

1 The authors gratefully acknowledge the anonymous reviewer for his/her insightful  
2 comments and constructive suggestions that lead to the significant improvement of  
3 the quality of this manuscript. The authors are also grateful to the editors for their  
4 very kind help and comments. We have checked the manuscript carefully and tried  
5 our best to address all the comments. Below boldface is used to indicate the  
6 comments from the reviewers and editor, and italics is for the point-by-point  
7 responses.

8

9 The main changes in the revised version are: add more detailed overview about the  
10 GCV method in the introduction and discussion sections. The English grammar is  
11 also checked.

12 Additionally, I wonder if it is possible to add a coauthor. His contribution is the more  
13 discussion of the experiment results and the English language polish.

14

1 Comments from editor:

2 **Both the referees have suggested a few additional changes, which they and I**  
3 **consider to be minor in nature. Thus I recommend this manuscript for**  
4 **publication subject to minor changes. If you choose not to implement some of**  
5 **these changes, please state clearly the reasons for this decision.**

6 **Response:** *Thanks for your review and positive recommendation. In this revised*  
7 *version, we mainly added more detailed overview about the GCV method to the*  
8 *introduction and discussion sections. (P6 L7-15 and P20 L8-18)*

9

1 Comments from anonymous reviewer

2 General Comment:

3 **This paper presents a new technique in estimating model error covariance**  
4 **inflation factor which is necessary to arrest the divergence of the filter. This**  
5 **paper relies on estimating the inflation factor from an objective function**  
6 **inspired from the domain of generalized cross validation (GCV) techniques**  
7 **widely used in the field of machine-learning. The author also shows that this**  
8 **method, in comparison to a no-inflation based method and a spatio-temporal**  
9 **uniform inflation based scheme, improves the root-mean-squared error and**  
10 **enhances the influence of observations on the analysis when applied to the**  
11 **Lorenz 96 model. The author also performs a series of sensitivity experiments to**  
12 **claim the superiority of his method. The computational cost involved in the new**  
13 **method is shown to be marginally larger than the existing methods.**

14 **Response:** *Thanks for your review and positive recommendation.*

15

16 **However, since this method is imported from a different field, it will be nice to**  
17 **present a little more detailed overview of the new method including its**  
18 **properties for the benefit of the data assimilation community at large. It will be**  
19 **good to highlight the limitations of this method, if any, as well.**

20 **Response:** *Thanks for your constructive comments. We added more detailed*  
21 *overview about GCV method to the introduction and discussion sections in the*  
22 *revised version as follows. (P6 L7-15 and P20 L8-18)*

23 *Actually the GCV criterion is based on a predictive mean-square error criterion*  
24 *that attempts to obtain a best estimate (Wahba et al. 1995). It has a rotation-invariant*  
25 *property that is relative to the orthogonal transformation of the observations and is a*  
26 *consistent estimate of the relative loss (Gu 2002). For the inverse problems in such*  
27 *fields as meteorological data assimilation, GCV method can choose parameters*  
28 *systematically by minimizing a given objective function that will improve the*  
29 *inversion results. It can particularly select parameters that reflect not only*

1 *measurement accuracies from different sources but also model capability (Krakauer*  
2 *et al. 2004).*

3 *It can be applied in inverse problems in such fields as meteorological data*  
4 *assimilation (Wahba et al. 1995). Specifically, GCV provides a well-characterized*  
5 *method, which can select a regularization parameter by minimizing the predictive*  
6 *data errors with rotation-invariant in a least squares solution (MacCarthy et al.*  
7 *2011). In data assimilation research fields, observation data such as in-situ*  
8 *observation and remote sensing data are usually from different sources. GCV is*  
9 *particularly useful for choosing relative parameters that reflect not only*  
10 *measurement accuracies from different sources but also model capability (Krakauer*  
11 *et al. 2004). Apparently, GCV method requires calculating the trace of a large matrix,*  
12 *which may be commonly computationally prohibitive for large inverse problems*  
13 *(MacCarthy et al. 2011).*

14

15 **It is understandable that English is not the native language of the author.**  
16 **However through multiple iterations the paper has come to a much better shape**  
17 **and needs some minor language corrections. I recommend publications after**  
18 **minor changes.**

19 **Response:** *Thanks for your comments. The manuscript language has been checked*  
20 *again in the revised version.*

21

22 Specific Comment:

23 **1) The author claims that the GCV has many favorable properties but didn't**  
24 **mention any except the rotation-invariant property. It will be nice to put down**  
25 **some of the properties and how it advantage to the field of data assimilation.**

26 **Response:** *Thanks for your comments. The more detailed overview about GCV*  
27 *method are added to the introduction and discussion sections in the revised version*  
28 *(P6 L7-15 and P20 L8-18). Please see the reply to the general comment.*

29

1 **2) P6 L6: It will be nice to elaborate on the claims made by the author since this**  
2 **technique is new to the data assimilation community.**

3 **Response:** *Thanks for your comments. For the data assimilation research fields,*  
4 *GCV method can choose parameters systematically by minimizing a given objective*  
5 *function that will improve the assimilation results. It can particularly select*  
6 *parameters that reflect not only measurement accuracies from different sources but*  
7 *also model capability. We have added this explanation to the revised version (P6*  
8 *L10-15).*

9

10 **3) In Fig 1, there are two “N”. One stands for the number of time-steps and the**  
11 **other for “NO”. It will be good to remove this confusion.**

12 **Response:** *Sorry for this confusion. It has been revised in Figure 1.*

13

14 **4) P18 L7-9: The author claims that a 30-member ensemble is necessary to**  
15 **estimate statistically robust result. However, this study is true only for**  
16 **Lorenz-96 model and should not be generalized for other systems. A system in**  
17 **which the errors grow in multiple directions (>30 ) will need more ensembles to**  
18 **produce statistically robust results.**

19 **Response:** *Thank you for point out this. In the revised version, we have added this*  
20 *explanation. (P18 L14-15)*

21

22 Technical Comment:

23 **1) P4 L8: “may be missed in...” may be changed to “may not have been captured**  
24 **by ...”.**

25 **Response:** *The phrase has been changed.*

26

27 **2) P4 L12: “by inflation factor ...” may be changed to “by an inflation factor...”.**

28 **Response:** *The phrase has been changed.*

29

1 **3) P4 L16: “It also seems ...”. This sentence is not clear and should be**  
2 **restructured.**

3 **Response:** *The sentence has changed to “It is not appropriate to ...”.*

4

5 **4) P4 L20: “the inflation factor can be...” may be changed to “the inflation factor**  
6 **is...”.**

7 **Response:** *The word has been changed.*

8

9 **5) P8 L17: “(observation minus forecast residual)” may be changed to**  
10 **“(observation minus forecast residual in observation space)”.**

11 **Response:** *The phrase has been changed.*

12

13 **6) P13 L15: “...was well defined...” should be changed to “... is well defined...”.**

14 **Response:** *The word has been changed.*

15

16 **7) P16 L16: “The variety of the analysis RMSE ...”. It is not clear what the**  
17 **author wants to convey.**

18 **Response:** *This confused sentence has been deleted, without affect the meaning of*  
19 *this paragraph.*

20

21 *Again, thanks for your thorough review and constructive comments.*

22 *The references in this reply are listed below.*

23

24 Gu, C., 2002: *Smoothing Spline ANOVA Models*. Springer-Verlag, 289 pp.

25 Krakauer, N. Y., T. Schneider, J. T. Randerson, and S. C. Olsen, 2004: Using generalized  
26 cross-validation to select parameters in inversions for regional carbon fluxes. *Geophysical Research*  
27 *Letters*, **31**.

28 MacCarthy, J. K., B. Borchers, and R. C. Aster, 2011: Efficient stochastic estimation of the model  
29 resolution matrix diagonal and generalized cross-validation for large geophysical inverse problems.  
30 *Journal of Geophysical Research*, **116**.

31 Wahba, G., D. R. Johnson, 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 Review*,  
33 **123**, 3358-3369.

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

3

4 Guocan Wu<sup>1,2</sup> and Xiaogu Zheng<sup>3</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 3 Key Laboratory of Regional Climate-Environment Research for East Asia, Institute  
11 of Atmospheric Physics, Chinese Academy of Sciences, Beijing, China

12

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 on the atmosphere-like Lorenz-96 model with spatially correlated observations,  
15 and is shown to reduce the analysis error and increase its sensitivity to the  
16 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 “true” underlying states. Data assimilation is a powerful mechanism for  
5 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 assimilation provides an analysis state that is usually a better estimate of the  
8 state variable because it 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 accomplished with the rapid development of  
15 supercomputers (Reichle 2008), although finding the appropriate analysis state is a  
16 much difficult problem when the models are nonlinear.

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

1 the model forecasts and excludes the observations. It can eventually result in the  
2 divergence of the filter (Anderson and Anderson 1999; Constantinescu et al. 2007;  
3 Wu et al. 2014).

4 The covariance inflation technique is used to mitigate filter divergence by  
5 inflating the empirical covariance in EnKF, and it can increase the weight of the  
6 observations in the analysis state (Xu et al. 2013). In reality, this method will perturb  
7 the subspace spanned by the ensemble vectors and better capture the sub-growing  
8 directions that may **not have been captured by** the original ensemble (Yang et al.  
9 2015). Therefore, using the inflation technique to enhance the estimate accuracy of  
10 the forecast error covariance matrix is increasingly important.

11 A widely used inflation technique involves multiplying the forecast error matrix  
12 by **an** inflation factor, which must be chosen appropriately. In early studies,  
13 researchers usually tuned the inflation factor by repeated assimilation experiments  
14 and selected the estimated inflation factor according to their experience and prior  
15 knowledge (Anderson and Anderson 1999). However, such methods are very  
16 empirical and subjective. It **is not appropriate** to use the same inflation factor during  
17 all the assimilation procedure. Too small or too large an inflation factor will cause the  
18 analysis state to over rely on the model forecasts or observations, and can seriously  
19 undermine the accuracy and stability of the filter.

20 In later studies, the inflation factor **is** estimated online based on the innovation  
21 statistic (observation-minus-forecast; (Dee 1995; Dee and Silva 1999)) with different  
22 conditions. Moment estimation can facilitate the calculation by solving an equation of

1 the innovation statistic and its realization (Li et al. 2009; Miyoshi 2011; Wang and  
2 Bishop 2003). Maximum likelihood approach can obtain a better estimate of the  
3 inflation factor than moment approach, although it must calculate a high dimensional  
4 matrix determinant (Liang et al. 2012; Zheng 2009). Bayesian approach assumes a  
5 prior distribution for the inflation factor but is limited by spatially independent  
6 observational errors (Anderson 2007, 2009). This study seeks to address the  
7 estimation of the inflation factor from the perspective of cross validation (CV).

8       The concept of CV was first introduced for linear regressions (Allen 1974) and  
9 spline smoothing (Wahba and Wold 1975), and it represents a common approach that  
10 can be applied to estimate tuning parameters in generalized additive models,  
11 nonparametric regressions and kernel smoothing (Eubank 1999; Gentle et al. 2004;  
12 Green and Silverman. 1994; Wand and Jones 1995). Usually, the data are divided into  
13 subsets some of which are used for modelling and analysis while others for  
14 verification and validation. The most widely used technique removes only one data  
15 point and uses the remainder to estimate the value at this point to test the estimation  
16 accuracy, which is also called the leave-one-out cross validation (Gu and Wahba  
17 1991).

18       The basic motivation behind CV is to minimize the prediction error at the  
19 sampling points. The generalised cross validation (GCV) is a modified form of  
20 ordinary CV, that has been found to possess several favourable properties and is more  
21 popular for selecting tuning parameters (Craven and Wahba 1979). For instance, Gu  
22 and Wahba (1991) applied the Newton's method to optimize the GCV score with

1 multiple smoothing parameters in a smoothing spline model. Wahba (1995) briefly  
2 reviewed the properties of the GCV and conducted an experiment to choose  
3 smoothing parameters in the context of variational data assimilation schemes with  
4 numerical weather prediction models. Zheng and Basher (1995) also applied the GCV  
5 in a thin-plate smoothing spline model of spatial climate data to deal with South  
6 Pacific rainfalls.

7       Actually, the GCV criterion is based on a predictive mean-square error criterion  
8 that attempts to obtain a best estimate (Wahba et al. 1995). It has a rotation-invariant  
9 property that is relative to the orthogonal transformation of the observations and is a  
10 consistent estimate of the relative loss (Gu 2002). For the inverse problems in such  
11 fields as meteorological data assimilation, GCV method can choose parameters  
12 systematically by minimizing a given objective function that will improve the  
13 assimilation results. It can particularly select parameters that reflect not only  
14 measurement accuracies from different sources but also model capability (Krakauer  
15 et al. 2004).

16       This study proposes a new method for choosing the inflation factor using GCV  
17 method. The suitability of this choice is assessed using a statistic known as the  
18 analysis sensitivity, which apportions uncertainty in the output to different sources of  
19 uncertainty in the input (Saltelli et al. 2004; Saltelli et al. 2008). In the context of  
20 statistical data assimilation, this quantity describes the sensitivity of the analysis to  
21 the observations, which is complementary to the sensitivity of the analysis to model  
22 forecasts (Cardinali et al. 2004; Liu et al. 2009).

1 This study focuses on a methodology that can be potentially applied to  
 2 geophysical applications of data assimilation in the near future. This paper consists of  
 3 four sections. The conventional EnKF scheme is summarized and the improved EnKF  
 4 with GCV inflation scheme is proposed in Section 2; the verification and validation  
 5 processes are conducted on an idealized model in Section 3; the discussions are  
 6 presented in Section 4 and conclusions are given in Section 5.

7

8

## 9 **2. Methodology**

10

### 11 **2.1. EnKF algorithm**

12 For consistency, a nonlinear discrete-time dynamical forecast model and linear  
 13 observation system can be expressed as follows (Ide et al. 1997):

$$14 \quad \mathbf{x}_i^t = M_{i-1}(\mathbf{x}_{i-1}^a) + \boldsymbol{\eta}_i, \quad (1)$$

$$15 \quad \mathbf{y}_i^o = \mathbf{H}_i \mathbf{x}_i^t + \boldsymbol{\varepsilon}_i, \quad (2)$$

16 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 the  
 17  $n$ -dimensional true state vector at the  $i$ -th time step;  $\mathbf{x}_{i-1}^a = \{x_{i-1,1}^a, x_{i-1,2}^a, \dots, x_{i-1,n}^a\}^T$   
 18 represents the  $n$ -dimensional analysis state vector, which is an estimate of  $\mathbf{x}_{i-1}^t$ ;  
 19  $M_{i-1}$  represents a nonlinear dynamical forecast operator such as a numerical weather  
 20 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 observation  
 21 vector;  $\mathbf{H}_i$  represents the observation operator matrix; and  $\boldsymbol{\eta}_i$  and  $\boldsymbol{\varepsilon}_i$  represent  
 22 the forecast and observation error vectors, which are assumed to be time-uncorrelated,

1 statistically independent of each other and have mean zero and covariance matrices  
2  $\mathbf{P}_i$  and  $\mathbf{R}_i$ , respectively. The EnKF assimilation result is a series of analysis states  
3  $\mathbf{x}_i^a$  that is an accurate estimate of the corresponding true states  $\mathbf{x}_i^t$  based on the  
4 information provided by  $M_i$  and  $\mathbf{y}_i^o$ .

5 Suppose the perturbed analysis state at a previous time step  $\mathbf{x}_{i-1}^{a(j)}$  has been  
6 estimated ( $1 \leq j \leq m$  and  $m$  is the ensemble size), the detailed EnKF assimilation  
7 procedure is summarized as the following forecast step and analysis step (Burgers et  
8 al. 1998; Evensen 1994).

9 Step 1. Forecast step.

10 The perturbed forecast states are generated by running dynamical model  
11 forward:

$$12 \quad \mathbf{x}_i^{f(j)} = M_{i-1}(\mathbf{x}_{i-1}^{a(j)}). \quad (3)$$

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

$$16 \quad \mathbf{P}_i = \frac{1}{m-1} \sum_{j=1}^m (\mathbf{x}_i^{f(j)} - \mathbf{x}_i^f)(\mathbf{x}_i^{f(j)} - \mathbf{x}_i^f)^T. \quad (4)$$

17 Step 2. Analysis step.

18 The analysis state is estimated by minimizing the following cost function:

$$19 \quad J(\mathbf{x}) = (\mathbf{x} - \mathbf{x}_i^f)^T \mathbf{P}_i^{-1} (\mathbf{x} - \mathbf{x}_i^f) + (\mathbf{y}_i^o - \mathbf{H}_i \mathbf{x})^T \mathbf{R}_i^{-1} (\mathbf{y}_i^o - \mathbf{H}_i \mathbf{x}), \quad (5)$$

20 which has the analytic form

$$21 \quad \mathbf{x}_i^a = \mathbf{x}_i^f + \mathbf{P}_i \mathbf{H}_i^T (\mathbf{H}_i \mathbf{P}_i \mathbf{H}_i^T + \mathbf{R}_i)^{-1} \mathbf{d}_i, \quad (6)$$

22 where

$$1 \quad \mathbf{d}_i = \mathbf{y}_i^o - \mathbf{H}_i \mathbf{x}_i^f \quad (7)$$

2 is the innovation statistic (observation-minus-forecast residual in observation space).

3 To complete the ensemble forecast, the perturbed analysis states are calculated using  
4 perturbed observations (Burgers et al. 1998):

$$5 \quad \mathbf{x}_i^{a(i)} = \mathbf{x}_i^{f(i)} + \mathbf{P}_i \mathbf{H}_i^T (\mathbf{H}_i \mathbf{P}_i \mathbf{H}_i^T + \mathbf{R}_i)^{-1} (\mathbf{d}_i + \boldsymbol{\varepsilon}_i^{(i)}), \quad (8)$$

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

7 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

8 Sherman-Morrison-Woodbury formula (Golub and Loan 1996; Liang et al. 2012;

9 Tippett et al. 2003). Finally, set  $i = i + 1$ , return to Step 1 for the model forecast at

10 the next time step and repeat until the model reaches the last time step  $N$ .

11

## 12 **2.2. Influence matrix and forecast error inflation**

13 The forecast error inflation procedure should be added to any ensemble-based

14 assimilation scheme to prevent the filter from diverging (Anderson and Anderson

15 1999; Constantinescu et al. 2007). Multiplicative inflation is one of the commonly

16 used inflation techniques, and it adjusts the initially estimated forecast error

17 covariance matrix  $\mathbf{P}_i$  to  $\lambda_i \mathbf{P}_i$  after estimating the inflation factors  $\lambda_i$  properly.

18 In this study, a new procedure for estimating multiplicative inflation factors  $\lambda_i$

19 is proposed based on the following GCV function (Craven and Wahba 1979)

$$20 \quad GCV_i(\lambda) = \frac{\frac{1}{p_i} \mathbf{d}_i^T \mathbf{R}_i^{-1/2} (\mathbf{I}_{p_i} - \mathbf{A}_i(\lambda))^2 \mathbf{R}_i^{-1/2} \mathbf{d}_i}{\left[ \frac{1}{p_i} \text{Tr}(\mathbf{I}_{p_i} - \mathbf{A}_i(\lambda)) \right]^2}, \quad (9)$$

21 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

1 matrix of  $\mathbf{R}_i$ ; and

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

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

4 The inflation factor  $\lambda_i$  is estimated by minimizing the GCV (Eq. (9)) as an

5 objective function, and it is implemented between Steps 1 and 2 in Section 2.1. Then,

6 the perturbed analysis states are modified to

$$7 \quad \mathbf{x}_i^{a(i)} = \mathbf{x}_i^{f(i)} + \lambda_i \mathbf{P}_i \mathbf{H}_i^T \left( \mathbf{H}_i \lambda_i \mathbf{P}_i \mathbf{H}_i^T + \mathbf{R}_i \right)^{-1} \left( \mathbf{d}_i + \boldsymbol{\varepsilon}_i^{(i)} \right). \quad (11)$$

8 The flowchart of the EnKF equipped with the proposed forecast error inflation based

9 on the GCV method is shown in Figure 1.

10

### 11 **2.3. Analysis sensitivity**

12 In the EnKF, the analysis state (Eq. (6)) is a weighted average of the observation

13 and forecast. That is:

$$14 \quad \mathbf{x}_i^a = \mathbf{K}_i \mathbf{y}_i^o + (\mathbf{I}_n - \mathbf{K}_i \mathbf{H}_i) \mathbf{x}_i^f \quad (12)$$

15 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

16 matrix with dimension  $n \times n$ . Then, the normalized analysis vector can be expressed

17 as follows:

$$18 \quad \tilde{\mathbf{y}}_i^a = \mathbf{R}_i^{-1/2} \mathbf{H}_i \mathbf{K}_i \mathbf{R}_i^{1/2} \tilde{\mathbf{y}}_i^o + \mathbf{R}_i^{-1/2} \left( \mathbf{I}_{p_i} - \mathbf{H}_i \mathbf{K}_i \right) \mathbf{R}_i^{1/2} \tilde{\mathbf{y}}_i^f \quad (13)$$

19 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

20 observation space. The sensitivities of the analysis to the observation and forecast are

21 defined by Eq. (14) and (15), respectively:

$$22 \quad \mathbf{S}_i^o = \frac{\partial \tilde{\mathbf{y}}_i^a}{\partial \tilde{\mathbf{y}}_i^o} = \mathbf{R}_i^{1/2} \mathbf{K}_i^T \mathbf{H}_i^T \mathbf{R}_i^{-1/2}, \quad (14)$$



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

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

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

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

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

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

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 and exclude useful information from the observations. This is captured by the fact  
 6 that the GAI values are rather small for the conventional EnKF scheme. 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 for the  
 14 conventional EnKF, which also dramatically decreases until the observations  
 15 ultimately have an irrelevant impact on the analysis states. The inflation technique  
 16 can effectively compensate for the underestimation of the forecast ensemble spread,  
 17 and thereby can improve 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) is 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, [respectively](#). The model can be considered  
 3 representative of an atmospheric quantity (e.g., zonal wind speed) distributed on a  
 4 latitude circle. Therefore, the Lorenz-96 model has been widely used as a test bed to  
 5 evaluate the 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 and 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 by  
 15 adding random noises that were multivariate-normally distributed with mean zero and  
 16 covariance matrix  $\mathbf{R}_i$  to the true states. The frequency was every 4 time steps, which  
 17 can be used to mimic daily observations in practical problems, such as satellite data.  
 18 The observation errors were assumed to be spatially correlated, which is common in  
 19 applications involving remote sensing and radiance data. The variance of the  
 20 observation at each grid point was set to  $\sigma_o^2 = 1$ , and the covariance of the  
 21 observations between the  $j$ -th and  $k$ -th grid points was as follows:

$$22 \quad \mathbf{R}_i(j, k) = \sigma_o^2 \times 0.5^{\min\{|j-k|, 40-|j-k|\}}. \quad (20)$$

1

### 2 ***3.2. Assimilation scheme comparison***

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

11       The Lorenz-96 model was run for 2000 time steps, which is equivalent to  
12 approximately 500 days in realistic problems. The synthetic observations were  
13 assimilated at every grid point and every 4 time steps using the conventional EnKF,  
14 the constant inflated EnKF and the improved EnKF schemes for comparisons. The  
15 time series of estimated inflation factors are shown in Figure 2. It can be seen that,  
16 the estimated inflation factors vary between 1 and 6 in most instances, although the  
17 values smaller than 1 are estimated in several assimilation time steps. The median of  
18 the estimated inflation factors was 1.88, which was used as the inflation factor in the  
19 constant inflated EnKF scheme. Since the median is a robust and highly efficient  
20 statistic of the central tendency, this can ensure a relative fair comparison between the  
21 constant inflated EnKF and the improved EnKF schemes.

22       The forecast ensemble spread of the conventional EnKF, constant inflated

1 EnKF and improved EnKF are plotted in Figure 3. For the conventional EnKF,  
2 because the forecast states usually shrink together, the forecast ensemble spread was  
3 quite small and had a mean value of 0.36. The mean spread value of the improved  
4 EnKF was 3.32, which was larger than that of the constant inflated EnKF (3.25).  
5 These findings illustrate that the underestimation of forecast ensemble spread can be  
6 effectively compensated for by the two EnKF schemes with forecast error inflation  
7 and that the improved EnKF is more effective than the constant inflated EnKF.

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

15 To evaluate the analysis estimate accuracy, the analysis RMSE (Eq. (18)) and  
16 the corresponding values of the GCV functions (Eq. (9)) were calculated and plotted  
17 in Figures 5 and 6, respectively. The results illustrate that the analysis RMSE and the  
18 values of the GCV functions decrease sharply for the two EnKF with forecast error  
19 inflation schemes. However, the GCV function and the RMSE values of the improved  
20 EnKF were about 15% smaller than those of the constant inflated EnKF, indicating  
21 that the online estimate method performs better than the simple multiplicative  
22 inflation techniques with a constant value. The correlation coefficient of the analysis

1 RMSE and the value of the GCV function at the assimilation time step were  
2 approximately 0.76, which indicates that the GCV function is a good criterion to  
3 estimate the inflation factor.

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

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

18

### 19 ***3.3 Influence of ensemble size and observation number***

20 Intuitively, for any ensemble-based assimilation scheme, a large ensemble size  
21 will lead to small analysis errors; however, the computational costs are high for  
22 practical problems. The ensemble size in the practical land surface assimilation

1 problem is usually several tens of members (Kirchgessner et al. 2014). The  
2 preferences of the proposed inflation method and the constant inflation method with  
3 respect to different ensemble sizes (10, 30 and 50) were evaluated, and the results are  
4 listed in Table 1. It shows that for each scheme, using a 10-member ensemble  
5 produced a threefold increase in the analysis RMSE, while using a 50-member  
6 ensemble reduced the analysis RMSE by 20% relative to the analysis RMSE obtained  
7 using a 30-member ensemble. The forecast ensemble spread increased slightly from a  
8 10-member ensemble to a 50-member ensemble. The GAI and GCV function values  
9 changed sharply from a 10-member ensemble to a 30-member ensemble, and they  
10 became relatively stable from a 30-member ensemble to a 50-member ensemble.  
11 Ensembles less than 10 were unstable, and no significant changes occurred for  
12 ensembles greater than 50. Considering the computational costs for practical  
13 problems, a 30-member ensemble may be necessary for Lorenz-96 model to estimate  
14 statistically robust results. In the realistic problem, a system in which the errors grow  
15 in multiple directions will need more ensembles to produce statistically robust results.

16 To evaluate the preferences of the inflation method with respect to different  
17 numbers of observations, synthetic observations were generated at every other grid  
18 point and for every 4 time steps. Hence, a total of 20 observations were performed at  
19 each observation step in this case. The assimilation results with ensemble sizes of 10,  
20 30 and 50 are listed in Table 2, which shows that the GAI values were larger than  
21 those with 40-observations in all assimilation schemes. This finding may be related to  
22 the relatively small denominator of the GAI statistic (Eq. (16)) in the 20-observation



1 experiments. The forecast ensemble spread does not change much but the GCV  
2 function and the RMSE values increase greatly in the 20-observation experiments  
3 with respect to those in the 40-observation experiments, which illustrates that more  
4 observations will lead to less analysis error.

5

6

## 7 **4. Discussions**

8

### 9 *4.1 Performance of the GCV inflation*

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

18 The use of multiplicative covariance inflation techniques can mitigate this  
19 problem to some extent. Several methods have been proposed in the literature, and  
20 each has different assumptions. For instance, the moment approach can be easily  
21 conducted based on the moment estimation of the innovation statistic. The maximum  
22 likelihood approach can obtain a more accurate inflation factor than the moment

1 approach, but requires computing high dimensional matrix determinants. The  
2 Bayesian approach assumes a prior distribution for the inflation factor but is limited  
3 to spatially independent observational errors. In this study, the inflation factor was  
4 estimated based on cross-validation and the analysis sensitivity was detected. The  
5 estimated inflation factor by minimizing the GCV function is not affected by the  
6 observation unit and can optimize the analysis sensitivity to the observation.

7 In fact, the GCV method can evaluate and compare learning algorithms and  
8 represents a widely used statistical method. It can be applied in inverse problems in  
9 such fields as meteorological data assimilation (Wahba et al. 1995). Specifically,  
10 GCV provides a well-characterized method, which can select a regularization  
11 parameter by minimizing the predictive data errors with rotation-invariant in a least  
12 squares solution (MacCarthy et al. 2011). In data assimilation research fields,  
13 observation data such as in-situ observation and remote sensing data are usually from  
14 different sources. GCV is particularly useful for choosing relative parameters that  
15 reflect not only measurement accuracies from different sources but also model  
16 capability (Krakauer et al. 2004). Apparently, GCV method requires calculating the  
17 trace of a large matrix, which may be commonly computationally prohibitive for  
18 large inverse problems (MacCarthy et al. 2011).

19 In this study, the GCV concept was adopted for the inflation factor estimation  
20 in the improved EnKF assimilation scheme and was validated with the Lorenz-96  
21 model. The assimilation results showed that inflating the conventional EnKF using  
22 the factor estimated by minimizing the GCV function can indeed reduce the analysis

1 RMSE. Therefore, the GCV function can accurately quantify the goodness of fit of  
 2 the error covariance matrix. The values of the GCV function obviously decreased in  
 3 the proposed approach compared the conventional EnKF and constant inflated EnKF  
 4 schemes. The analysis RMSE of the proposed approach was also much smaller than  
 5 those of the conventional EnKF and constant inflated EnKF schemes, which suggests  
 6 that the GCV criterion works well for estimating the inflation factor.

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

14

#### 15 **4.2 Computational cost**

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

$$19 \quad GCV_i(\lambda) = \frac{p_i \mathbf{d}_i^T (\mathbf{H}_i \lambda \mathbf{P}_i \mathbf{H}_i^T + \mathbf{R}_i)^{-1} \mathbf{R}_i (\mathbf{H}_i \lambda \mathbf{P}_i \mathbf{H}_i^T + \mathbf{R}_i)^{-1} \mathbf{d}_i}{\left[ \text{Tr} \left( (\mathbf{H}_i \lambda \mathbf{P}_i \mathbf{H}_i^T + \mathbf{R}_i)^{-1} \mathbf{R}_i \right) \right]^2}. \quad (21)$$

20 Because both the numerator and denominator of the GCV function are scalars, the  
 21 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

1 calculated using the Sherman-Morrison-Woodbury formula. Furthermore, the inverse  
2 matrix calculation and the multiplication process are also indispensable for the  
3 conventional EnKF (Eq. (6)). Essentially, no additional computational burden is  
4 associated with the improved EnKF for the inverse matrix. Therefore, the total  
5 computational costs of the improved EnKF are feasible.

6 For the Lorenz-96 experiments in this study, the conventional EnKF, constant  
7 inflated EnKF and proposed improved EnKF assimilation schemes were conducted  
8 using R language on a computer with Intel Core i5 CPU and 8 GB RAM. The  
9 running times with different observation numbers and ensemble sizes were listed in  
10 Tables 1 and 2. It shows that for each assimilation scheme, the computational cost  
11 increases as the ensemble size grows. For the fixed observation number and ensemble  
12 size, the conventional EnKF, which does not involve the forecast error inflation, has  
13 the least running time but at a cost of losing assimilation accuracy. The proposed  
14 EnKF scheme is about 15% smaller in analysis RMSE, but only about 5% longer in  
15 running time than the constant inflated EnKF scheme. For the operational  
16 meteorological/ocean models, the most computational cost is in the ensemble model  
17 integrations (Ravazzani et al. 2016). Therefore, the proposed EnKF scheme does not  
18 significantly increase computational cost.

19

### 20 **4.3 Notes**

21 It is worth noting that the inflation factor is assumed to be constant in space in  
22 this study, which may be not the case in realistic assimilation problems. Forcing all

1 components of the state vector to use the same inflation factor could systematically  
2 overinflate the ensemble variances in sparsely observed areas, especially when the  
3 observations are unevenly distributed. In the presence of sparse observations, the  
4 state that is not observed can be improved only by the physical mechanism of the  
5 forecast model, although this improvement is limited. Therefore, a multiplicative  
6 inflation may not be sufficiently effective to enhance the assimilation accuracy. In  
7 this case, the additive inflation and the localization technique can be applied to  
8 further improve the assimilation quality in the presence of sparse observations  
9 (Miyoshi and Kunii 2011; Yang et al. 2015).

10

11

## 12 **5. Conclusions**

13

14 In this study, the approach for using GCV as a metric to estimate the covariance  
15 inflation factor was proposed. In the case studies conducted in Section 3, the  
16 observations were relatively evenly distributed and the assimilation accuracy could  
17 indeed be improved by the forecast error inflation technique. These findings provide  
18 insights on the methodology and validation of the Lorenz-96 model and illustrate the  
19 feasibility of our approach. In the near future, methods of modifying the adaptive  
20 procedure to suit the system with unevenly distributed observations and applying to  
21 more sophisticated dynamic and observation systems will be investigated.

22

1 **Appendix A**

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

3 
$$\tilde{\mathbf{y}}_i^o = \mathbf{R}_i^{-1/2} \mathbf{H}_i \mathbf{x}_i^t + \tilde{\boldsymbol{\varepsilon}}_i, \quad (\text{A1})$$

4 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  
 5 the identity matrix with the dimensions  $p_i \times p_i$ . Similarly, the normalized analysis  
 6 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  
 7 observation vector to the normalized analysis vector, thereby ignoring the normalized  
 8 forecast state in the observation space (Gu 2002):

9 
$$\tilde{\mathbf{y}}_i^a - \mathbf{R}_i^{-1/2} \mathbf{H}_i \mathbf{x}_i^f = \mathbf{A}_i \left( \tilde{\mathbf{y}}_i^o - \mathbf{R}_i^{-1/2} \mathbf{H}_i \mathbf{x}_i^f \right). \quad (\text{A2})$$

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

12 
$$\mathbf{A}_i = \mathbf{I}_{p_i} - \mathbf{R}_i^{1/2} \left( \mathbf{H}_i \mathbf{P}_i \mathbf{H}_i^T + \mathbf{R}_i \right)^{-1} \mathbf{R}_i^{1/2}. \quad (\text{A3})$$

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

15 
$$\mathbf{A}_i(\lambda) = \mathbf{I}_{p_i} - \mathbf{R}_i^{1/2} \left( \mathbf{H}_i \lambda \mathbf{P}_i \mathbf{H}_i^T + \mathbf{R}_i \right)^{-1} \mathbf{R}_i^{1/2}, \quad (\text{A4})$$

16 The principle of CV is to minimize the estimated error at the observation grid  
 17 point. Lacking an independent validation data set, a common alternative strategy is to  
 18 minimize the squared distance between the normalized observation value and the  
 19 analysis value while not using the observation on the same grid point, which is the  
 20 following objective function:

21 
$$V_i(\lambda) = \frac{1}{p_i} \sum_{k=1}^{p_i} \left( \tilde{\mathbf{y}}_{i,k}^o - \left( \mathbf{R}_i^{-1/2} \mathbf{H}_i \mathbf{x}_i^{a[k]} \right)_k \right)^2, \quad (\text{A5})$$

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

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

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

$$5 \quad V_i(\lambda) = \frac{1}{p_i} \sum_{k=1}^{p_i} \frac{(\tilde{\mathbf{y}}_{i,k}^o - (\mathbf{R}_i^{-1/2} \mathbf{H}_i \mathbf{x}_i^a)_k)^2}{(1 - a_{k,k})^2}, \quad (\text{A7})$$

6 where  $a_{k,k}$  is the element at the site pair  $(k, k)$  of the influence matrix  $\mathbf{A}_i(\lambda)$ . Then,  
7  $a_{k,k}$  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  
8 ignored to obtain the following GCV statistic (Gu 2002):

$$9 \quad \text{GCV}_i(\lambda) = \frac{\frac{1}{p_i} \mathbf{d}_i^\top \mathbf{R}_i^{-1/2} (\mathbf{I}_{p_i} - \mathbf{A}_i(\lambda))^2 \mathbf{R}_i^{-1/2} \mathbf{d}_i}{\left[ \frac{1}{p_i} \text{Tr}(\mathbf{I}_{p_i} - \mathbf{A}_i(\lambda)) \right]^2}. \quad (\text{A8})$$

10

## 11 **Appendix B**

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

$$13 \quad \mathbf{S}_i^o = \frac{\partial \tilde{\mathbf{y}}_i^a}{\partial \tilde{\mathbf{y}}_i^o} = \mathbf{R}_i^{1/2} \mathbf{K}_i^\top \mathbf{H}_i^\top \mathbf{R}_i^{-1/2}, \quad (\text{B1})$$

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

$$\begin{aligned} 15 \quad \mathbf{S}_i^o &= \mathbf{R}_i^{1/2} \mathbf{K}_i^\top \mathbf{H}_i^\top \mathbf{R}_i^{-1/2} \\ 16 \quad &= \mathbf{R}_i^{1/2} (\mathbf{H}_i \mathbf{P}_i \mathbf{H}_i^\top + \mathbf{R}_i)^{-1} \mathbf{H}_i \mathbf{P}_i \mathbf{H}_i^\top \mathbf{R}_i^{-1/2} \\ 17 \quad &= \mathbf{R}_i^{1/2} (\mathbf{H}_i \mathbf{P}_i \mathbf{H}_i^\top + \mathbf{R}_i)^{-1} (\mathbf{H}_i \mathbf{P}_i \mathbf{H}_i^\top + \mathbf{R}_i - \mathbf{R}_i) \mathbf{R}_i^{-1/2} \\ 18 \quad &= \mathbf{R}_i^{1/2} (\mathbf{H}_i \mathbf{P}_i \mathbf{H}_i^\top + \mathbf{R}_i)^{-1} (\mathbf{H}_i \mathbf{P}_i \mathbf{H}_i^\top + \mathbf{R}_i) \mathbf{R}_i^{-1/2} - \mathbf{R}_i^{1/2} (\mathbf{H}_i \mathbf{P}_i \mathbf{H}_i^\top + \mathbf{R}_i)^{-1} \mathbf{R}_i \mathbf{R}_i^{-1/2} \end{aligned}$$

$$\begin{aligned}
&= \mathbf{I}_{p_i} - \mathbf{R}_i^{1/2} \left( \mathbf{H}_i \lambda \mathbf{P} \mathbf{H}_i^T + \mathbf{R}_i \right)^{-1} \mathbf{R}_i^{1/2} \\
&= \mathbf{A}_i
\end{aligned} \tag{B2}$$

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

4

**Acknowledgements.** This work is supported by the National Natural Science Foundation of China (Grant No. 91647202), the National Basic Research Program of China (Grant No. 2015CB953703) and the National Natural Science Foundation of China (Grant No. 41405098). [The authors would like to gratefully acknowledge the two anonymous reviewers and the editor for their constructive comments, which helped significantly in improving the quality of this manuscript.](#)

11



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, as well as the running times  
4 (second) for different assimilation schemes. The observation number is 40 and the  
5 ensemble size is selected as 10, 30 and 50, respectively.

6

Scheme	Ensemble Size	Spread	GAI	GCV	RMSE	Running Time
Conventional EnKF	10	0.23	4.56%	36.38	4.50	70.73
	30	0.36	10.78%	31.14	4.01	215.92
	50	0.41	13.58%	25.21	3.52	346.69
Constant inflated EnKF	10	3.15	4.78%	35.91	4.38	77.41
	30	3.25	27.48%	5.56	1.41	238.25
	50	3.27	19.67%	5.03	1.14	384.63
Improved EnKF	10	3.26	5.24%	35.56	3.74	81.31
	30	3.32	29.21%	3.29	1.10	251.06
	50	3.45	35.63%	2.30	0.88	405.68

7

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

2

Scheme	Ensemble Size	Spread	GAI	GCV	RMSE	Running Time
Conventional EnKF	10	0.41	10.77%	33.64	4.85	67.75
	30	0.59	20.92%	22.89	4.10	181.27
	50	0.68	26.41%	14.97	3.29	295.92
Constant inflated EnKF	10	3.03	11.73%	33.39	4.64	71.22
	30	3.18	30.07%	17.12	3.92	203.64
	50	3.27	39.51%	12.74	3.37	322.29
Improved EnKF	10	3.33	13.25%	32.17	4.39	74.84
	30	3.36	35.09%	14.99	3.46	213.81
	50	3.48	41.28%	5.19	2.86	339.41

3

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. The median of the estimated inflation factors is 1.88.

5 Figure 3. Forecast ensemble spread of the conventional EnKF (black line), the  
6 constant inflated EnKF (red line) and the improved EnKF (blue line) for the  
7 Lorenz-96 experiment with 40-observation and 30-ensemble member. The constant  
8 multiplicative inflation factor is set as 1.88.

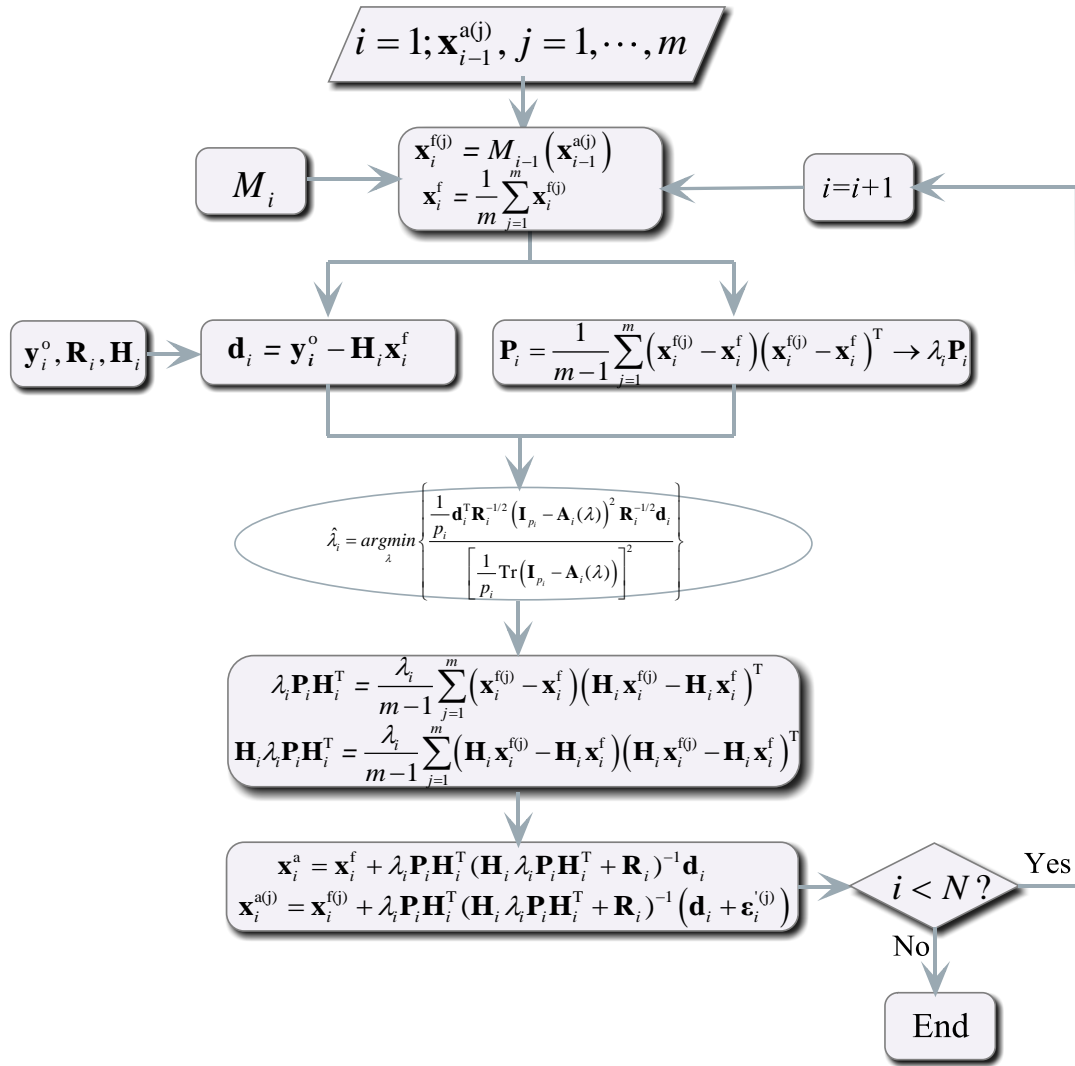
9 Figure 4. GAI statistics of the conventional EnKF (black line), the constant inflated  
10 EnKF (red line) and the improved EnKF (blue line) for the Lorenz-96 experiment  
11 with 40-observation and 30-ensemble member. The constant multiplicative inflation  
12 factor is set as 1.88.

13 Figure 5. Analysis RMSE of the conventional EnKF (black line), the constant inflated  
14 EnKF (red line) and the improved EnKF (blue line) for the Lorenz-96 experiment  
15 with 40-observation and 30-ensemble member. The constant multiplicative inflation  
16 factor is set as 1.88.

17 Figure 6. GCV function values of the conventional EnKF (black line), the constant  
18 inflated EnKF (red line) and the improved EnKF (blue line) for the Lorenz-96  
19 experiment with 40-observation and 30-ensemble member. The constant  
20 multiplicative inflation factor is set as 1.88.

21 Figure 7. Ensemble analysis state members of the conventional EnKF (black line), the  
22 constant inflated EnKF (red line) and the improved EnKF (blue line) for the

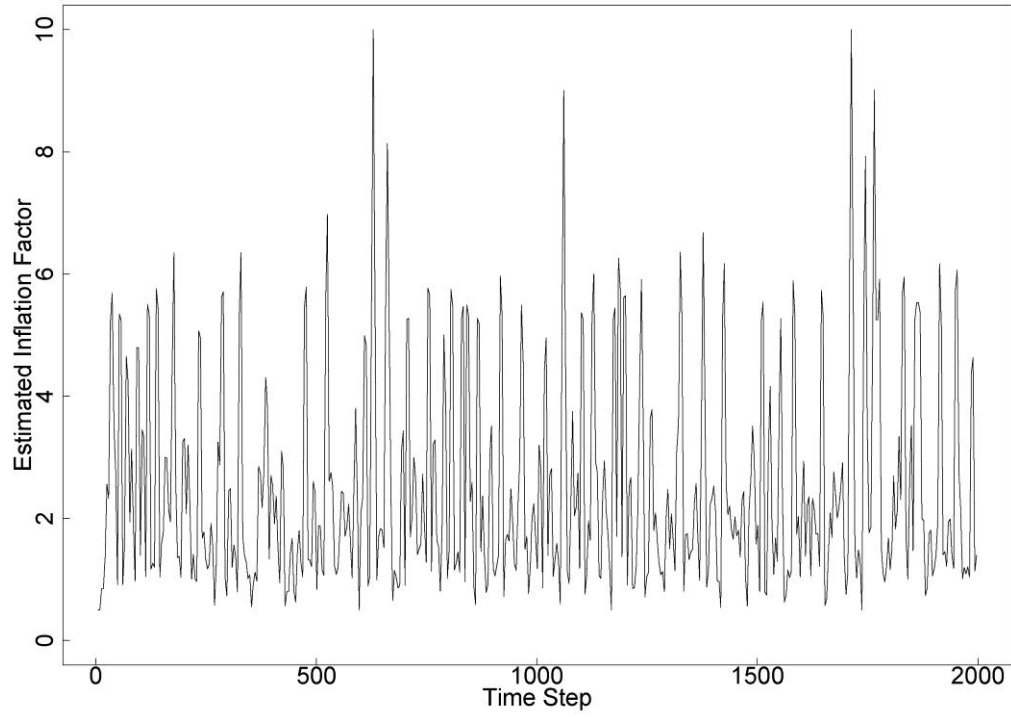
- 1 Lorenz-96 experiment with 40-observation and 30-ensemble member. The constant
- 2 multiplicative inflation factor is set as 1.88. The green line refers to the true trajectory
- 3 obtained by the numerical solution.
- 4



1

2 Figure 1. Flowchart of the proposed assimilation scheme.

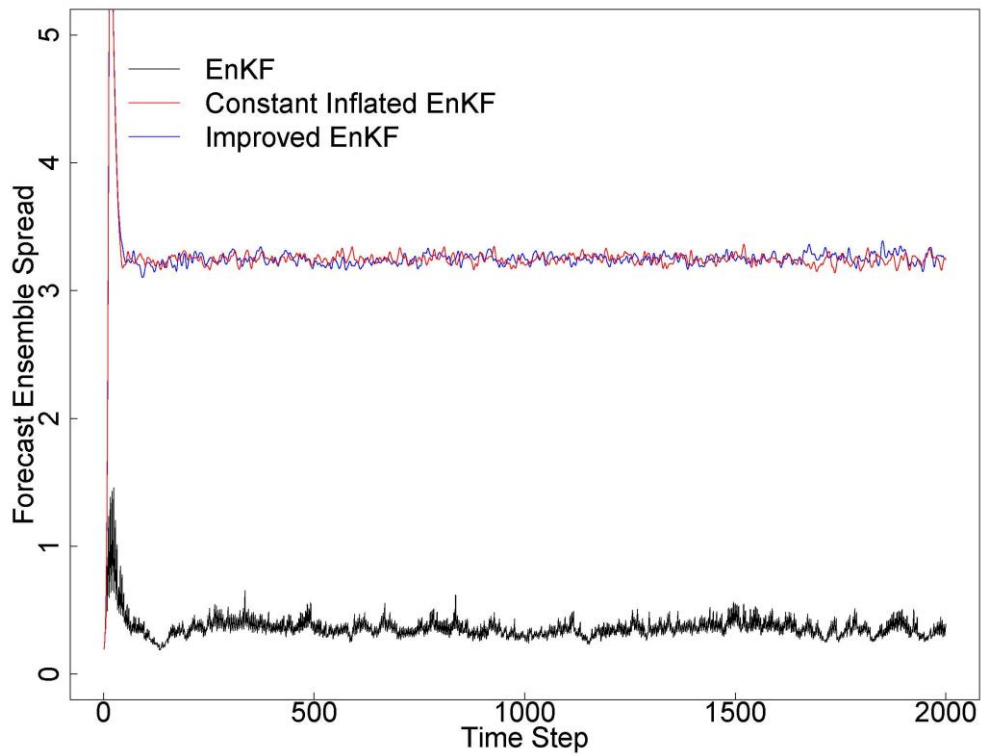
3



1

2 Figure 2. Time series of the estimated inflation factors by minimizing the GCV  
3 function. The median of the estimated inflation factors is 1.88.

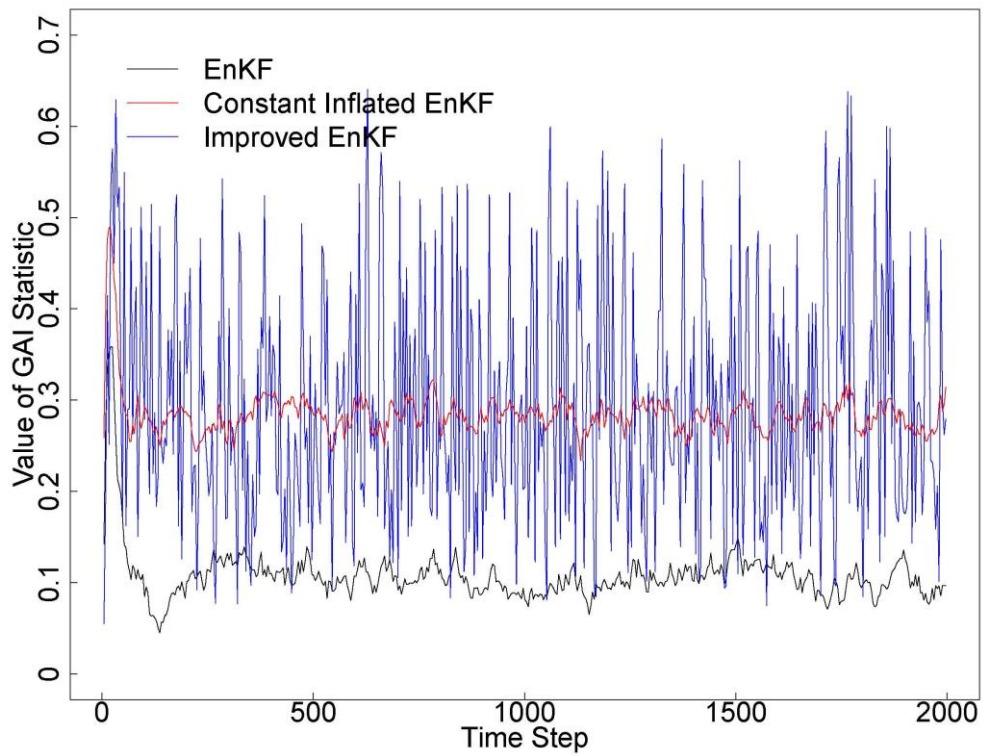
4



1

2 Figure 3. Forecast ensemble spread of the conventional EnKF (black line), the  
3 constant inflated EnKF (red line) and the improved EnKF (blue line) for the  
4 Lorenz-96 experiment with 40-observation and 30-ensemble member. The constant  
5 multiplicative inflation factor is set as 1.88.

6

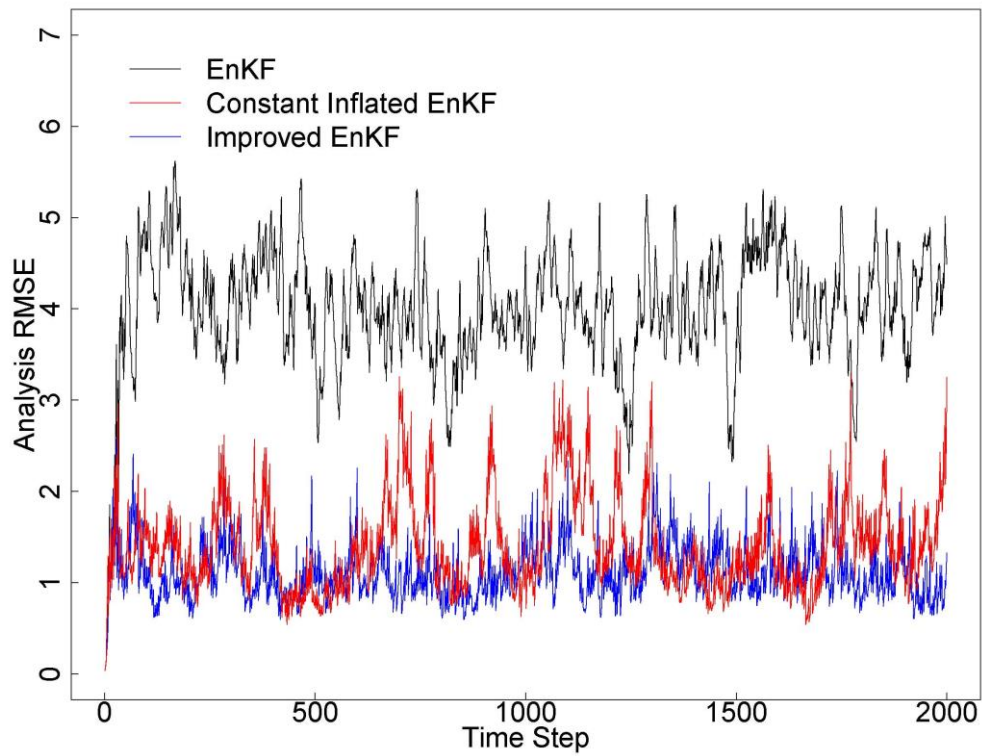


1

2 Figure 4. GAI statistics of the conventional EnKF (black line), the constant inflated  
3 EnKF (red line) and the improved EnKF (blue line) for the Lorenz-96 experiment  
4 with 40-observation and 30-ensemble member. The constant multiplicative inflation  
5 factor is set as 1.88.

6

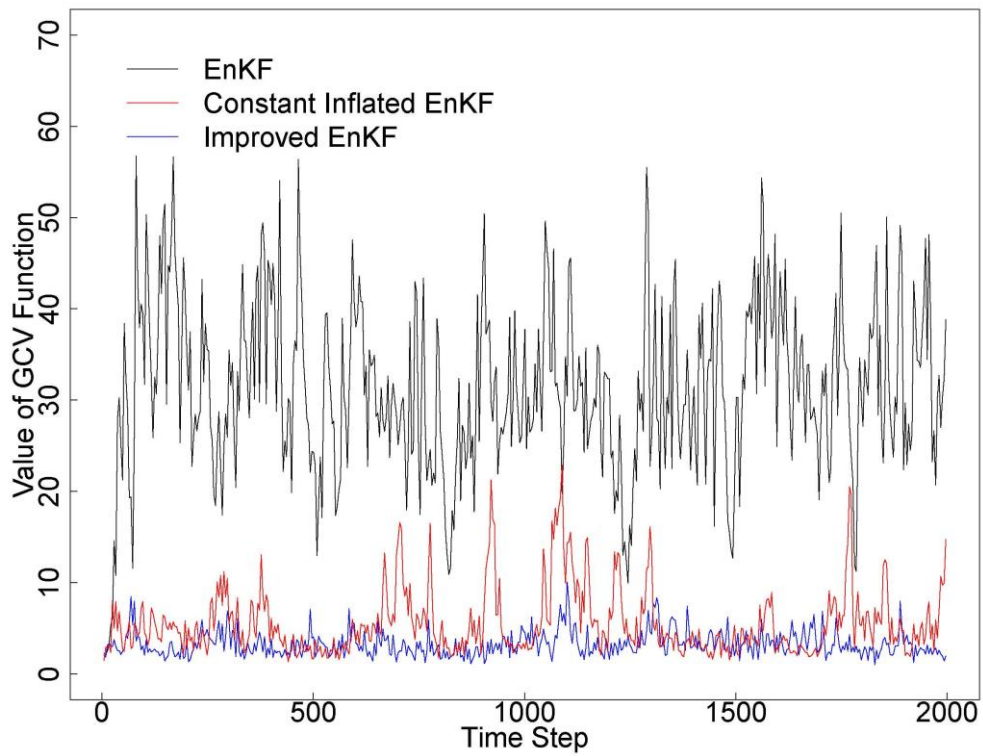




1

2 Figure 5. Analysis RMSE of the conventional EnKF (black line), the constant inflated  
3 EnKF (red line) and the improved EnKF (blue line) for the Lorenz-96 experiment  
4 with 40-observation and 30-ensemble member. The constant multiplicative inflation  
5 factor is set as 1.88.

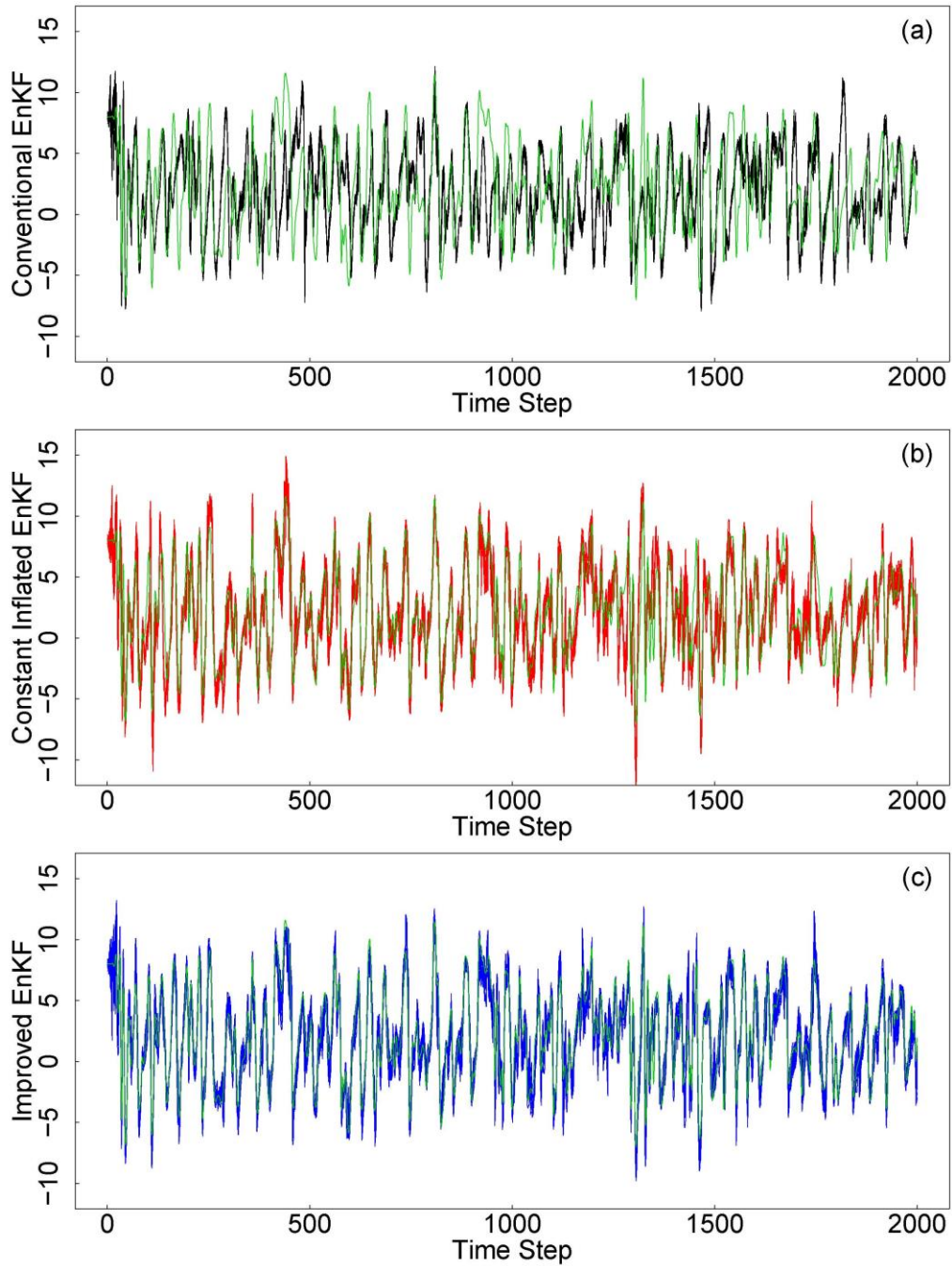
6



1

2 Figure 6. GCV function values of the conventional EnKF (black line), the constant  
 3 inflated EnKF (red line) and the improved EnKF (blue line) for the Lorenz-96  
 4 experiment with 40-observation and 30-ensemble member. The constant  
 5 multiplicative inflation factor is set as 1.88.

6



1

2 Figure 7. Ensemble analysis state members of the conventional EnKF (black line), the  
 3 constant inflated EnKF (red line) and the improved EnKF (blue line) for the  
 4 Lorenz-96 experiment with 40-observation and 30-ensemble member. The constant  
 5 multiplicative inflation factor is set as 1.88. The green line refers to the true trajectory  
 6 obtained by the numerical solution.

7

## 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 Krakauer, N. Y., T. Schneider, J. T. Randerson, and S. C. Olsen, 2004: Using generalized cross-validation

1 to select parameters in inversions for regional carbon fluxes. *Geophysical Research Letters*, **31**.

2 Li, H., E. Kalnay, and T. Miyoshi, 2009: Simultaneous estimation of covariance inflation and  
3 observation errors within an ensemble Kalman filter. *Quarterly Journal of the Royal Meteorological*  
4 *Society*, **135**, 523-533.

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

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

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

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

13 MacCarthy, J. K., B. Borchers, and R. C. Aster, 2011: Efficient stochastic estimation of the model  
14 resolution matrix diagonal and generalized cross-validation for large geophysical inverse problems.  
15 *Journal of Geophysical Research*, **116**.

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

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

20 Miyoshi, T., and M. Kunii, 2011: The Local Ensemble Transform Kalman Filter with the Weather  
21 Research and Forecasting Model: Experiments with Real Observations. *Pure & Applied Geophysics*,  
22 **169**, 321-333.

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

25 Ravazzani, G., A. Amengual, A. Ceppi, V. Homar, R. Romero, G. Lombardi, and M. Mancini, 2016:  
26 Potentialities of ensemble strategies for flood forecasting over the Milano urban area. *Journal of*  
27 *Hydrology*, **539**, 237-253.

28 Reichle, R. H., 2008: Data assimilation methods in the Earth sciences. *Advances in Water Resources*,  
29 **31**, 1411-1418.

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

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

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

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

37 Wahba, G., and S. Wold, 1975: A completely automatic french curve. *Communications in Statistics*, **4**,  
38 1-17.

39 Wahba, G., D. R. Johnson, F. Gao, and J. Gong, 1995: Adaptive tuning of numerical weather prediction  
40 models randomized GCV in three- and four-dimensional data assimilation. *Monthly Weather Review*,  
41 **123**, 3358-3369.

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

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

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

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

7 Xu, T., J. J. Gómez-Hernández, H. Zhou, and L. Li, 2013: The power of transient piezometric head data  
8 in inverse modeling: An application of the localized normal-score EnKF with covariance inflation in a  
9 heterogenous bimodal hydraulic conductivity field. *Advances in Water Resources*, **54**, 100-118.

10 Yang, S.-C., E. Kalnay, and T. Enomoto, 2015: Ensemble singular vectors and their use as additive  
11 inflation in EnKF. *Tellus A*, **67**.

12 Zheng, X., 2009: An adaptive estimation of forecast error statistic for Kalman filtering data  
13 assimilation. *Advances in Atmospheric Sciences*, **26**, 154-160.

14 Zheng, X., and R. Basher, 1995: Thin-plate smoothing spline modeling of spatial climate data and its  
15 application to mapping south Pacific rainfall. *Monthly Weather Review*, **123**, 3086-3102.

16

17