A local ensemble Kalman filter for atmospheric data assimilation. *Tellus*, **56A**, 415–428, <https://doi.org/10.1111/j.1600-0870.2004.00076.x>.
- Rogers, R. F., J. A. Zhang, J. Zawislak, H. Jiang, G. R. Alvey, E. J. Zipser, and S. N. Stevenson, 2016: Observations of the structure and evolution of Hurricane Edouard (2014) during intensity change. Part II: Kinematic structure and the distribution of deep convection. *Mon. Wea. Rev.*, **144**, 3355–3376, <https://doi.org/10.1175/MWR-D-16-0017.1>.
- Saad, Y., 1992: Analysis of some Krylov subspace approximations to the matrix exponential operator. *SIAM J. Numer. Anal.*, **29**, 209–228, <https://doi.org/10.1137/0729014>.
- Sakov, P., and L. Bertino, 2011: Relation between two common localisation methods for the EnKF. *Comput. Geosci.*, **15**, 225–237, <https://doi.org/10.1007/s10596-010-9202-6>.
- Schmit, T. J., P. Griffith, M. M. Gunshor, J. M. Daniels, S. J. Goodman, and W. J. Lebar, 2017: A closer look at the ABI on the GOES-R series. *Bull. Amer. Meteor. Soc.*, **98**, 681–698, <https://doi.org/10.1175/BAMS-D-15-00230.1>.
- Stewart, J. L., A. Aksoy, and Z. S. Haddad, 2017: Parallel direct solution of the ensemble square root Kalman filter equations with observation principal components. *J. Atmos. Oceanic Technol.*, **34**, 1867–1884, <https://doi.org/10.1175/JTECH-D-16-0140.1>.
- Tippett, M. K., J. L. Anderson, C. H. Bishop, T. M. Hamill, and J. S. Whitaker, 2003: Ensemble square root filters. *Mon. Wea. Rev.*, **131**, 1485–1490, [https://doi.org/10.1175/1520-0493\(2003\)131<1485:ESRF>2.0.CO;2](https://doi.org/10.1175/1520-0493(2003)131<1485:ESRF>2.0.CO;2).

- Van Der Vorst, H. A., 1987: An iterative solution method for solving  $f(A)x = b$ , using Krylov subspace information obtained for the symmetric positive definite matrix  $A$ . *J. Comput. Appl. Math.*, **18**, 249–263, [https://doi.org/10.1016/0377-0427\(87\)90020-3](https://doi.org/10.1016/0377-0427(87)90020-3).
- Vukicevic, T., A. Aksoy, P. Reasor, S. D. Aberson, K. J. Sellwood, and F. Marks, 2013: Joint impact of forecast tendency and state error biases in ensemble Kalman filter data assimilation of inner-core tropical cyclone observations. *Mon. Wea. Rev.*, **141**, 2992–3006, <https://doi.org/10.1175/MWR-D-12-00211.1>.
- Wang, Y., Y. Jung, T. A. Supinie, and M. Xue, 2013: A hybrid MPI–OpenMP parallel algorithm and performance analysis for an ensemble square root filter designed for multiscale observations. *J. Atmos. Oceanic Technol.*, **30**, 1382–1397, <https://doi.org/10.1175/JTECH-D-12-00165.1>.
- Whitaker, J. S., and T. M. Hamill, 2002: Ensemble data assimilation without perturbed observations. *Mon. Wea. Rev.*, **130**, 1913–1924, [https://doi.org/10.1175/1520-0493\(2002\)130<1913:EDAWPO>2.0.CO;2](https://doi.org/10.1175/1520-0493(2002)130<1913:EDAWPO>2.0.CO;2).
- Zawislak, J., H. Jiang, G. R. Alvey, E. J. Zipser, R. F. Rogers, J. A. Zhang, and S. N. Stevenson, 2016: Observations of the structure and evolution of Hurricane Edouard (2014) during intensity change. Part I: Relationship between the thermodynamic structure and precipitation. *Mon. Wea. Rev.*, **144**, 3333–3354, <https://doi.org/10.1175/MWR-D-16-0018.1>.
- Zhang, S., M. J. Harrison, A. T. Wittenberg, A. Rosati, J. L. Anderson, and V. Balaji, 2005: Initialization of an ENSO forecast system using a parallelized ensemble filter. *Mon. Wea. Rev.*, **133**, 3176–3201, <https://doi.org/10.1175/MWR3024.1>.