# Using orthogonal vectors to improve the ensemble space of the EnKF and its effect on data assimilation and forecasting

Yung-Yun Cheng[1], Shu-Chih Yang[2,3], Zhe-Hui Lin[2], and Yung-An Lee[2]
[1]Center for Weather Climate and Disaster Research, National Taiwan University, Taipei, Taiwan
[2]Department of Atmospheric Sciences, National Central University, Taoyuan, Taiwan
[3]GPS Science and Application Research Center, National Central University, Taoyuan, Taiwan

*Correspondence to*: Dr. Shu-Chih Yang (shuchih.yang@atm.ncu.edu.tw)

**Abstract.** The space spanned by the background ensemble provides a basis for correcting forecast errors in the ensemble Kalman filter. However, the ensemble space may not fully capture the forecast errors due to the limited ensemble size and systematic model errors, which affect the assimilation performance. This study proposes a new algorithm to generate pseudo members to properly expand the ensemble space during the analysis step. The pseudomembers adopt vectors orthogonal to the original ensemble and are included in the ensemble using the centered spherical simplex ensemble method. The new algorithm is investigated with a six-member ensemble Kalman filter implemented in the Lorenz 40-variable model. Our results suggest that orthogonal vectors with the ensemble singular vector or ensemble mean vector can serve as effective pseudo members for improving the analysis accuracy, especially when the background has large errors.

## 1 Introduction

The ensemble Kalman filter (EnKF) has the great advantage of using flow-dependent background error covariance (BEC) and has been widely applied to state estimation in geophysics. The BEC is sampled by the background ensemble, and its characteristic is crucial since it determines how the observations are spread out to correct the model state. The space spanned by the background ensemble members (the ensemble space) is expected to capture the dynamically growing errors, and the ensemble space provides a basis for corrections.

However, the use of a finite ensemble size can cause the ensemble space to be underestimated; hence, the growing errors may not be well captured, and corrections from the EnKF are less optimal. Bocquet and Carrassi (2017) indicated that the stability of the EnKF relies on the subspace spanned by the ensemble members that represent the unstable-neutral subspace. In other words, maintaining the ensemble space is important for the EnKF performance. Strategies such as additive covariance inflation (Whitaker et al., 2008) or hybrid methods (Hamill and Snyder, 2000) are commonly used to increase the dimensionality of the ensemble. These methods expand the overall ensemble space but are operated empirically without a particular direction.

Previous studies have suggested that vectors stimulate the growing modes can improve the dimensionality of the ensemble space and the performance of the EnKF. For example, Carrassi et al. (2008) used bred vectors as the direction of the analysis increment to update the analysis states in an unstable area. Yang et al. (2015) used a two-sided method to apply the initial ensemble singular vectors (IESV) as additive inflation to correct the fastest-growing errors in a quasi-geostrophic model.

Chang et al. (2020), under the framework of a hybrid-gain data assimilation framework (Penny, 2014), used part of the variational information orthogonal to the EnKF analysis perturbation to correct the EnKF means. These studies emphasize on the importance of generating additional effective correction. Inspired by these works, this study proposes generating

35 pseudoensemble members to increase the ensemble space to better capture forecast errors without increasing the computational cost. We investigate whether the use of pseudomembers could improve the analysis and forecast and which type of pseudomember is most effective in this regard.

This paper is organized as follows. Section 2 introduces the generation of pseudomembers. Section 3 presents the impact of using pseudomembers for EnKF analysis. Section 4 provides the summary and discussion of this work.

## 40   2 Methodology and experimental design

### 2.1 General setup

This study uses the 40-variable Lorenz-96 model (Lorenz 1996, Lorenz and Emanuel 1998)) to conduct data assimilation experiments (See Appendix A for the detailed procedure). Observation simulation system experiments are conducted with the local ensemble transform Kalman filter algorithm (LETKF; Hunt et al., 2007). An observation is provided every two grid

45 points with an observation error variance of 1.0. Table 1 lists the details of the assimilation parameters.

Before performing data assimilation, a new vector is included as the extra ensemble member. The traditional double-sized method (Toth and Kalnay 1993) requires an even number of included members and can lead to ill-conditioned problems during EnKF computation. Therefore, we adopt the centered spherical simplex ensemble (CSSE; Wang et al., 2004) method, which can add any number of members without modifying the ensemble mean and spread. More importantly, the CSSE method

50 avoids ill-conditioned problems.

### 2.2 Deriving the vectors for pseudomembers

An added member is referred to as a pseudomember given that it is generated at the analysis time and is not used during the forecast stage. Three types of vectors are used for generating pseudomembers, including the initial ensemble singular vector (IESV), ensemble mean vector (EMV), and random singular vector (RSV). Given a set of ensemble forecasts, IESV finds fast-

55 growing perturbations within a period by linearly combining the ensemble perturbations (Enomoto et al., 2015, Yang et al., 2015). In ensemble data assimilation, EMV is used to define the ensemble perturbations (as the deviation) and is not accounted for in the degrees of freedom, although the perturbations evolve upon the mean state. However, it is likely that the forecast errors carry a component with the structure of the mean, such as the large-scale pattern, and will be not be represented in the ensemble perturbations. Therefore, we propose using the EMV to expand the ensemble space. The RSV is generated by

60 randomly selecting a vector from the null singular vectors, which are orthogonal to the original ensemble space but with zero singular values.

Nonlinear Processes
in Geophysics
Discussions

After the vector is generated, its component orthogonal to the ensemble space is added as the pseudomember for EnKF computation. To obtain the orthogonal component, first, we apply SVD to the ensemble space to orthogonalize the ensemble. Second, we use Equations (1) and (2) to obtain the orthogonal component of the generated vector (orthogonal IESV1 or orthogonal EMV).

$$\mathbf{v}_{\text{final proj}} = \sum_{i=1}^{K-1} \frac{(\tilde{\mathbf{v}} \cdot \mathbf{v}_i)}{|\mathbf{v}_i|^2} \mathbf{v}_i \qquad (1)$$

$$\mathbf{v}_{\text{orth}} = \text{normalize}(\tilde{\mathbf{v}} - \mathbf{v}_{\text{final proj}}) \qquad (2)$$

where $\tilde{\mathbf{v}}$ is the normalized generated vector, $\mathbf{v}_i$ is the $i^{\text{th}}$ orthogonal vector of the ensemble space, $\mathbf{v}_{\text{final proj}}$ is the projection of the generated vector at the ensemble space, and $\mathbf{v}_{\text{orth}}$ is the orthogonal component of the generated vector. Finally, the orthogonal component of the generated vector is taken as the new ensemble perturbation for expanding the ensemble space. The orthogonal vector is rescaled to have the amplitude of the background ensemble spread. Different experiments are designed. The control experiment (CNTL) conducts standard LETKF assimilation with six members and is taken as the baseline. We conduct four experiments with the added pseudomembers. The first three experiments use the RSV, the orthogonal components of global IESV1 and EMV. The last one adds two orthogonal vectors from IESV1 and EMV.

## 2. 3 Setup of the increased-size EnKF system

Figure 1 shows the flow chart of our experiments. With the $K$-member background ensemble, $M$-member pseudovectors (RSV, global IESV1 or global EMV) are generated. With CSSE, the ensemble size becomes $(K+M)$, and the LETKF analysis is performed with the new ensemble. Based on Fig. 1, we conducted our experiments with offline and online frameworks. The offline framework, in which the LEKTF analysis is not cycled, is used to investigate how the ensemble space varies after adding the pseudovector and to understand the benefits of the increased-size EnKF system by clean comparisons. The background ensemble is provided by the background ensemble of the CNTL at each analysis step. In contrast, the analysis is cycled in the online experiment framework to evaluate the accumulated feedback from the increased-size EnKF system. However, $K$ members need to be selected from the new $(K+M)$ members so that the following ensemble forecast is done without the need for extra computational costs. To do so, we remove the last $M$ members with Equation 3 to keep the ensemble mean and the ensemble spread the same as when using the $(K+M)$ members.

$$\mathbf{v}_i^{new} = \bar{\mathbf{v}}_i|_{(K+M)} + \left(\frac{1}{K}\sum_{n=K+1}^{K+M} \mathbf{v}'_n + \mathbf{v}'_i\right)\left(\frac{\sigma_{K+M}}{\sigma_K}\right) \qquad i = 1,2 \dots K \qquad (3)$$

where $\mathbf{v}_i^{new}$ is the new $i^{\text{th}}$ member, $\bar{\mathbf{v}}_i|_{(K+M)}$ is the mean of the $(K+M)$ members, $\mathbf{v}'_i$ is the $i^{\text{th}}$ member perturbation, and $\sigma_{K+M}$ and $\sigma_K$ represent the ensemble spread of the $(K+M)$ and $K$ members, respectively. In Equation (3), $\frac{1}{K}\sum_{n=K+1}^{K+M} \mathbf{v}'_n$ is used to ensure that the sum of the new perturbations of the first $K$ members is equal to zero. In this study, $M$ is set to one or two as a demonstration of proof of concept.

## 3 Results

### 3.1 Results of the offline experiments

95 This subsection illustrates how including the orthogonal vector modifies the ensemble space. The results are obtained from the offline setting, in which the orthogonal vector is directly included in the background ensemble of the standard LETKF experiment (CNTL) without cycling the impact. In this study, the orthogonal vector is computed and added globally, and can enhance the local BEC matrix, the basis for local analysis correction. Here, we focus on the area with the maximum forecast error. A local maximum forecast error (LME) area is defined as a local area spanning seven grids, whose center has the largest

100 forecast error.

The characteristics of the ensemble space in the LME area is represented by the eigen spectrum of the local ensemble perturbation. Figure 2 shows the percentage of eigenvalues averaged for 550 DA cycles. The ensemble space of CNTL (black) has 5 nonzero eigenvalues provided 6 members. The newly added orthogonal vector successfully provides an independent mode into the ensemble space. Moreover, the orthogonal vectors from both the IESV1 (Fig. 2, orange) and the EMV (Fig. 2,

105 green) provide a larger expansion of the ensemble space than the RSV (Fig. 2, gray) in the LME area. This indicates that the orthogonal IESV1 and orthogonal EMV can contribute an additional independent mode. Additionally, the 8-member (Fig. 2, blue) experiment, which includes two orthogonal vectors, can increase two independent modes, so the ensemble space is expanded into seven modes.

The expansion of the ensemble space can further modify the analysis correction. This can be illustrated by the projection of

110 forecast errors onto the ensemble space and the improvement in the analysis errors with the additional orthogonal vector in the LME area. The calculation of the projection is similar to Eq. 2, except $\tilde{\mathbf{v}}$ is replaced with the error of the background ensemble mean and the amplitude of $\mathbf{v}_{\text{final proj}}$ is calculated at each analysis cycle in each experiment. We compare how well the modification of the ensemble space can help to capture the error of the background ensemble mean. The projection of CNTL (Fig. 3b, black) decreases when the LETKF performs poorly in reducing the background error (Fig. 3a). This also confirms

115 that LETKF assimilation is less successful when the ensemble cannot capture the forecast error well. All experiments with the additional orthogonal vector can increase the projection, i.e., the ensemble can better capture the forecast error and provide a more effective correction to reduce the error (Fig. 3c), especially at the analysis times when the CNTL analysis errors are larger than the background errors (highlighted by the red dashed boxes in Fig. 3). The orthogonal EMV is most effective in increasing the projection of the forecast errors among the three orthogonal vectors. Furthermore, the 8-member experiment

120 has the best performance in capturing forecast errors and reducing analysis errors. Note that the projection of the forecast errors only illustrates how well the ensemble space encompasses the forecast errors. It should be noted that the forecast errors in the LME area, spanning 7 grids, can be fully represented by 7 orthogonal vectors (the 8-member experiment) (Fig. 3b, blue), and the background errors may not be completely removed due to issues such as the underrepresentation of background error variance. Nevertheless, the offline experiments confirm the potential of adding the orthogonal vector to provide more effective

125 corrections, and the improvement is highly flow dependent.

Figure 4 sorts the projection of the CNTL background ensemble on the background error according to the CNTL analysis RMSE and the projection of the added orthogonal vector on the forecast error residual. The forecast error residual is the unexplained forecast errors after removing the projection of the original background ensemble from the original forecast error. All calculations are performed on the whole domain. First, the large CNTL analysis RMSE corresponds to the low

130  projection of the original background ensemble on forecast errors (i.e., a large RMSE vs. the blue dots). In general, the larger the CTRL RMSE is, the higher the orthogonal vector capturing the residual error. Moreover, compared to the orthogonal IESV1, the orthogonal EMV projects more onto the forecast error residual at most analysis times, indicating that it is more useful to increase the ensemble space to reduce the analysis errors. This justifies why adding the orthogonal EMV is most effective in improving the performance of LETKF assimilation.

135  **3.2 Results of the online experiments**

In this subsection, we compare the results of the online experiments, in which the impact of using the additional pseudovector is cycled during the analysis step and further feedbacks into the next background ensemble through analysis cycling. To highlight this accumulated impact, we evaluate the analysis error based on different criteria with the RMSE of the analysis ensemble mean. The first group includes the analyses at all analysis cycles for comparing the general performance of

140  all experiments. The second and third groups focus on the analysis cycles, whose analysis RMSE values in the CNTL are between one and two standard deviations and two standard deviations larger than the mean RMSE, respectively. These two groups highlight the effect of adding the new vectors in LETKF assimilation when the original background ensemble is less capable of reducing forecast errors.

First, we confirm that adding a pseudomember is useful in improving the analysis accuracy in general (Fig. 5a).

145  Among the four types of orthogonal vectors, adding the orthogonal EMV is most effective when the standard LETKF assimilation has poor performance (Fig. 5c). The analysis RMSE is reduced by approximately 55%, and the 30 time-step forecast RMSE is reduced by approximately 38%. The advantage of using the EMV is more significant in the online experiment than in the offline experiment, indicating that the additional correction is beneficial for correcting the growing error and thus results in positive feedback. Notably, adding the IESV1 without orthogonalization still improves the analysis with large errors

150  (Figs. 5b and 5c). In other words, in the situation that the original ensemble space cannot provide effective correction, adding the IESV1 is still useful for expanding the ensemble space to better capture the forecast error. The reason that adding the IESV1 without orthogonalization degrades the general performance is attributed to overcorrections at the cycles when the ensemble space is already good at representing the forecast error.

We provide further comparison experiments using one or two orthogonal vectors with the standard LETKF with 7 members.

155  While using more members introduces more computation during the ensemble forecast, experiments using the pseudovectors for assimilation do not have extra computation for performing ensemble forecasting. The 8-member experiment successfully combines the advantages of using the orthogonal IESV1 and the orthogonal EMV, and thus, it has a better performance than both single-pseudomember experiments, especially for groups with large analysis errors (Fig. 6b,c). More importantly, the 8-

member experiment has a comparable performance with the 7-member standard LETKF in general (Fig. 6a) and even
160    outperforms in the group with mildly large analysis errors (Fig. 6b). However, when the analysis error grows to a certain range,
the 7-member standard LETKF experiment has the best performance. Such an improvement is still valid and even more evident
when the model is imperfect (Fig. 6d-f). (See Appendix B for the detailed procedure.)

**4 Summary and conclusion**

This study proposed a new idea of adding cost-free pseudomembers, which are orthogonal to the original ensemble space, to
165    expand the ensemble space and to improve the performance of EnKF assimilation and forecast. Based on the Lorenz-96 model,
this idea is investigated with offline and online framework. Three types of pseudoensemble members are compared, including
the RSV, global orthogonal IESV1 and orthogonal EMV. Compared to the RSV, both the IESV1 and EMV are very effective
in expanding the ensemble space in sensitive areas while the orthogonal EMV is most effective in improving the analysis
accuracy. This implies that the large-scale errors in the ensemble mean may not be well captured by the ensemble perturbations.
170    Furthermore, adding more orthogonal vectors can further improve the analysis and forecast accuracy. The 8-member analysis,
which uses two additional cost-free orthogonal vectors during LETKF assimilation, has an accuracy comparable to that of the
standard 7-member LETKF analysis, which has a higher computational cost. The conclusion holds in imperfect model settings
as well.

In the current operational ensemble DA system, additive covariance inflation is commonly adopted to expand the ensemble
175    space to improve the analysis corrections. However, how to choose additive inflation is ad hoc and may even break the balance
of the structure of the original ensemble space and degrade the performance of the data assimilation system. The newly
proposed simple idea could be a gentle alternative. In the future, we plan to explore the feasibility of this idea in a real model;
in particular when the forecast errors highly depend on the dynamics of the mean state and including the orthogonal EMV as
pseudomember for ocean ensemble data assimilation.

**Appendix A: Introduction of Lorenz-96 model**

180

The Lorenz-96 model has been used to study the issue of error growth and the probability of atmosphere and weather
forecasting (Kekem et al. 2018). The governing equation of the Lorenz-96 model is:

$$\frac{dx_j}{dt} = x_{j-1}(x_{j+1} - x_{j-2}) - x_j + F \qquad j = 1, \dots, n \qquad \text{(A.1)}$$

where $n$ is the dimension of the system, j is the index of the analysis grids, and $F$ is the external forcing term. $\frac{dx_j}{dt}$ can be
185    interpreted as some atmospheric quantity measured along the same latitude of the earth (Lorenz, 2006a).

**Appendix B: The extra supplement of imperfect model experiments.**

The imperfect model online experiments are conducted as follows. We change the value of $F$ used in the governing equations of the Lorenz-96 model from 8 to 7.8, and the inflation in each experiment is 2.3 (except 1.9 for in the 7-member standard
190 LETKF experiment). The orthogonal EMV experiment shows similar performance as the orthogonal IESV1 experiment, so we only display the result of the orthogonal EMV experiment. The orthogonal EMV experiment shows that the ensemble space can capture more forecast errors and improve the performance of EnKF assimilation. Furthermore, the 8-member experiment combines the advantages of both the orthogonal IESV1 and orthogonal EMV experiments, it performs better than either of them, and it performs similarly to the 7-member standard LETKF experiment.

195

**Author contributions**

**Competing interests**

The authors declare that they have no conflict of interest.

205 **Acknowledgments**

**Code and data availability**

210 Details of the experiment data could be reproduced by the code. The control run and code will be released at Github (https://github.com/yungyun0721/orthogonal_vector_code) after this manuscript is accepted.

**References**

Bocquet, M., and A. Carrassi: Four-dimensional ensemble variational data assimilation and the unstable subspace. Tellus A: Dynamic Meteorology and Oceanography, 69(1), 1304504, 2017.

215 Carrassi, A., A. Trevisan, L. Descamps, O. Talagrand, and F. Uboldi: Controlling instabilities along a 3DVar analysis cycle by assimilating in the unstable subspace: A comparison with the EnKF, Nonlinear Process. Geophys, 15, 503-521, 2008.

Chang, S. G. Penny, and S.-C. Yang: Hybrid Gain Data Assimilation using Variational Corrections in the Subspace Orthogonal to the Ensemble. Mon. Wea. Rev., 148, 2331–2350, doi:10.1175/MWR-D-19-0128.1, 2020.

Enomoto, T., Yamane, S. and Ohfuchi, W.: Simple sensitivity using ensemble forecasts. J. Meteorol. Soc. Jpn. 93, 199-213.
220   DOI: 10.2151/jmsj.2015-011, 2015.

Hamill, T. M., and C. Snyder: A Hybrid Ensemble Kalman Filter–3D Variational Analysis Scheme. Mon. Weather Rev., 128, 2905–2919, doi:10.1175/1520-0493(2000)1282.0.CO;2, 2000.

Hunt, B. R., and E. J. Kostelich, and I. Szunyogh: Efficient data assimilation for spatiotemporal chaos: A local ensemble Kalman filter, Physica D, 230, 112–126, 2007.

225   Kekem, D. L. van, and A. E. Sterk: Dynamics of the Lorenz-96 model: Bifurcations, symmetries and waves. Rijksuniversiteit Groningen, 2018.

Lorenz, E. N.: Predictability: A problem partly solved. In Seminar on Predictability, volume Vol. I, ECMWF, Reading, UK, 1–18, 1996.

Lorenz, E. N., and K. A. Emanuel: Optimal sites for supplementary weather observations: Simulation with a small model. J.
230   Atmos.Sci., 55, 399–41, 1998.

Lorenz, E.N.: 'Predictability—A Problem partly solved', in: Palmer, T.N. & Hagedorn, R. (eds.), Predictability of Weather and Climate, Cambridge: Cambridge University Press, chap. 3, pp. 40–58, 2006a.

Penny, S. G: The hybrid local ensemble transform Kalman filter. Mon. Wea. Rev., 142, 2139–2149, doi:10.1175/MWR-D-13-00131.1, 2014.

235   Toth, Z., and Kalnay, E.: Ensemble forecasting at NMC: The generation of perturbations. Bull. Amer. Meteor. Soc.,74, 2317–2330, 1993.

Wang, X., C. H. Bishop, and S. J. Julier: Which is better, an ensemble of positive-negative pairs or a centered spherical simplex ensemble?: Monthly Weather Review, 132, 1590–1605, 2004.

Whitaker, J. S., Hamill, T. M., Wei, X., Song, Y. and Toth, Z.: Ensemble data assimilation with the NCEP Global Forecast
240   System. Mon. Weather Rev. 136, 463-481, doi: https://doi.org/10.1175/2007MWR2018.1, 2008.

Yang, S.-C., E. Kalnay, and T. Enomoto: Ensemble singular vecotrs and their use as additive inflation in ENKF. Tellus A, 67, 26536, doi: http://dx.doi.org/10.3402/tellusa.v67.26536, 2015.

| Model settings | |
|---|---|
| **Model** | Lorenz 96 (N=40) |
| **Observation** | Every 2 grid points |
| **DA system** | LETKF |
| **Truncation localization** | 45 degrees (5 grid points) |
| **R-localization** | 12.5 degree |
| **Multiplicative inflation** | 1.8 |

| DA interval | Every 30 time steps |
|---|---|
| Total time step | 600 DA cycles |
| Ensemble size | 6 members |

245    **Table 1 The settings of the model, observations, data assimilation and ensemble.**
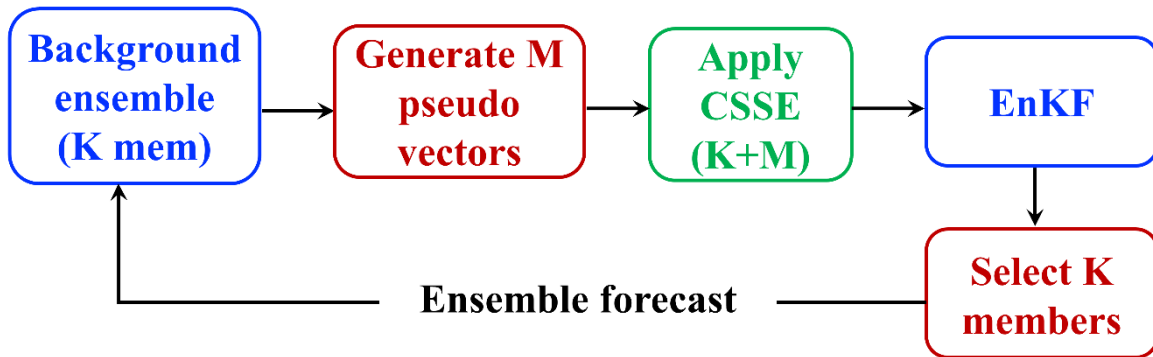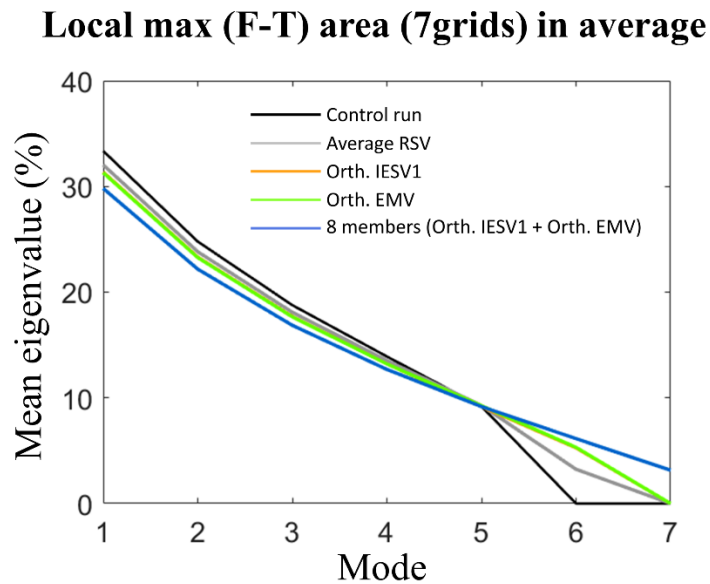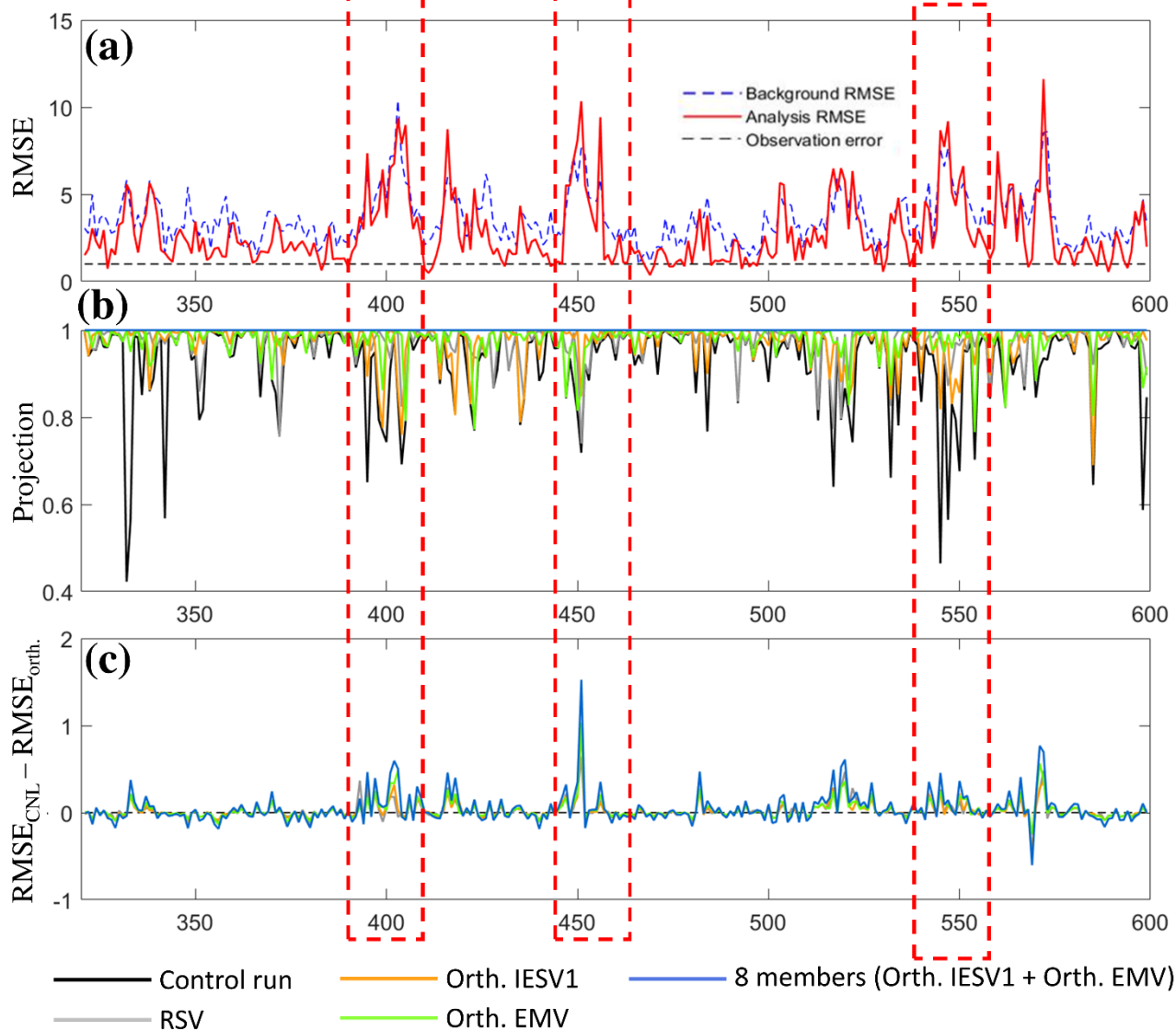


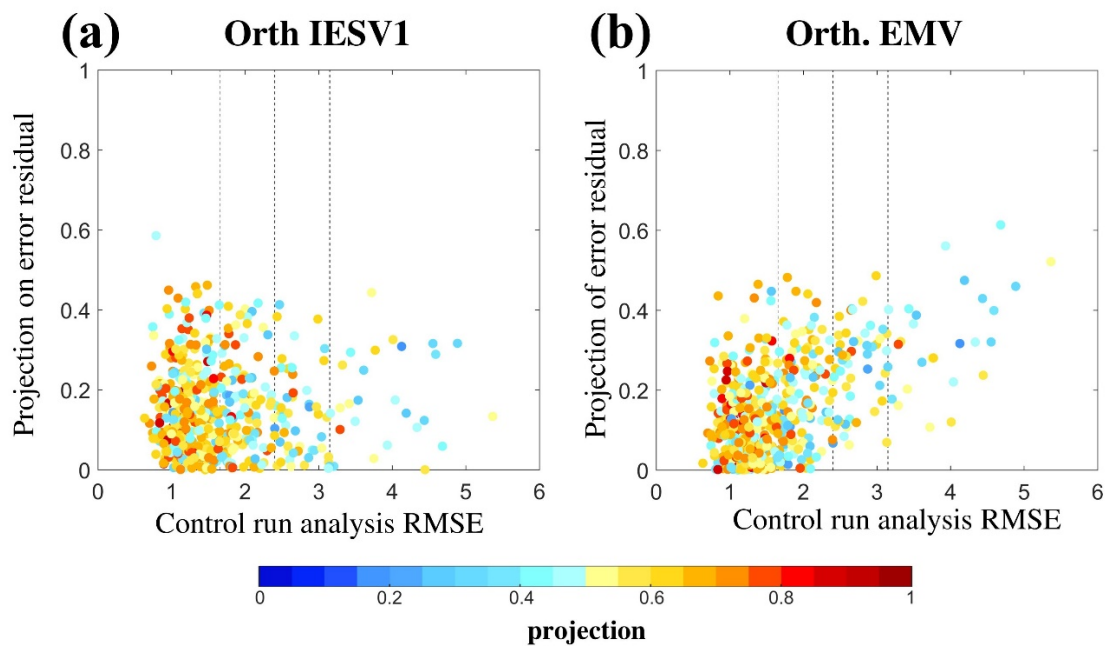Figure 1: The framework of the increased-size EnKF system



250    **Figure 2: The mean eigenvalue percentage (Y-axis) of the control run (black), the average RSV (gray), the orthogonal IESV1 (orange), the orthogonal EMV (green), and the 8-member experiments (blue) in each eigenmode (X-axis). The calculation is done using the ensemble perturbation in the local area centered at the grid with the maximum forecast error.**
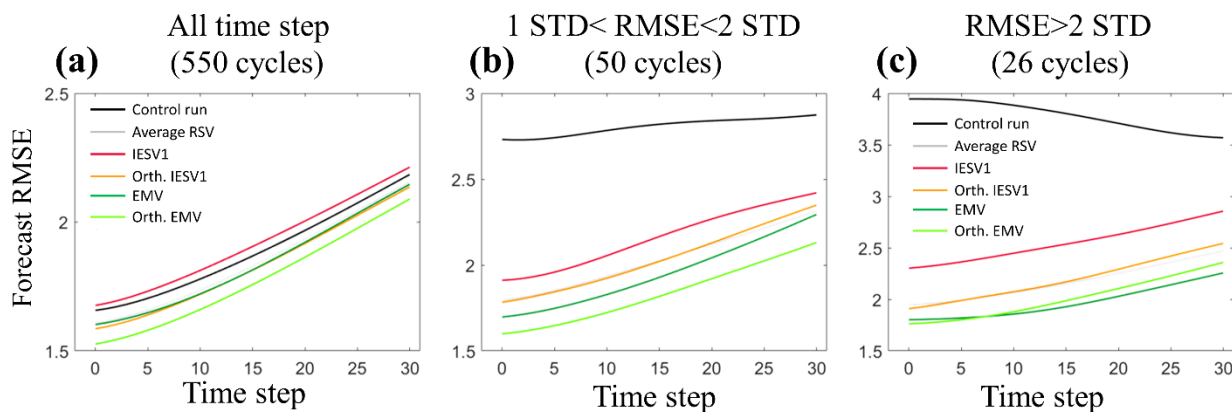
Figure 3: Time series (X-axis) of (a) the background (blue dashed line) and analysis (red) RMS errors, (b) the projection on the background error, and (c) the RMSE differences between the CNTL and experiments using pseudomember (as in Figure 2).

**Figure 4: Scatterplots of the projection on the background error (shading) according according to, the control analysis RMSE (X-axis) and the projection of the added pseudomember on the forecast errors residual: (a) orthogonal IESV1 and (b) orthogonal EMV experiments.**

**Figure 5: In the online experiment, the patterns of the mean analysis and forecast RMS errors (Y-axis) in each time step (X-axis) with one pseudovector experiment: the control run (black), the average RSV (gray), the IESV1 (red), the orthogonal IESV1 (orange), the EMV (dark green), and the orthogonal EMV (light green). (a) is in the all-time step, (b) considers the larger error time steps of the control run's analysis RMSE values, and (c) includes the largest error time steps of the control run's analysis RMSE values.**
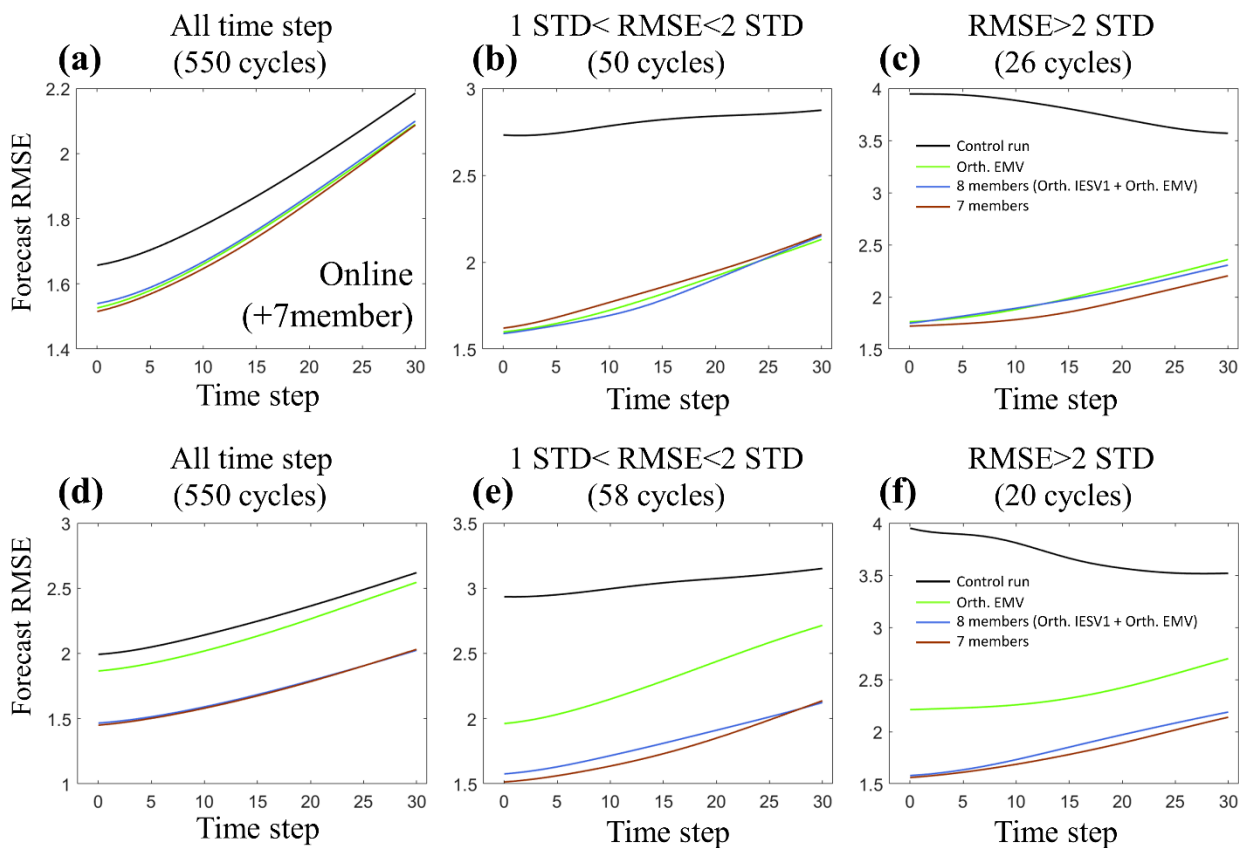
270

**Figure 6: As in Figure 5, except for different experiments: the 8-member experiment (blue), the 7-member standard LETKF experiment (brown); (a)-(c) the perfect model experiment and (d)-(e) the imperfect model experiment.**

275