

1 **An Estimate of the Inflation Factor and Analysis Sensitivity**  
2 **in the Ensemble Kalman Filter**

3

4 Guocan Wu<sup>1,2</sup>

5

6

7 1 College of Global Change and Earth System Science, Beijing Normal University,

8 Beijing, China

9 2 Joint Center for Global Change Studies, Beijing, China

10

1 **Abstract**

2

3 The Ensemble Kalman Filter is a widely used ensemble-based assimilation  
4 method, which estimates the forecast error covariance matrix using a Monte Carlo  
5 approach that involves an ensemble of short-term forecasts. While the accuracy of the  
6 forecast error covariance matrix is crucial for achieving accurate forecasts, the  
7 estimate given by the EnKF needs to be improved using inflation techniques.  
8 Otherwise, the sampling covariance matrix of perturbed forecast states will  
9 underestimate the true forecast error covariance matrix because of the limited  
10 ensemble size and large model errors, which may eventually result in the divergence  
11 of the filter.

12 In this study, the forecast error covariance inflation factor is estimated using a  
13 generalized cross-validation technique. The improved EnKF assimilation scheme is  
14 tested with the atmosphere-like Lorenz-96 model with spatially correlated  
15 observations, and is shown to reduce the analysis error and increase its sensitivity to  
16 the observations.

17 **Key words:** data assimilation; ensemble Kalman filter; forecast error inflation;  
18 analysis sensitivity; cross validation

19

# 1. Introduction

2

3 For state variables in geophysical research fields, a common assumption is that  
4 systems have a “true” underlying state. Data assimilation is a powerful mechanism  
5 for estimating the true trajectory based on the effective combination of a dynamic  
6 forecast system (such as a numerical model) and observations (Miller et al. 1994).  
7 Data assimilations provide an analysis state that is usually a better estimate of the  
8 state variable because it fully considers all of the information provided by the model  
9 forecasts and observations. In fact, the analysis state can generally be treated as the  
10 weighted average of the model forecasts and observations, while the weights are  
11 approximately proportional to the inverse of the corresponding covariance matrices  
12 (Talagrand 1997). Therefore, the performance of a data assimilation method relies  
13 significantly on whether the error covariance matrices are estimated accurately. If this  
14 is the case, the assimilation can be attributed to technical aspect and can be  
15 accomplished with the rapid development of supercomputers (Reichle 2008),  
16 although finding the global minimum is a much more difficult problem when the  
17 models are nonlinear.

18 The ensemble Kalman filter (EnKF) is a practical ensemble-based assimilation  
19 scheme that estimates the forecast error covariance matrix using a Monte Carlo  
20 method with the short-term ensemble forecast states (Burgers et al. 1998; Evensen  
21 1994). Because of the limited ensemble size and large model errors, the sampling  
22 covariance matrix of the ensemble forecast states usually underestimates the true

1 forecast error covariance matrix. This finding indicates that the filter is over reliant on  
2 the model forecasts and excludes the observations, and can eventually result in the  
3 divergence of the filter (Anderson; Anderson 1999; Constantinescu et al. 2007; Wu et  
4 al. 2014).

5         The covariance inflation technique is used to mitigate filter divergence by  
6 inflating the empirical covariance in EnKF, and it can increase the weight of the  
7 observations in the analysis state (Xu et al. 2013). In reality, this method will perturb  
8 the subspace spanned by the ensemble vectors and better capture the sub-growing  
9 directions that may be missed in the original ensemble (Yang et al. 2015). Therefore,  
10 using the inflation technique to enhance the estimate accuracy of the forecast error  
11 covariance matrix is increasingly important.

12         In early studies on forecast error inflation, researchers usually tuned the inflation  
13 factor by repeated assimilation experiments and selected the estimated inflation factor  
14 according to their experience and prior knowledge (Anderson; Anderson 1999).  
15 However, such methods are very empirical and subjective. In later studies, the  
16 inflation factor can be estimated online based on the innovation statistic  
17 (observation-minus-forecast; (Dee 1995; Dee; Silva 1999)) with different conditions.  
18 The moment estimation can facilitate the calculation by solving an equation of the  
19 innovation statistic and its realization (Li et al. 2009; Miyoshi 2011; Wang; Bishop  
20 2003). The maximum likelihood approach can obtain a better estimate of the inflation  
21 factor, although it must calculate a high dimensional matrix determinant (Liang et al.  
22 2012; Zheng 2009). The Bayesian approach assumes a prior distribution for the

1 inflation factor but is limited by spatially independent observational errors (Anderson  
2 2007, 2009). This study seeks to address the estimation of the inflation factor from  
3 the perspective of cross validation (CV).

4 The concept of CV was first introduced for linear regressions (Allen 1974) and  
5 spline smoothing (Wahba; Wold 1975), and it represents a common approach that can  
6 be applied to estimate tuning parameters in generalized additive models,  
7 nonparametric regressions and kernel smoothing (Eubank 1999; Gentle et al. 2004;  
8 Green; Silverman. 1994; Wand; Jones 1995). In CV, the data are divided into subsets  
9 some of which are used for modelling and analysis while others for verification and  
10 validation. The most widely used technique removes only one data point and uses the  
11 remainder to estimate the value at this point to test the estimation accuracy. This most  
12 commonly used form is also called the leave-one-out cross validation (Gu; Wahba  
13 1991).

14 The basic motivation behind CV is to minimize the prediction error at the  
15 sampling points. The generalised cross validation (GCV) is a modified form of  
16 ordinary CV, that has been found to possess several favourable properties and is more  
17 popular for selecting tuning parameters (Craven; Wahba 1979). For instance, Gu and  
18 Wahba applied the Newton method to optimize the GCV score with multiple  
19 smoothing parameters in a smoothing spline model (Gu; Wahba 1991). Wahba briefly  
20 reviewed the properties of the GCV and conducted an experiment to choose  
21 smoothing parameters in the context of variational data assimilation schemes with  
22 numerical weather prediction models (Wahba et al. 1995). Zheng and Basher also

1 used a GCV in a thin-plate smoothing spline model of spatial climate data and  
2 applied it to South Pacific rainfalls (Zheng; Basher 1995). The GCV criterion has a  
3 rotation-invariant property that is relative to the orthogonal transformation of the  
4 observations and is a consistent estimate of the relative loss (Gu 2002).

5 In the covariance inflation scheme, the forecast error matrix is multiplied by an  
6 appropriate inflation factor. Usually, the inflation factor is larger than 1. Too small or  
7 too large an inflation factor will cause the analysis state to be over reliant on model  
8 forecasts or observations. Hence, the inflation factor should be estimated accurately.  
9 After this, the weights of the model forecasts and observations in the analysis state  
10 can be reassigned. Therefore, the analysis sensitivity was also investigated in this  
11 study. Generally speaking, analysis sensitivity is used to apportion uncertainty in the  
12 output can to different sources of uncertainty in the input (Saltelli et al. 2004; Saltelli  
13 et al. 2008). In the context of statistical data assimilation, this quantity describes the  
14 sensitivity of the analysis to the observations, which is complementary to the  
15 sensitivity of the analysis to model forecasts (Cardinali et al. 2004; Liu et al. 2009)..

16 This study focuses on a methodology that can be potentially applied to  
17 geophysical applications of data assimilation in the near future. This paper consists of  
18 four sections. The conventional EnKF scheme is summarized and the improved EnKF  
19 with a forecast error inflation scheme is proposed in Section 2; the verification and  
20 validation processes are conducted on an idealized model in Section 3; and the  
21 discussion and conclusions are given in Section 4.

22

1

2 **2. Methodology**

3

4 **2.1. EnKF algorithm**

5 For consistency, a nonlinear discrete-time dynamical forecast model and linear

6 observation system can be expressed as follows (Ide et al. 1997):

7 
$$\mathbf{x}_i^t = M_{i-1}(\mathbf{x}_{i-1}^a) + \boldsymbol{\eta}_i, \quad (1)$$

8 
$$\mathbf{y}_i^o = \mathbf{H}_i \mathbf{x}_i^t + \boldsymbol{\varepsilon}_i, \quad (2)$$

9 where  $i$  represents the time index;  $\mathbf{x}_i^t = \{x_{i,1}^t, x_{i,2}^t, \dots, x_{i,n}^t\}^T$  represents the10  $n$ -dimensional true state vector at the  $i$ -th time step;  $\mathbf{x}_{i-1}^a = \{x_{i-1,1}^a, x_{i-1,2}^a, \dots, x_{i-1,n}^a\}^T$ 11 represents the  $n$ -dimensional analysis state vector, which is an estimate of  $\mathbf{x}_{i-1}^t$ ;  $M_{i-1}$ 

12 represents a nonlinear dynamical forecast operator such as a numeric weather

13 prediction model;  $\mathbf{y}_i^o = \{y_{i,1}^o, y_{i,2}^o, \dots, y_{i,p_i}^o\}^T$  represents a  $p_i$ -dimensional observation14 vector;  $\mathbf{H}_i$  represents the observation operator matrix; and  $\boldsymbol{\eta}_i$  and  $\boldsymbol{\varepsilon}_i$  represent

15 the forecast and observation error vectors, which are assumed to be time-uncorrelated,

16 statistically independent of each other and have mean zero and covariance matrices

17  $\mathbf{P}_i$  and  $\mathbf{R}_i$ , respectively. The EnKF assimilation result is a series of analysis states18  $\mathbf{x}_i^a$  that is an accurate estimate of the corresponding true states  $\mathbf{x}_i^t$  based on the19 information provided by  $M_i$  and  $\mathbf{y}_i^o$ .20 Supposing the perturbed analysis state at a previous time step  $\mathbf{x}_{i-1}^{a(j)}$  has been21 estimated ( $1 \leq j \leq m$  and  $m$  is the ensemble size), the detailed EnKF assimilation

22 procedure is summarized as the following forecast step and analysis step (Burgers et

1 al. 1998; Evensen 1994).

2 Step 1. Forecast step.

3 The perturbed forecast states are generated by dynamical model forecast  
4 forward:

$$5 \quad \mathbf{x}_i^{f(j)} = M_{i-1}(\mathbf{x}_{i-1}^{a(j)}). \quad (3)$$

6 The forecast state  $\mathbf{x}_i^f$  is defined as the ensemble mean of  $\mathbf{x}_i^{f(j)}$ , and the forecast error  
7 covariance matrix is initially estimated as the sampling covariance matrix of  
8 perturbed forecast states:

$$9 \quad \mathbf{P}_i = \frac{1}{m-1} \sum_{j=1}^m (\mathbf{x}_i^{f(j)} - \mathbf{x}_i^f)(\mathbf{x}_i^{f(j)} - \mathbf{x}_i^f)^T. \quad (4)$$

10 Step 2. Analysis step.

11 The analysis state is estimated by minimizing the following cost function

$$12 \quad J(\mathbf{x}) = (\mathbf{x} - \mathbf{x}_i^f)^T \mathbf{P}_i^{-1} (\mathbf{x} - \mathbf{x}_i^f) + (\mathbf{y}_i^o - \mathbf{H}_i \mathbf{x})^T \mathbf{R}_i^{-1} (\mathbf{y}_i^o - \mathbf{H}_i \mathbf{x}), \quad (5)$$

13 which has the analytic form

$$14 \quad \mathbf{x}_i^a = \mathbf{x}_i^f + \mathbf{P}_i \mathbf{H}_i^T (\mathbf{H}_i \mathbf{P}_i \mathbf{H}_i^T + \mathbf{R}_i)^{-1} \mathbf{d}_i, \quad (6)$$

15 where

$$16 \quad \mathbf{d}_i = \mathbf{y}_i^o - \mathbf{H}_i \mathbf{x}_i^f, \quad (7)$$

17 is the innovation statistic (observation-minus-forecast residual). To complete the  
18 ensemble forecast, the perturbed analysis states are calculated using perturbed  
19 observations (Burgers et al. 1998):

$$20 \quad \mathbf{x}_i^{a(j)} = \mathbf{x}_i^f + \mathbf{P}_i \mathbf{H}_i^T (\mathbf{H}_i \mathbf{P}_i \mathbf{H}_i^T + \mathbf{R}_i)^{-1} (\mathbf{d}_i + \boldsymbol{\varepsilon}_i^{(j)}), \quad (8)$$

21 where  $\boldsymbol{\varepsilon}_i^{(j)}$  is a normally distributed random variable with mean zero and covariance



1 matrix  $\mathbf{R}_i$ . Here,  $(\mathbf{H}_i \mathbf{P}_i \mathbf{H}_i^T + \mathbf{R}_i)^{-1}$  can be easily calculated using the  
 2 Sherman-Morrison-Woodbury formula (Liang et al. 2012; Tippett et al. 2003).  
 3 Finally, set  $i = i + 1$  and return to Step 1 for the model forecast at the next time step  
 4 and repeat until the model reaches the last time step  $N$ .

## 6 **2.2. Influence matrix and forecast error inflation**

7 The forecast error inflation scheme should be included in any ensemble-based  
 8 assimilation scheme or else to prevent the filter from diverging (Anderson; Anderson  
 9 1999; Constantinescu et al. 2007). Multiplicative inflation is one of the commonly  
 10 used inflation techniques, and it adjusts the initially estimated forecast error  
 11 covariance matrix  $\mathbf{P}_i$  to  $\lambda_i \mathbf{P}_i$  by estimating the inflation factors  $\lambda_i$  properly.

12 In previous studies, a number of methods were used to estimate the inflation  
 13 factor, such as the maximum likelihood approach (Liang et al. 2012; Zheng 2009),  
 14 moment approach (Li et al. 2009; Miyoshi 2011; Wang; Bishop 2003) and Bayesian  
 15 approach (Anderson 2007, 2009). In this study, a new procedure for estimating  
 16 multiplicative inflation factors  $\lambda_i$  is proposed based on the following GCV function  
 17 (Craven; Wahba 1979)

$$18 \quad GCV_i(\lambda) = \frac{\frac{1}{p_i} \mathbf{d}_i^T \mathbf{R}_i^{-1/2} (\mathbf{I}_{p_i} - \mathbf{A}_i(\lambda))^2 \mathbf{R}_i^{-1/2} \mathbf{d}_i}{\left[ \frac{1}{p_i} \text{Tr}(\mathbf{I}_{p_i} - \mathbf{A}_i(\lambda)) \right]^2}, \quad (9)$$

19 where  $\mathbf{I}_{p_i}$  is the identity matrix with dimension  $p_i \times p_i$ ;  $\mathbf{R}_i^{-1/2}$  is the square root  
 20 matrix of  $\mathbf{R}_i$ ; and

$$1 \quad \mathbf{A}_i(\lambda) = \mathbf{I}_{p_i} - \mathbf{R}_i^{1/2} \left( \mathbf{H}_i \lambda \mathbf{P}_i \mathbf{H}_i^T + \mathbf{R}_i \right)^{-1} \mathbf{R}_i^{1/2} \quad (10)$$

2 is the influence matrix (see Appendix for details).

3 The inflation factor  $\lambda_i$  is estimated by minimizing the GCV (Eq. (9)) as an  
 4 objective function, and it is implemented between Steps 1 and 2 in Section 2.1. Then,  
 5 the perturbed analysis states are modified to

$$6 \quad \mathbf{x}_i^{a(i)} = \mathbf{x}_i^{f(i)} + \lambda_i \mathbf{P}_i \mathbf{H}_i^T \left( \mathbf{H}_i \lambda_i \mathbf{P}_i \mathbf{H}_i^T + \mathbf{R}_i \right)^{-1} \left( \mathbf{d}_i + \boldsymbol{\varepsilon}_i^{(i)} \right). \quad (11)$$

7 The flowchart of the EnKF equipped with the forecast error inflation based on the  
 8 GCV method is shown in Figure 1.

9

### 10 **2.3. Analysis sensitivity**

11 In the EnKF, the analysis state (Eq. (6)) is a weighted average of the observation  
 12 and forecast. That is:

$$13 \quad \mathbf{x}_i^a = \mathbf{K}_i \mathbf{y}_i^o + \left( \mathbf{I}_n - \mathbf{K}_i \mathbf{H}_i \right) \mathbf{x}_i^f \quad (12)$$

14 where  $\mathbf{K}_i = \mathbf{P}_i \mathbf{H}_i^T \left( \mathbf{H}_i \mathbf{P}_i \mathbf{H}_i^T + \mathbf{R}_i \right)^{-1}$  is the Kalman gain matrix and  $\mathbf{I}_n$  is the identity  
 15 matrix with dimension  $n \times n$ . Then, the normalized analysis vector can be expressed  
 16 as follows:

$$17 \quad \tilde{\mathbf{y}}_i^a = \mathbf{R}_i^{-1/2} \mathbf{H}_i \mathbf{K}_i \mathbf{R}_i^{1/2} \tilde{\mathbf{y}}_i^o + \mathbf{R}_i^{-1/2} \left( \mathbf{I}_{p_i} - \mathbf{H}_i \mathbf{K}_i \right) \mathbf{R}_i^{1/2} \tilde{\mathbf{y}}_i^f \quad (13)$$

18 where  $\tilde{\mathbf{y}}_i^f = \mathbf{R}_i^{-1/2} \mathbf{H}_i \mathbf{x}_i^f$  is the normalized projection of the forecast on the  
 19 observation space. The sensitivities of the analysis to the observation and forecast are  
 20 defined by Eq. (14) and (15), respectively, as follows:

$$21 \quad \mathbf{S}_i^o = \frac{\partial \tilde{\mathbf{y}}_i^a}{\partial \tilde{\mathbf{y}}_i^o} = \mathbf{R}_i^{1/2} \mathbf{K}_i^T \mathbf{H}_i^T \mathbf{R}_i^{-1/2}, \quad (14)$$

$$1 \quad \mathbf{S}_i^f = \frac{\partial \tilde{\mathbf{y}}_i^a}{\partial \tilde{\mathbf{y}}_i^f} = \mathbf{R}_i^{1/2} (\mathbf{I}_{p_i} - \mathbf{K}_i^T \mathbf{H}_i^T) \mathbf{R}_i^{-1/2}, \quad (15)$$

2 which satisfy  $\mathbf{S}_i^o + \mathbf{S}_i^f = \mathbf{I}_{p_i}$ .

3 The elements of the matrix  $\mathbf{S}_i^o$  reflect the sensitivity of the normalized analysis  
 4 state to the normalized observations; its diagonal elements are the analysis  
 5 self-sensitivities and the off-diagonal elements are the cross-sensitivities. On the  
 6 other hand, the elements of the matrix  $\mathbf{S}_i^f$  reflect the sensitivity of the normalized  
 7 analysis state to the normalized forecast vector. The two quantities are  
 8 complementary, and the GCV function can be interpreted as minimizing the  
 9 normalized forecast sensitivity because the inflation scheme will increase the  
 10 observation weight appropriately.

11 In fact, the sensitivity matrix  $\mathbf{S}_i^o$  is equal to the influence matrix  $\mathbf{A}_i$  (see  
 12 Appendix B for detailed proof), whose trace can be used to measure the “equivalent  
 13 number of parameters” or “degrees of freedom for the signal” (Gu 2002; Pena; Yohai  
 14 1991). Similarly, the sensitivity matrix  $\mathbf{S}_i^o$  can be interpreted as a measurement of  
 15 the amount of information extracted from the observations (Ellison et al. 2009). Trace  
 16 diagnostic can be used to analyse the sensitivities to observations or forecast vectors  
 17 (Cardinali et al. 2004). The Global Average Influence (GAI) at the  $i$ -th time step is  
 18 defined as the globally averaged observation influence:

$$19 \quad GAI = \frac{\text{Tr}(\mathbf{S}_i^o)}{p_i} \quad (16)$$

20 where  $p_i$  is the total number of observations at the  $i$ -th time step.

21 In the conventional EnKF, the forecast error covariance matrix  $\mathbf{P}_i$  is initially

1 estimated using a Monte Carlo method with short-term ensemble forecast states.  
 2 However, because of the limited ensemble size and large model errors, the sampling  
 3 covariance matrix of perturbed forecast states usually underestimate the true forecast  
 4 error covariance matrix. This will cause the analysis to over rely on the forecast state,  
 5 excluding useful information from the observations. This is captured by the fact that  
 6 for the conventional EnKF scheme the GAI values are rather small. Adjusting the  
 7 inflation of the forecast error covariance matrix alleviates this problem to some extent,  
 8 as will be shown in the following simulations.

9

#### 10 ***2.4 Forecast ensemble spread and analysis RMSE***

11 The spread of the forecast ensemble at the  $i$ -th step is defined as follows:

$$12 \quad \text{Spread} = \sqrt{\frac{1}{n(m-1)} \sum_{j=1}^m \|\mathbf{x}_{i,j}^f - \mathbf{x}_i^f\|^2}. \quad (17)$$

13 Roughly speaking, the forecast ensemble spread is usually underestimated in  
 14 the conventional EnKF, which also dramatically decreases until the observations  
 15 ultimately have an irrelevant impact on the analysis states. Applying the inflation  
 16 technique, an underestimation of the forecast ensemble spread can be effectively  
 17 compensated for, thereby improving the assimilation results.

18 In the following experiments, the “true” state  $\mathbf{x}_i^t$  is non-dimensional and can  
 19 be obtained by a numerical solution of partial differential equations. In this case, the  
 20 distance of the analysis state to the true state can be defined as the analysis  
 21 root-mean-square error (RMSE), which is used to evaluate the accuracy of the  
 22 assimilation results. The RMSE at the  $i$ -th time step is defined as follows:

1 
$$\text{RMSE} = \sqrt{\frac{1}{n} \sum_{k=1}^n (x_{i,k}^a - x_{i,k}^t)^2} . \quad (18)$$

2 where  $x_{i,k}^a$  and  $x_{i,k}^t$  are the  $k$ -th components of the analysis state and true state at  
 3 the  $i$ -th time step. In principle, a smaller RMSE indicates a better performance of the  
 4 assimilation scheme.

5

6

### 7 **3. Numerical Experiments**

8

9 The proposed data assimilation scheme was tested using the Lorenz-96 model  
 10 (Lorenz 1996) with model errors and a linear observation system as a test bed. The  
 11 performances of the assimilation schemes described in Section 2 were evaluated via  
 12 the following experiments.

13

#### 14 **3.1. Dynamical forecast model and observation systems**

15 The Lorenz-96 model (Lorenz 1996) is a quadratic nonlinear dynamical system  
 16 that has properties relevant to realistic forecast problems and is governed by the  
 17 equation:

18 
$$\frac{d\mathbf{X}_k}{dt} = (\mathbf{X}_{k+1} - \mathbf{X}_{k-2})\mathbf{X}_{k-1} - \mathbf{X}_k + F , \quad (19)$$

19 where  $k = 1, 2, \dots, 40$ . The cyclic boundary conditions  $\mathbf{X}_{-1} = \mathbf{X}_{K-1}$ ,  $\mathbf{X}_0 = \mathbf{X}_K$ , and  
 20  $\mathbf{X}_{K+1} = \mathbf{X}_1$  were applied to ensure that Eq. (19) was well defined for all values of  $k$ .

21 The Lorenz-96 model is “atmosphere-like” because the three terms on the right-hand

1 side of Eq. (19) are analogous to a nonlinear advection-like term, a damping term,  
2 and an external forcing term. The model can be considered representative of an  
3 atmospheric quantity (e.g., zonal wind speed) distributed on a latitude circle.  
4 Therefore, the Lorenz-96 model has been widely used as a test bed to evaluate the  
5 performance of assimilation schemes in many studies (Wu et al. 2013).

6 The true state is derived by a fourth-order Runge-Kutta time integration scheme  
7 (Butcher 2003). The time step for generating the numerical solution was set at 0.05  
8 non-dimensional units, which is roughly equivalent to 6 hours in real time assuming  
9 that the characteristic time-scale of the dissipation in the atmosphere is 5 days  
10 (Lorenz 1996). The forcing term was set as  $F = 8$ , so that the leading Lyapunov  
11 exponent implies an error-doubling time of approximately 8 time steps and the fractal  
12 dimension of the attractor was 27.1 (Lorenz; Emanuel 1998). The initial value was  
13 chosen to be  $X_k = F$  when  $k \neq 20$  and  $X_{20} = 1.001F$ .

14 In this study, the synthetic observations were assumed to be generated at all of  
15 the 40 model grids by adding random noises that were multivariate-normally  
16 distributed with mean zero and covariance matrix  $\mathbf{R}_i$  to the true states. The  
17 frequency was set as every 4 time steps, which can be used to mimic daily  
18 observations in practical problems, such as satellite data. The observation errors were  
19 assumed to be spatially correlated, which is common in applications involving remote  
20 sensing and radiance data. The variance of the observation on each grid point was set  
21 to  $\sigma_o^2 = 1$ , and the covariance of the observations between the  $j$ -th and  $k$ -th grid  
22 points was as follows:

$$\mathbf{R}_i(j, k) = \sigma_o^2 \times 0.5^{\min\{|j-k|, 40-|j-k|\}}. \quad (20)$$

### 3.2. Assimilation scheme comparison

Because model errors are inevitable in practical dynamical forecast models, it is reasonable for us to add model error to the Lorenz-96 model in the assimilation process. The Lorenz-96 model is a forced dissipative model with a parameter  $F$  that controls the strength of the forcing (Eq. (19)). Modifying the forcing strength  $F$  changes the model forecast considerably. For values of  $F$  that are larger than 3, the system is chaotic (Lorenz; Emanuel 1998). To simulate model errors, the forcing term for the forecast was set to 7, while using  $F=8$  to generate the “true” state. The initially selected ensemble size was 30.

The Lorenz-96 model was run for 2000 time steps, which is equivalent to approximately 500 days in realistic problems. The synthetic observations were assimilated every 4 time steps using the conventional EnKF and the improved EnKF with forecast error inflation. The time series of estimated inflation factors are shown in Figure 2, which vary between 1 and 6 with greatly majority. The median was 1.88, which was used in the following comparison of the improved EnKF and the simple multiplicative inflation techniques like setting a constant inflation factor.

The forecast ensemble spread of the conventional EnKF, improved EnKF and constant inflated EnKF are plotted in Figure 3. For the conventional EnKF, because the forecast states usually shrink together, the forecast ensemble spread was quite small and had a mean value of 0.36. The mean spread value of the improved EnKF was 3.32, which was slightly larger than that of the constant inflated EnKF (3.25).

1 These findings illustrate that the underestimation of forecast ensemble spread can be  
2 effectively compensated for by the two EnKF schemes with forecast error inflation  
3 and that the improved EnKF is more effective than the constant inflated EnKF.

4 To evaluate the analysis sensitivity, the GAI statistics (Eq. (16)) were calculated,  
5 and the results are plotted in Figure 4. The increase in GAI from 10% for the  
6 conventional EnKF to 30% for the EnKF with forecast error inflation indicates that  
7 the latter relies more on the observations. This finding is important because the  
8 observations can play a more significant role in combining the results with the model  
9 forecasts to generate the analysis state. In addition to small fluctuations, the mean  
10 GAI value of the constant inflated EnKF was 27.80%, which was smaller than that of  
11 the improved EnKF.

12 To evaluate the resulting estimate, the analysis RMSE (Eq. (18)) and the  
13 corresponding values of the GCV functions (Eq. (9)) were calculated and plotted in  
14 Figures 5 and 6, respectively. The results illustrate that the analysis RMSE and the  
15 values of the GCV functions decrease sharply for the two EnKF with forecast error  
16 inflation schemes. However, the GCV function and the RMSE values of the improved  
17 EnKF were smaller than those of the constant inflated EnKF, indicating that the  
18 online estimate method performs better than the simple multiplicative inflation  
19 techniques with a constant value. The variability of the analysis RMSE was  
20 consistent with that of the GCV function for the EnKF with the forecast error  
21 inflation scheme. The correlation coefficient of the analysis RMSE and the value of  
22 the GCV function at the assimilation time step were approximately 0.76, which



1 indicates that the GCV function is a good criterion to estimate the inflation factor.

2       The ensemble analysis state members of the conventional EnKF, improved  
3 EnKF and constant inflated EnKF are shown in Figure 7, and the results indicate the  
4 uncertainty of the analysis state to some extent. The true trajectory obtained by the  
5 numerical solution is also plotted, and it illustrates that a larger difference occurred  
6 between the true trajectory and the ensemble analysis state members for the  
7 conventional EnKF than for the improved EnKF and constant inflated EnKF. In  
8 addition, the analysis state was more consistent with the true trajectory for the  
9 improved EnKF than that for the constant inflated EnKF. Therefore, the forecast error  
10 inflation can lead to a more accurate analysis state than the constant inflated EnKF.

11       The time-mean values of the forecast ensemble spread, the GAI statistics, the  
12 GCV functions and the analysis RMSE over 2000 time steps are listed in Table 1.  
13 These results illustrate that the forecast error inflation technique using the GCV  
14 function performs better than the constant inflated EnKF, which can indeed increase  
15 the analysis sensitivity to the observations and reduce the analysis RMSE.

16

### 17 ***3.3 Influence of ensemble size and observation number***

18       Intuitively, for any ensemble-based assimilation scheme, a large ensemble size  
19 will lead to small analysis errors; however, the computational costs are high for  
20 practical problems. The ensemble size in the practical land surface assimilation  
21 problem is usually several tens of members (Kirchgessner et al. 2014). The  
22 preferences of the proposed inflation method with respect to different ensemble sizes

1 (10, 30 and 50) were evaluated, and the results are listed in Table 1, which shows that  
2 using a 10-member ensemble produced a threefold increase in the analysis RMSE,  
3 while using a 50-member ensemble reduced the analysis RMSE by 20% relative to  
4 the analysis RMSE obtained using a 30-member ensemble. The forecast ensemble  
5 spread increased slightly from a 10-member ensemble to a 50-member ensemble. The  
6 GAI and GCV function values changed sharply from a 10-member ensemble to a  
7 30-member ensemble, and they became relatively stable from a 30-member ensemble  
8 to a 50-member ensemble. Ensembles less than 10 were unstable, and no significant  
9 changes occurred for ensembles greater than 50. Considering the computational costs  
10 for practical problems, a 30-member ensemble may be necessary to estimate  
11 statistically robust results.

12 To evaluate the preferences of the inflation method with respect to different  
13 numbers of observations, synthetic observations were generated at every other grid  
14 point and for every 4 time steps. Hence, a total of 20 observations were performed at  
15 each observation step in this case. The assimilation results with ensemble sizes of 10,  
16 30 and 50 are listed in Table 2, which shows that the GAI values were larger than  
17 those with 40-observations in all assimilation schemes. This finding may be related to  
18 the relatively small denominator of the GAI statistic (Eq. (16)) in the 20-observation  
19 experiments. The forecast ensemble spread does not change much but the GCV  
20 function and the RMSE values increase greatly in the 20-observation experiments  
21 with respect to those in the 40-observation experiments, which illustrates that more  
22 observations will lead to less analysis error.

1

2

### 3 **4. Discussion and Conclusions**

4

5       Accurate estimates of the forecast error covariance matrix are crucial to the  
6 success of any data assimilation scheme. In the conventional EnKF assimilation  
7 scheme, the forecast error covariance matrix is estimated as the sampling covariance  
8 matrix of the ensemble forecast states. However, a limited ensemble size and large  
9 model errors often cause the matrix to be underestimated, which produces an analysis  
10 state that over relies on the forecast and excludes observations, which can eventually  
11 cause the filter to diverge. Therefore, the forecast error inflation with proper inflation  
12 factors is increasingly important.

13       The use of multiplicative covariance inflation techniques can mitigate this  
14 problem to some extent. Several methods have been proposed in the literature, and  
15 each has different assumptions. For instance, the moment approach can be easily  
16 conducted based on the moment estimation of the innovation statistic. The maximum  
17 likelihood approach can obtain a more accurate inflation factor than the moment  
18 approach, but requires computing high dimensional matrix determinants. The  
19 Bayesian approach assumes a prior distribution for the inflation factor but is limited  
20 to spatially independent observational errors. In this study, the inflation factor was  
21 estimated using a GCV and the analysis sensitivity was detected.

22       The assimilation results showed that inflating the conventional EnKF using the

1 factor estimated by minimizing the GCV function can indeed reduce the analysis  
2 RMSE. Therefore, the GCV function can accurately quantify the goodness of fit of  
3 the error covariance matrix. In fact, the CV method can evaluate and compare  
4 learning algorithms and represents a widely used statistical method. In this study, the  
5 CV concept was adopted for the inflation factor estimation in the improved EnKF  
6 assimilation scheme and was validated with the Lorenz-96 model. The values of the  
7 GCV function obviously decreased in the proposed approach compared the  
8 conventional EnKF scheme. The analysis RMSE in the proposed approach was also  
9 much smaller than that in the conventional EnKF scheme, which suggests that this  
10 method of minimizing the GCV works well for estimating the inflation factor.

11 The highest computational cost when minimizing the GCV function is related  
12 to calculating the influence matrix  $\mathbf{A}_i(\lambda)$ . Because the matrix multiplication is  
13 commutative for the trace, the GCV function can be easily re-expressed as follows:

$$14 \quad GCV_i(\lambda) = \frac{p_i \mathbf{d}_i^T (\mathbf{H}_i \lambda \mathbf{P}_i \mathbf{H}_i^T + \mathbf{R}_i)^{-1} \mathbf{R}_i (\mathbf{H}_i \lambda \mathbf{P}_i \mathbf{H}_i^T + \mathbf{R}_i)^{-1} \mathbf{d}_i}{\left[ \text{Tr} \left( (\mathbf{H}_i \lambda \mathbf{P}_i \mathbf{H}_i^T + \mathbf{R}_i)^{-1} \mathbf{R}_i \right) \right]^2}. \quad (21)$$

15 Because both the numerator and denominator of the GCV function are scalars, the  
16 inverse matrix is needed only in  $(\mathbf{H}_i \lambda \mathbf{P}_i \mathbf{H}_i^T + \mathbf{R}_i)^{-1}$ , which can be effectively  
17 calculated using the Sherman–Morrison–Woodbury formula (Golub; Loan 1996).  
18 Furthermore, the inverse matrix calculation and the multiplication process are also  
19 indispensable for the conventional EnKF (Eq. (6)). Essentially, no additional  
20 computational burden is associated with the improved EnKF by minimizing the GCV  
21 function. Therefore, the total computational costs of the improved EnKF are feasible.

1           The analysis sensitivities in the proposed approach and in the conventional  
2 EnKF scheme were also investigated in this study. The time-averaged GAI statistic  
3 increases from about 10% in the conventional EnKF scheme to about 30% using the  
4 proposed inflation method. This illustrates that the inflation mitigates the problem of  
5 the analysis depending excessively on the forecast and excluding the observations.  
6 The relationship of the analysis state to the forecast state and the observations are  
7 more reasonable.

8           It is also worth noting that the inflation factor is assumed to be constant in  
9 space in this study, which may be not the case in realistic assimilation problems.  
10 Forcing all components of the state vector to use the same inflation factor could  
11 systematically overinflate the ensemble variances in sparsely observed areas,  
12 especially when the observations are unevenly distributed. In the presence of sparse  
13 observations, the state that is not observed can be improved only by the physical  
14 mechanism of the forecast model, although this improvement is limited. Therefore, a  
15 multiplicative inflation may not be sufficiently effective to enhance the assimilation  
16 accuracy. In this case, the additive inflation and the localization technique can be  
17 applied to further improve the assimilation quality in the presence of sparse  
18 observations (Miyoshi; Kunii 2011; Yang et al. 2015).

19           The examples shown here using the Lorenz-96 model illustrate the feasibility of  
20 this approach for using GCV as a metric to estimate the covariance inflation factor. In  
21 the case studies conducted in Section 3, the observations were relatively evenly  
22 distributed and the assimilation accuracy could indeed be improved by the forecast

1 error inflation technique. These findings provide insights on the methodology and  
 2 validation of the Lorenz-96 model and illustrate the feasibility of our approach. In the  
 3 near future, methods of modifying the adaptive procedure to suit the system with  
 4 unevenly distributed observations and applying the proposed methodologies using  
 5 more sophisticated dynamic and observation systems will be investigated.

6

## 7 **Appendix A**

8 From Eq. (2), the normalized observation equation can be defined as follows:

$$9 \quad \tilde{\mathbf{y}}_i^o = \mathbf{R}_i^{-1/2} \mathbf{H}_i \mathbf{x}_i^t + \tilde{\boldsymbol{\varepsilon}}_i, \quad (\text{A1})$$

10 where  $\tilde{\mathbf{y}}_i^o = \mathbf{R}_i^{-1/2} \mathbf{y}_i^o$  is the normalized observation vector and  $\tilde{\boldsymbol{\varepsilon}}_i \sim N(\mathbf{0}, \mathbf{I})$ ;  $\mathbf{I}_{p_i}$  is  
 11 the identity matrix with the dimensions  $p_i \times p_i$ . Similarly, the normalized analysis  
 12 vector is  $\tilde{\mathbf{y}}_i^a = \mathbf{R}_i^{-1/2} \mathbf{H}_i \mathbf{x}_i^a$  and the influence matrix  $\mathbf{A}_i$  relates the normalized  
 13 observation vector to the normalized analysis vector, thereby ignoring the normalized  
 14 forecast state in the observation space (Gu 2002):

$$15 \quad \tilde{\mathbf{y}}_i^a - \mathbf{R}_i^{-1/2} \mathbf{H}_i \mathbf{x}_i^f = \mathbf{A}_i \left( \tilde{\mathbf{y}}_i^o - \mathbf{R}_i^{-1/2} \mathbf{H}_i \mathbf{x}_i^f \right). \quad (\text{A2})$$

16 Because the analysis state  $\mathbf{x}_i^a$  is given by Eq. (5), the influence matrix  $\mathbf{A}_i$  can be  
 17 verified as follows:

$$18 \quad \mathbf{A}_i = \mathbf{I}_{p_i} - \mathbf{R}_i^{1/2} \left( \mathbf{H}_i \mathbf{P}_i \mathbf{H}_i^T + \mathbf{R}_i \right)^{-1} \mathbf{R}_i^{1/2}. \quad (\text{A3})$$

19 If the initial forecast error covariance matrix is inflated as described in Section 2.2,  
 20 then the influence matrix is treated as the following function of  $\lambda$

$$21 \quad \mathbf{A}_i(\lambda) = \mathbf{I}_{p_i} - \mathbf{R}_i^{1/2} \left( \mathbf{H}_i \lambda \mathbf{P}_i \mathbf{H}_i^T + \mathbf{R}_i \right)^{-1} \mathbf{R}_i^{1/2}, \quad (\text{A4})$$

22 The principle of CV is to minimize the estimated error at the observation grid

1 point. Lacking an independent validation data set, a common alternative strategy is to  
 2 minimize the squared distance between the normalized observation value and the  
 3 analysis value while not using the observation on the same grid point, which is the  
 4 following objective function:

$$5 \quad V_i(\lambda) = \frac{1}{p_i} \sum_{k=1}^{p_i} \left( \tilde{\mathbf{y}}_{i,k}^o - \left( \mathbf{R}_i^{-1/2} \mathbf{H}_i \mathbf{x}_i^{a[k]} \right)_k \right)^2, \quad (\text{A5})$$

6 where  $\mathbf{x}_i^{a[k]}$  is the minima of the following “delete-one” objective function:

$$7 \quad \left( \mathbf{x} - \mathbf{x}_i^f \right)^T (\lambda \mathbf{P}_i)^{-1} \left( \mathbf{x} - \mathbf{x}_i^f \right) + \left( \mathbf{y}_i^o - \mathbf{H}_i \mathbf{x} \right)_{-k}^T \mathbf{R}_{i,-k}^{-1/2} \left( \mathbf{y}_i^o - \mathbf{H}_i \mathbf{x} \right)_{-k}. \quad (\text{A6})$$

8 The subscript  $-k$  indicates a vector (matrix) with its  $k$ -th element ( $k$ -th row and  
 9 column) deleted. Instead of minimizing Eq. (A6)  $p_i$  times, the objective function  
 10 (Eq. (A5)) has another more simple expression (Gu 2002):

$$11 \quad V_i(\lambda) = \frac{1}{p_i} \sum_{k=1}^{p_i} \frac{\left( \tilde{\mathbf{y}}_{i,k}^o - \left( \mathbf{R}_i^{-1/2} \mathbf{H}_i \mathbf{x}_i^a \right)_k \right)^2}{(1 - a_{k,k})^2}, \quad (\text{A7})$$

12 where  $a_{k,k}$  is the element at the site pair  $(k, k)$  of the influence matrix  $\mathbf{A}_i(\lambda)$ . Then,

13  $a_{k,k}$  is substituted with the average  $\frac{1}{p_i} \sum_{k=1}^{p_i} a_{k,k} = \frac{1}{p_i} \text{Tr}(\mathbf{A}_i(\lambda))$  and the constant is

14 ignored to obtain the following GCV statistic (Gu 2002):

$$15 \quad GCV_i(\lambda) = \frac{\frac{1}{p_i} \mathbf{d}_i^T \mathbf{R}_i^{-1/2} \left( \mathbf{I}_{p_i} - \mathbf{A}_i(\lambda) \right)^2 \mathbf{R}_i^{-1/2} \mathbf{d}_i}{\left[ \frac{1}{p_i} \text{Tr} \left( \mathbf{I}_{p_i} - \mathbf{A}_i(\lambda) \right) \right]^2}. \quad (\text{A8})$$

16

## 17 **Appendix B**

18 The sensitivities of the analysis to the observation are defined as follows:

$$1 \quad \mathbf{S}_i^o = \frac{\partial \tilde{\mathbf{y}}_i^a}{\partial \tilde{\mathbf{y}}_i^o} = \mathbf{R}_i^{1/2} \mathbf{K}_i^T \mathbf{H}_i^T \mathbf{R}_i^{-1/2}, \quad (\text{B1})$$

2 Substitute the Kalman gain matrix  $\mathbf{K}_i = \mathbf{P}_i \mathbf{H}_i^T (\mathbf{H}_i \mathbf{P}_i \mathbf{H}_i^T + \mathbf{R}_i)^{-1}$  into  $\mathbf{S}_i^o$ , then:

$$\begin{aligned}
3 \quad \mathbf{S}_i^o &= \mathbf{R}_i^{1/2} \mathbf{K}_i^T \mathbf{H}_i^T \mathbf{R}_i^{-1/2} \\
4 \quad &= \mathbf{R}_i^{1/2} (\mathbf{H}_i \mathbf{P}_i \mathbf{H}_i^T + \mathbf{R}_i)^{-1} \mathbf{H}_i \mathbf{P}_i \mathbf{H}_i^T \mathbf{R}_i^{-1/2} \\
5 \quad &= \mathbf{R}_i^{1/2} (\mathbf{H}_i \mathbf{P}_i \mathbf{H}_i^T + \mathbf{R}_i)^{-1} (\mathbf{H}_i \mathbf{P}_i \mathbf{H}_i^T + \mathbf{R}_i - \mathbf{R}_i) \mathbf{R}_i^{-1/2} \\
6 \quad &= \mathbf{R}_i^{1/2} (\mathbf{H}_i \mathbf{P}_i \mathbf{H}_i^T + \mathbf{R}_i)^{-1} (\mathbf{H}_i \mathbf{P}_i \mathbf{H}_i^T + \mathbf{R}_i) \mathbf{R}_i^{-1/2} - \mathbf{R}_i^{1/2} (\mathbf{H}_i \mathbf{P}_i \mathbf{H}_i^T + \mathbf{R}_i)^{-1} \mathbf{R}_i \mathbf{R}_i^{-1/2} \\
7 \quad &= \mathbf{I}_{p_i} - \mathbf{R}_i^{1/2} (\mathbf{H}_i \mathbf{P}_i \mathbf{H}_i^T + \mathbf{R}_i)^{-1} \mathbf{R}_i \mathbf{R}_i^{-1/2} \\
8 \quad &= \mathbf{A}_i \quad (\text{B2})
\end{aligned}$$

9 Therefore, the sensitivity matrix  $\mathbf{S}_i^o$  is equal to the influence matrix  $\mathbf{A}_i$ .

10

11 **Acknowledgements.** This work is supported by the National Natural Science  
12 Foundation of China (Grant No. 91647202), the National Basic Research Program of  
13 China (Grant No. 2015CB953703) and the National Natural Science Foundation of  
14 China (Grant No. 41405098).

15



1 **References**

2 Table 1. Time-mean values of the forecast ensemble spread, GAI statistics, GCV  
 3 functions and analysis RMSE over 2000 time steps. The ensemble size is selected as  
 4 10, 30 and 50, respectively.

Ensemble Size	Conventional EnKF			EnKF with forecast inflation		
	10	30	50	10	30	50
Spread	0.23	0.36	0.41	3.26	3.32	3.45
GAI	4.56%	10.78%	13.58%	5.24%	29.21%	35.63%
GCV	36.38	31.14	25.21	35.56	3.29	2.30
RMSE	4.50	4.01	3.52	3.74	1.10	0.88

5

6 Table 2. Same as in Table 1 but for 20 observations.

Ensemble Size	Conventional EnKF			EnKF with forecast inflation		
	10	30	50	10	30	50
Spread	0.41	0.59	0.68	3.33	3.36	3.48
GAI	10.77%	20.92%	26.41%	13.25%	35.09%	41.28%
GCV	33.64	22.89	14.97	32.17	14.99	5.19
RMSE	4.85	4.10	3.29	4.39	3.46	2.86

7

1 **Figure captions**

2 Figure 1. Flowchart of the proposed assimilation scheme.

3 Figure 2. Time series of the estimated inflation factors by minimizing the GCV  
4 function.

5 Figure 3. Forecast ensemble spread of the conventional EnKF (black line), the  
6 improved EnKF (blue line) and the constant inflated EnKF (red line).

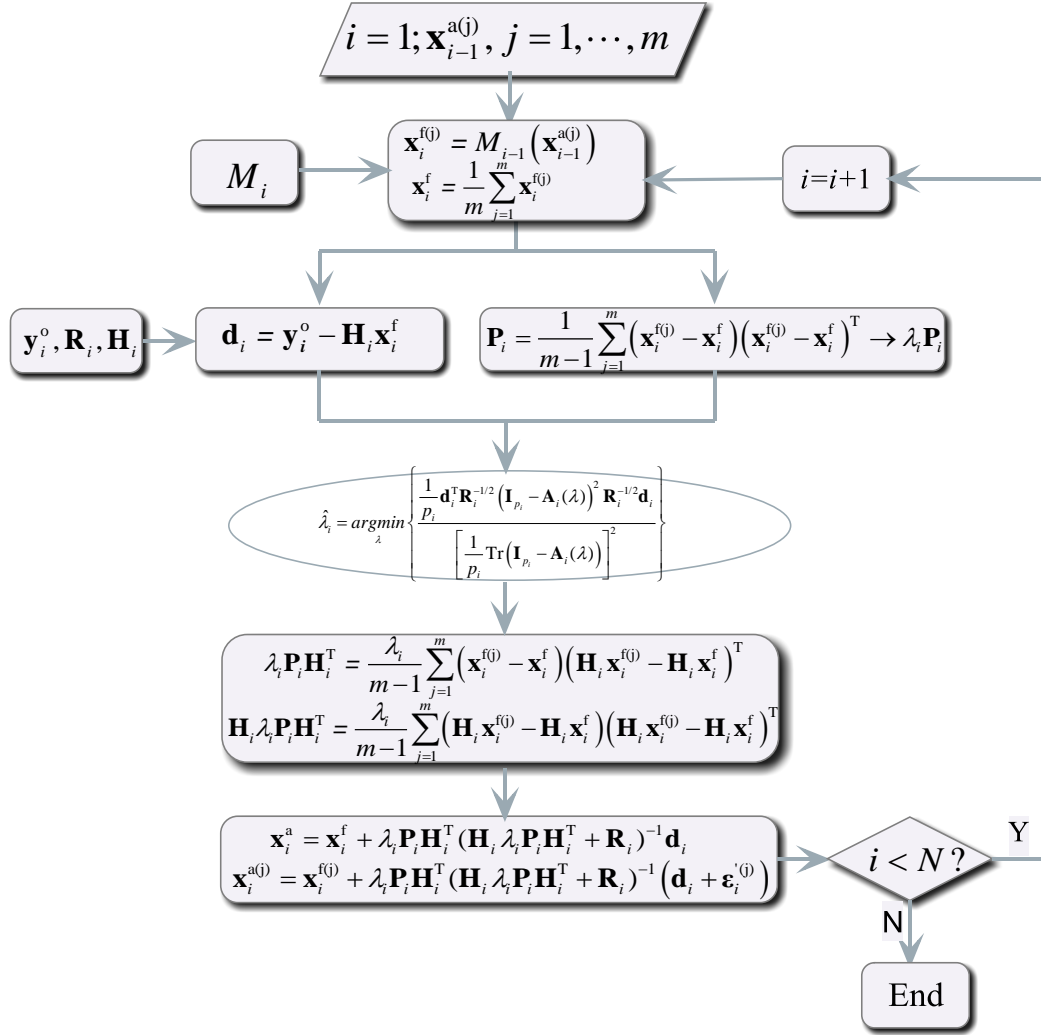
7 Figure 4. GAI statistics of the conventional EnKF (black line), the improved EnKF  
8 (blue line) and the constant inflated EnKF (red line).

9 Figure 5. Analysis RMSE of the conventional EnKF (black line), the improved EnKF  
10 (blue line) and the constant inflated EnKF (red line).

11 Figure 6. GCV function values of the conventional EnKF (black line), the improved  
12 EnKF (blue line) and the constant inflated EnKF (red line).

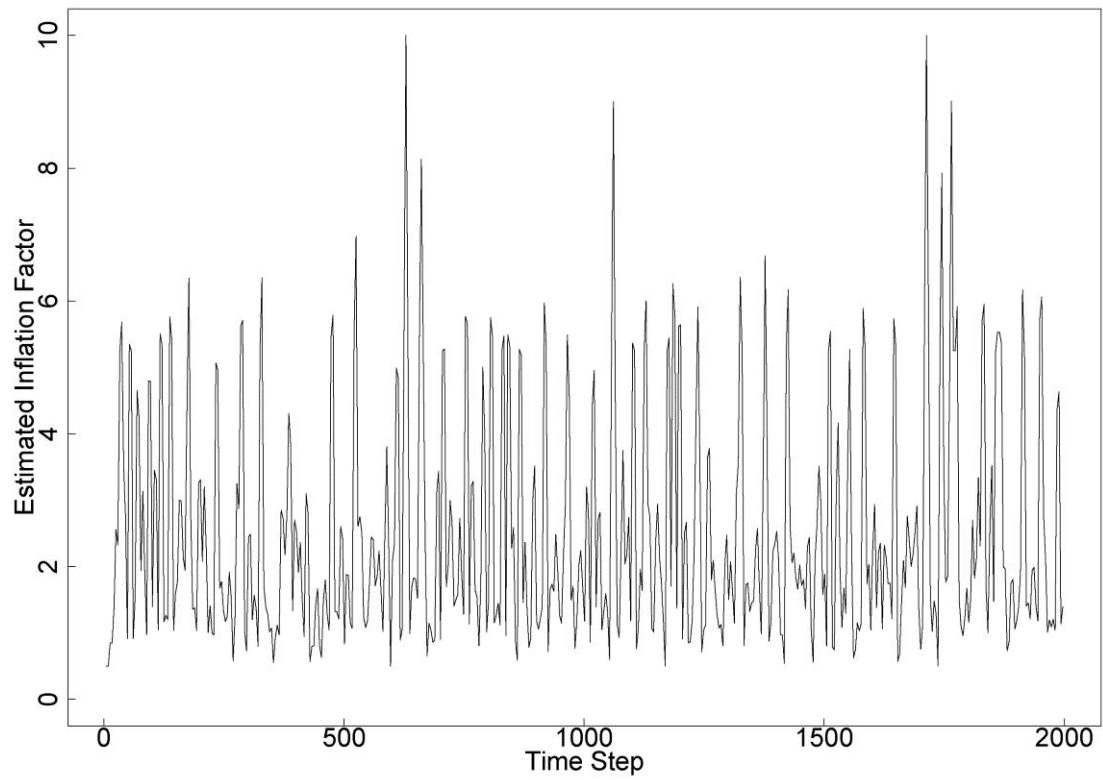
13 Figure 7. Ensemble analysis state members of the conventional EnKF (black line), the  
14 improved EnKF (blue line) and the constant inflated EnKF (red line). The green line  
15 refers to the true trajectory obtained by the numerical solution.

16



1  
2  
3

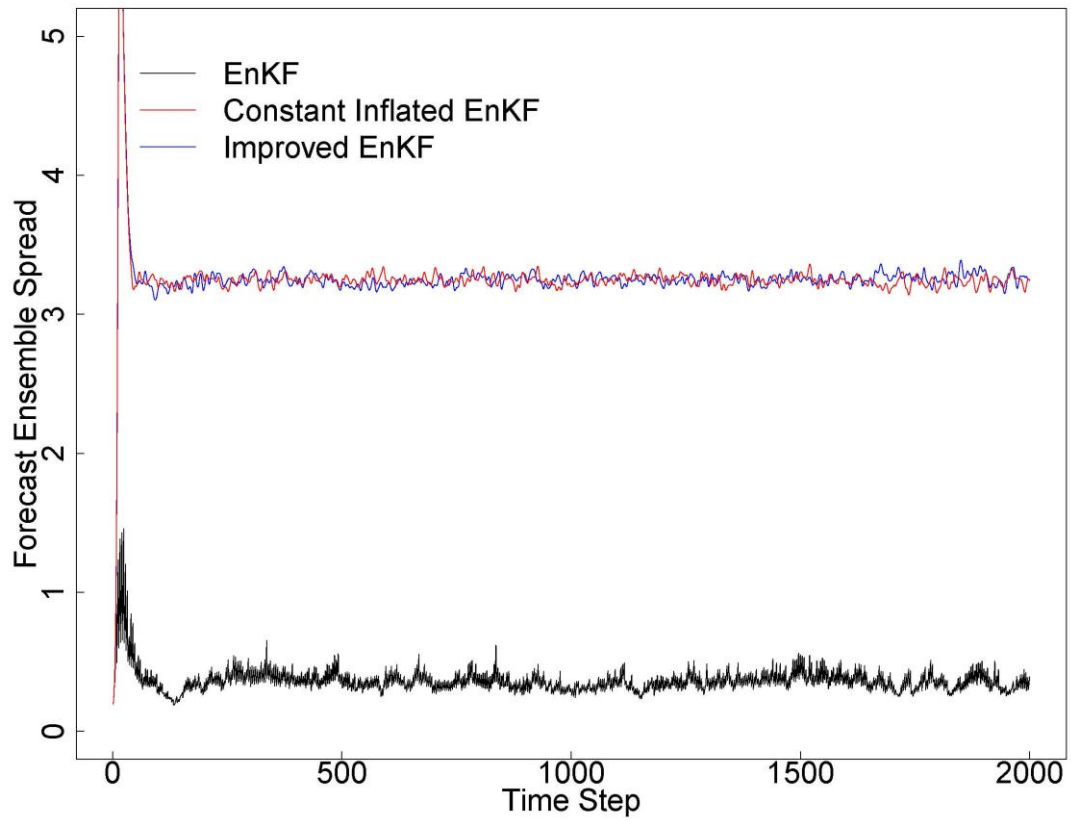
Figure 1. Flowchart of the proposed assimilation scheme.



1

2 Figure 2. Time series of the estimated inflation factors by minimizing the GCV  
3 function.

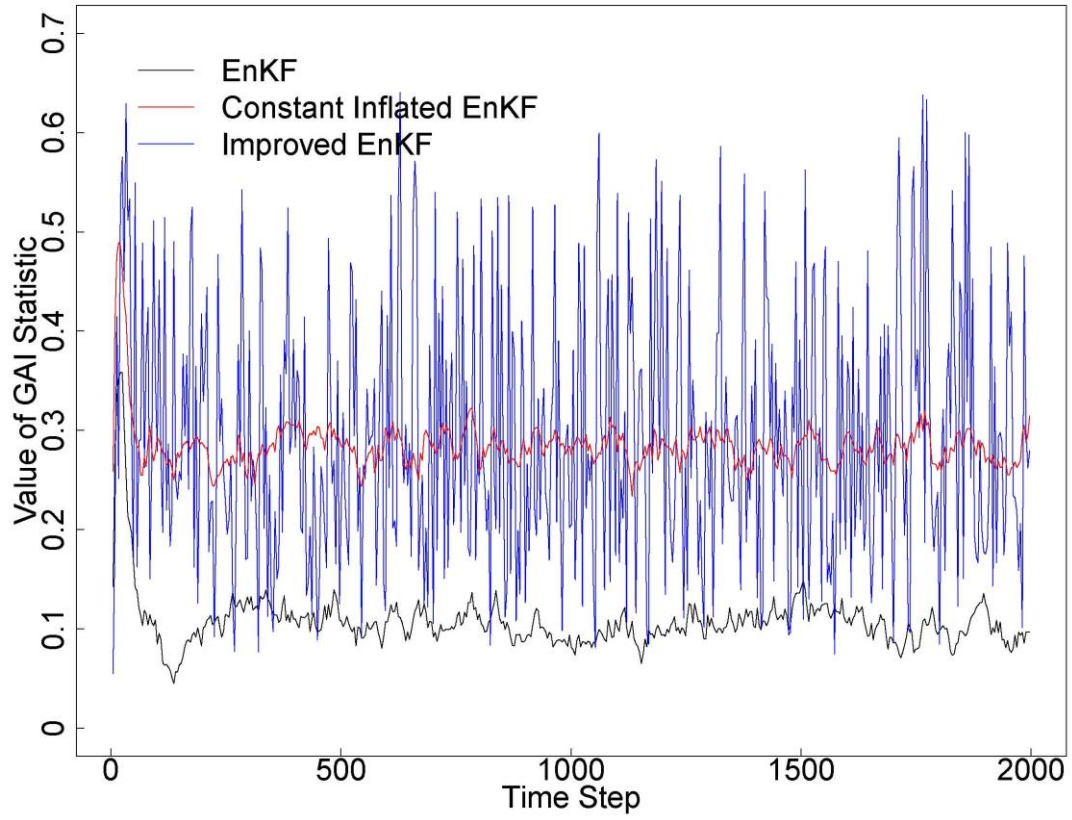
4



1

2 Figure 3. Forecast ensemble spread of the conventional EnKF (black line), the  
3 improved EnKF (blue line) and the constant inflated EnKF (red line).

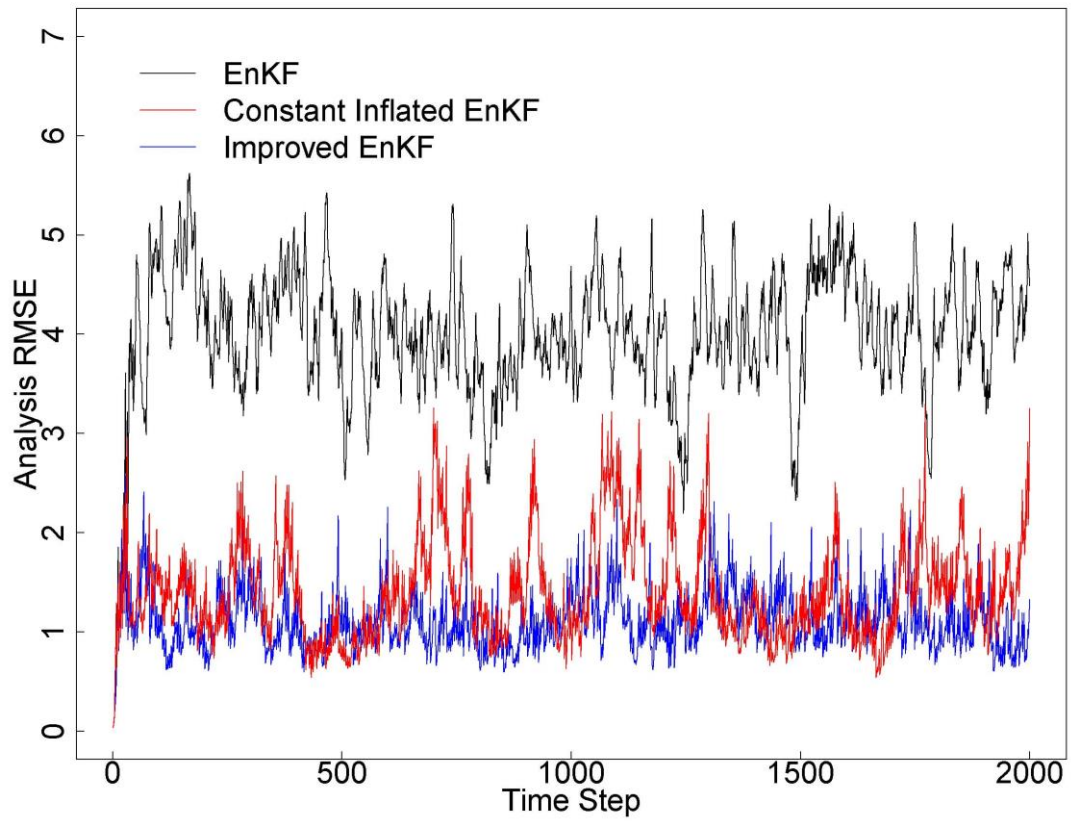
4



1

2 Figure 4. GAI statistics of the conventional EnKF (black line), the improved EnKF  
3 (blue line) and the constant inflated EnKF (red line).

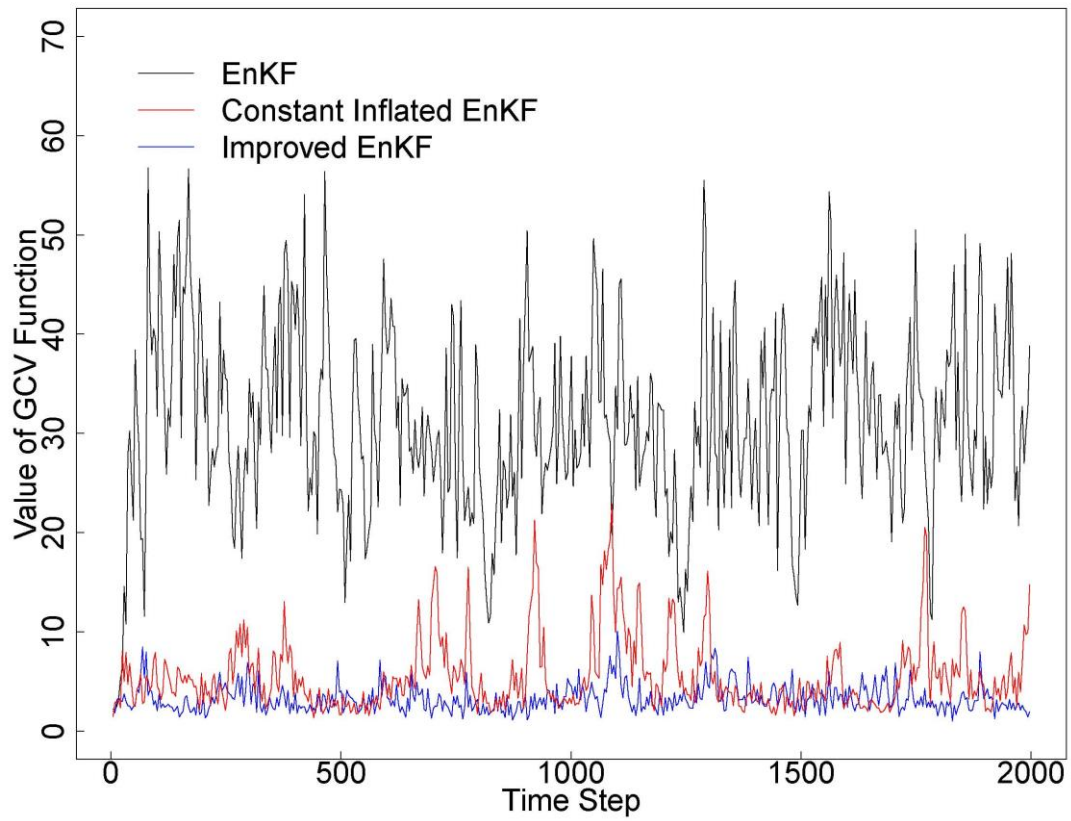
4



1

2 Figure 5. Analysis RMSE of the conventional EnKF (black line), the improved EnKF  
3 (blue line) and the constant inflated EnKF (red line).

4

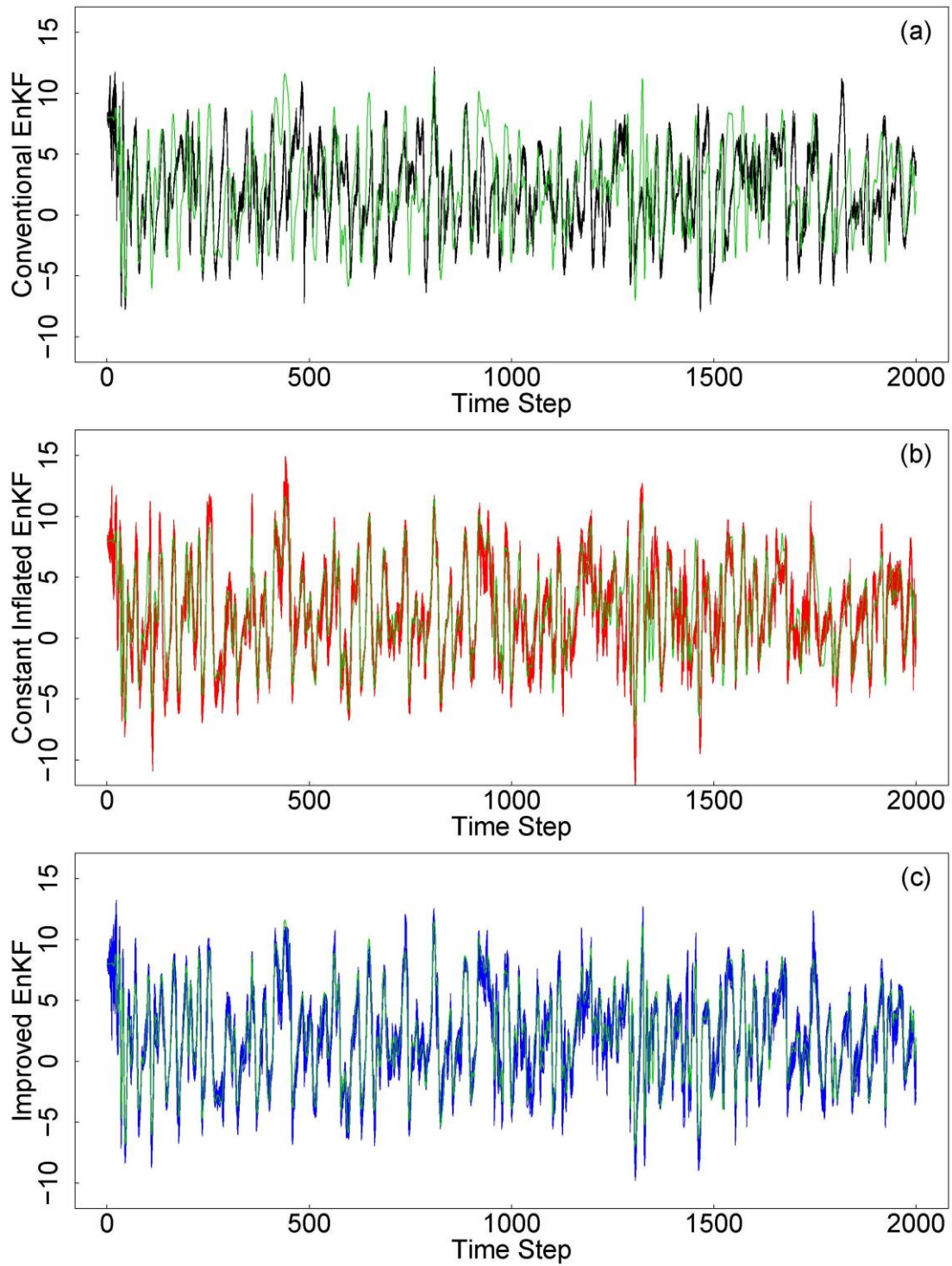


1

2 Figure 6. GCV function values of the conventional EnKF (black line), the improved  
3 EnKF (blue line) and the constant inflated EnKF (red line).

4





1

2 Figure 7. Ensemble analysis state members of the conventional EnKF (black line), the  
 3 improved EnKF (blue line) and the constant inflated EnKF (red line). The green line  
 4 refers to the true trajectory obtained by the numerical solution.

5

## 1   **References**

- 2   Allen, D. M., 1974: The relationship between variable selection and data augmentation and a method  
3   for prediction. *Technometrics*, **16**, 125-127.
- 4   Anderson, J. L., 2007: An adaptive covariance inflation error correction algorithm for ensemble filters.  
5   *Tellus*, **59A**, 210-224.
- 6   Anderson, J. L., 2009: Spatially and temporally varying adaptive covariance inflation for ensemble  
7   filters. *Tellus*, **61A**, 72-83.
- 8   Anderson, J. L., and S. L. Anderson, 1999: A Monte Carlo implementation of the nonlinear filtering  
9   problem to produce ensemble assimilations and forecasts. *Monthly Weather Review*, **127**, 2741-2758.
- 10   Burgers, G., P. J. Leeuwen, and G. Evensen, 1998: Analysis scheme in the ensemble kalman filter.  
11   *Monthly Weather Review*, **126**, 1719-1724.
- 12   Butcher, J. C., 2003: *Numerical methods for ordinary differential equations*.   JohnWiley & Sons, 425  
13   pp.
- 14   Cardinali, C., S. Pezzulli, and E. Andersson, 2004: Influence - matrix diagnostic of a data assimilation  
15   system. *Quarterly Journal of the Royal Meteorological Society*, **130**, 2767-2786.
- 16   Constantinescu, E. M., A. Sandu, T. Chai, and G. R. Carmichael, 2007: Ensemble-based chemical data  
17   assimilation I: general approach. *Quarterly Journal of the Royal Meteorological Society*, **133**,  
18   1229-1243.
- 19   Craven, P., and G. Wahba, 1979: Smoothing noisy data with spline functions. *Numerische Mathematik*,  
20   **31**, 377-403.
- 21   Dee, D. P., 1995: On-line estimation of error covariance parameters for atmospheric data assimilation.  
22   *Monthly Weather Review*, **123**, 1128-1145.
- 23   Dee, D. P., and A. M. Silva, 1999: Maximum-likelihood estimation of forecast and observation error  
24   covariance parameters part I: methodology. *Monthly Weather Review*, **127**, 1822-1834.
- 25   Ellison, C. J., J. R. Mahoney, and J. P. Crutchfield, 2009: Prediction, Retrodiction, and the Amount of  
26   Information Stored in the Present. *Journal of Statistical Physics*, **136**, 1005-1034.
- 27   Eubank, R. L., 1999: *Nonparametric regression and spline smoothing*.   Marcel Dekker, Inc., 338 pp.
- 28   Evensen, G., 1994: Sequential data assimilation with a nonlinear quasi-geostrophic model using  
29   Monte Carlo methods to forecast error statistics. *Journal of Geophysical Research*, **99**, 10143-10162.
- 30   Gentle, J. E., W. Hardle, and Y. Mori, 2004: *Handbook of computational statistics: concepts and*  
31   *methods*.   Springer, 1070 pp.
- 32   Golub, G. H., and C. F. V. Loan, 1996: *Matrix Computations*.   The Johns Hopkins University Press:  
33   Baltimore.
- 34   Green, P. J., and B. W. Silverman., 1994: *Nonparametric Regression and Generalized Additive Models*.  
35   Vol. 182, Chapman and Hall,.
- 36   Gu, C., 2002: *Smoothing Spline ANOVA Models*.   Springer-Verlag, 289 pp.
- 37   Gu, C., and G. Wahba, 1991: Minimizing GCV/GML scores with multiple smoothing parameters via the  
38   Newton method. *SIAM Journal on Scientific and Statistical Computation*, **12**, 383-398.
- 39   Ide, K., P. Courtier, M. Ghil, and A. C. Lorenc, 1997: Unified notation for data assimilation operational  
40   sequential and variational. *Journal of the Meteorological Society of Japan*, **75**, 181-189.
- 41   Kirchgeßner, P., L. Berger, and A. B. Gerstner, 2014: On the choice of an optimal localization radius in  
42   ensemble Kalman filter methods. *Monthly Weather Review*, **142**, 2165-2175.
- 43   Li, H., E. Kalnay, and T. Miyoshi, 2009: Simultaneous estimation of covariance inflation and

1 observation errors within an ensemble Kalman filter. *Quarterly Journal of the Royal Meteorological*  
2 *Society*, **135**, 523-533.

3 Liang, X., X. Zheng, S. Zhang, G. Wu, Y. Dai, and Y. Li, 2012: Maximum Likelihood Estimation of Inflation  
4 Factors on Error Covariance Matrices for Ensemble Kalman Filter Assimilation. *Quarterly Journal of the*  
5 *Royal Meteorological Society*, **138**, 263-273.

6 Liu, J., E. Kalnay, T. Miyoshi, and C. Cardinali, 2009: Analysis sensitivity calculation in an ensemble  
7 Kalman filter. *Quarterly Journal of the Royal Meteorological Society*, **135**, 1842-1851.

8 Lorenz, E. N., 1996: Predictability a problem partly solved.

9 Lorenz, E. N., and K. A. Emanuel, 1998: Optimal sites for supplementary weather observations  
10 simulation with a small model. *Journal of the Atmospheric Sciences*, **55**, 399-414.

11 Miller, R. N., M. Ghil, and F. Gauthiez, 1994: Advanced data assimilation in strongly nonlinear  
12 dynamical systems. *Journal of the Atmospheric Sciences*, **51**, 1037-1056.

13 Miyoshi, T., 2011: The Gaussian approach to adaptive covariance inflation and its implementation  
14 with the local ensemble transform Kalman filter. *Monthly Weather Review*, **139**, 1519-1534.

15 Miyoshi, T., and M. Kunii, 2011: The Local Ensemble Transform Kalman Filter with the Weather  
16 Research and Forecasting Model: Experiments with Real Observations. *Pure & Applied Geophysics*,  
17 **169**, 321-333.

18 Pena, D., and V. J. Yohai, 1991: The detection of influential subsets in linear regression using an  
19 influence matrix. *Journal of the Royal Statistical Society*, **57**, 145-156.

20 Reichle, R. H., 2008: Data assimilation methods in the Earth sciences. *Advances in Water Resources*,  
21 **31**, 1411-1418.

22 Saltelli, A., S. Tarantola, F. Campolongo, and M. Ratto, 2004: *Sensitivity Analysis in Practice: A Guide to*  
23 *Assessing Scientific Models*. JohnWiley & Sons, 219 pp.

24 Saltelli, A., and Coauthors, 2008: *Global Sensitivity Analysis: The Primer*. John Wiley & Sons, 292 pp.

25 Talagrand, O., 1997: Assimilation of Observations, an Introduction. *Journal of the Meteorological*  
26 *Society of Japan*, **75**, 191-209.

27 Tippett, M. K., J. L. Anderson, C. H. Bishop, T. M. Hamill, and J. S. Whitaker, 2003: Notes and  
28 correspondence ensemble square root filter. *Monthly Weather Review*, **131**, 1485-1490.

29 Wahba, G., and S. Wold, 1975: A completely automatic french curve. *Communications in Statistics*, **4**,  
30 1-17.

31 Wahba, G., R. J. Donald, F. Gao, and J. Gong, 1995: Adaptive Tuning of Numerical Weather Prediction  
32 Models: Randomized GCV in Three- and Four-Dimensional Data Assimilation. *Monthly Weather*  
33 *Review*, **123**, 3358-3369.

34 Wand, M. P., and M. C. Jones, 1995: *Kernel Smoothing*. Chapman and Hall, 212 pp.

35 Wang, X., and C. H. Bishop, 2003: A comparison of breeding and ensemble transform kalman filter  
36 ensemble forecast schemes. *Journal of the Atmospheric Sciences*, **60**, 1140-1158.

37 Wu, G., X. Zheng, L. Wang, S. Zhang, X. Liang, and Y. Li, 2013: A New Structure for Error Covariance  
38 Matrices and Their Adaptive Estimation in EnKF Assimilation. *Quarterly Journal of the Royal*  
39 *Meteorological Society*, **139**, 795-804.

40 Wu, G., X. Yi, X. Zheng, L. Wang, X. Liang, S. Zhang, and X. Zhang, 2014: Improving the Ensemble  
41 Transform Kalman Filter Using a Second-order Taylor Approximation of the Nonlinear Observation  
42 Operator. *Nonlinear Processes in Geophysics*, **21**, 955-970.

43 Xu, T., J. J. Gómez-Hernández, H. Zhou, and L. Li, 2013: The power of transient piezometric head data  
44 in inverse modeling: An application of the localized normal-score EnKF with covariance inflation in a

- 1 heterogenous bimodal hydraulic conductivity field. *Advances in Water Resources*, **54**, 100-118.
- 2 Yang, S.-C., E. Kalnay, and T. Enomoto, 2015: Ensemble singular vectors and their use as additive
- 3 inflation in EnKF. *Tellus A*, **67**.
- 4 Zheng, X., 2009: An adaptive estimation of forecast error statistic for Kalman filtering data
- 5 assimilation. *Advances in Atmospheric Sciences*, **26**, 154-160.
- 6 Zheng, X., and R. Basher, 1995: Thin-plate smoothing spline modeling of spatial climate data and its
- 7 application to mapping south Pacific rainfall. *Monthly Weather Review*, **123**, 3086-3102.