

1  
2 **An Ensemble Recentering Kalman Filter with an Application to Argo**  
3 **Temperature Data Assimilation into the NASA GEOS-5 Coupled Model**  
4

5  
6 Christian L. Keppenne<sup>1,2</sup>  
7

8 *<sup>1</sup>Global Modeling and Assimilation Office*  
9 *Code 610.1, NASA Goddard Space Flight Center, Greenbelt, MD 20771 USA*  
10

11 *<sup>2</sup>Science Systems and Applications inc.*  
12 *10210 Greenbelt Road, Suite 600, Lanham, Maryland 20706, USA*  
13  
14

15  
16  
17 **Correspondence address:**

18 Christian Keppenne

19 email: christian.keppenne@nasa.gov

20 telephone: 011-1-301-6145874

21 mail: Code 610.1. NASA Goddard Space Flight Center, Greenbelt, Maryland 20771, USA  
22

23  
24 **October 14, 2013**  
25

26

27

### Abstract

28 A two-step ensemble recentering Kalman filter (ERKF) analysis scheme is introduced. The  
29 algorithm consists of a recentering step followed by an ensemble Kalman filter (EnKF) analysis  
30 step. The recentering step is formulated such as to adjust the prior distribution of an ensemble of  
31 model states so that the deviations of individual samples from the sample mean are unchanged  
32 but the original sample mean is shifted to the prior position of the most likely particle, where the  
33 likelihood of each particle is measured in terms of closeness to a chosen subset of the  
34 observations. The computational cost of the ERKF is essentially the same as that of a same size  
35 EnKF.

36

37 The ERKF is applied to the assimilation of Argo temperature profiles into the OGCM component  
38 of an ensemble of NASA GEOS-5 coupled models. Unassimilated Argo salt data are used for  
39 validation. A surprisingly small number (16) of model trajectories is sufficient to significantly  
40 improve model estimates of salinity over estimates from an ensemble run without assimilation.  
41 The two-step algorithm also performs better than the EnKF although its performance is degraded  
42 in poorly observed regions.

43

### Keywords:

45 Data assimilation; Kalman filter; ensemble Kalman filter; Particle filter;

46 Coupled data assimilation

## 1. Introduction

Since its introduction by *Evensen* [1994, 1996] in the context of a quasigeostrophic ocean model, the ensemble Kalman filter (EnKF) has gained wide acceptance among the atmosphere and ocean modeling communities as a viable data assimilation technique. In an EnKF, the prohibitive cost associated with evolving the model-background error-covariance matrix according to the traditional Kalman filter [*Kalman*, 1960] formulation is avoided. Rather, the background-error covariance propagation is replaced with the concurrent integration of an ensemble of model trajectories. The forecast-error covariance statistics needed to compute the Kalman gain are then estimated from the statistical distribution of the ensemble of model states. The EnKF uses the ensemble mean to estimate the true state of the dynamical system. The posterior ensemble variance serves to estimate the analysis error variance.

The sample mean is the most likely forecast if the ensemble is normally distributed. A normal distribution is often assumed in climate system data analysis, either for convenience or for lack of a better assumption. However, it is known from nonlinear dynamical system theory and confirmed from observations that clustering around several likely forecast solutions can occur. For example, the Kuroshio is in a multiple equilibrium state in which meandering and straight paths are equally plausible but the intermediate mean path is unstable. Because of this, the central forecast (defined as the ensemble member closest to the mean in terms of RMS difference) is sometimes used to estimate the true state.

Another problem associated with using the sample mean to estimate the true state is that while the individual ensemble members are dynamically balanced states, their mean generally is not. In a multivariate EnKF in which certain prognostic model variables are updated as observations of other variables are assimilated, the resultant analysis increments can lead to one or more ensemble members becoming dynamically unstable when the analyzed ensemble is advanced further in time.

Particle filters are a class of estimation methods that do not estimate the true state with the sample mean. Instead, the true-state estimates obtained with particle filters are statistically plausible, dynamically balanced states. The original particle filter algorithm, sample importance resampling [SIR: *Gordon et al.*, 1993], approximates the true state with a weighted sum of particles (ensemble members) where the weights are proportional to the respective probability of each particle. In nonlinear systems, SIR can estimate the true system state more accurately than the EnKF when the number of particles is large enough. Nevertheless, for the small ensemble sizes typically used in GCM ensemble prediction, *i.e.* a few tens to a few hundred ensemble members, the EnKF is generally more accurate [*Weerts and El Serafy*, 2006]. Besides, the weighted-particle SIR estimate is just as likely to be unstable as the EnKF ensemble mean. If, however, the most likely particle is used to estimate the state, the estimate will generally be stable and balanced. Note however that an even larger number of particles will be required for such a scheme to be accurate.

The purpose of this note is to introduce a two-step ensemble recentering Kalman filter (ERKF) analysis procedure that combines advantages from the EnKF and SIR schemes. It is shown that the ERKF can produce a more accurate state estimate than a same-size EnKF with dramatically less particles than a typical SIR filter would require. Our demonstration uses the NASA Global

Modeling and Assimilation Office (GMAO) global earth observing system (GEOS) integrated ocean data assimilation system (iODAS) applied to the assimilation of Argo temperature profiles into the OGCM component of the GEOS-5 coupled system. The ERKF is introduced in Section 2. The next two Sections contain overviews of the CGCM and of the data assimilation system. The validation experiments are discussed in Section 5. Our conclusions follow in Section 6.

## 2. The Two-Step Analysis Procedure

Let us consider a nonlinear model,  $\mathbf{M}$ , that evolves an approximation,  $\mathbf{x}$ , of the true system state,  $\mathbf{x}_t$ , subject to a forcing term,  $\mathbf{f}$ . We define an observation operator,  $\mathbf{H}$ , that maps the true system state to an observation vector,  $\mathbf{y}$ . Then, the system that defines the estimated state evolution can be written:

$$\frac{d\mathbf{x}}{dt} = \mathbf{M}(\mathbf{x}, \mathbf{f}) + \mathbf{q}, \quad (1a)$$

$$\mathbf{y} = \mathbf{H}(\mathbf{x}_t) + \mathbf{r}, \quad (1b)$$

$$E((\mathbf{x} - \mathbf{x}_t)(\mathbf{x} - \mathbf{x}_t)^T) = \mathbf{P}, \quad E(\mathbf{q}\mathbf{q}^T) = \mathbf{Q}, \quad E(\mathbf{r}\mathbf{r}^T) = \mathbf{R}, \quad (1c)$$

where  $E$  is the usual expectation operator.

We want to optimally estimate  $\mathbf{x}_t$  given information about  $\mathbf{x}$ ,  $\mathbf{y}$  and their error distributions. The Kalman filter [Kalman, 1960] tells us that the estimate that minimizes the analysis error variance is

$$\mathbf{x}^a = \mathbf{x}^f + \mathbf{P}\mathbf{H}^T(\mathbf{H}\mathbf{P}\mathbf{H}^T + \mathbf{R})^{-1}(\mathbf{y} - \mathbf{H}\mathbf{x}^f), \quad (2)$$

where  $\mathbf{x}^f$  is the prior state estimate (model forecast).

In an EnKF,  $\mathbf{P}$  is estimated from the statistical distribution of an ensemble of model forecasts,  $\mathbf{x}_i, i = 1, \dots, n$ , each evolved according to (1a), so that

$$\mathbf{S} = \{\mathbf{x}_1 - \bar{\mathbf{x}}, \dots, \mathbf{x}_n - \bar{\mathbf{x}}\}, \quad (3a)$$

$$\mathbf{P} \approx \frac{1}{n-1} \mathbf{S}\mathbf{S}^T, \quad (3b)$$

where the overbar denotes the sample mean. The analysis for  $\mathbf{x}_i$  follows from substituting  $\mathbf{x}_i$  for  $\mathbf{x}^f$  in (2).

The first (recentering) step of the ERKF works by finding the most likely sample,  $\mathbf{x}_p$ , such that

$$\sigma_p = \sum_k \left( \frac{y_k - \mathbf{H}_k \mathbf{x}_p}{r_k} \right)^2, \quad (4)$$

is minimized, where the  $k$  subscript runs over a chosen subset of the latest batch of observations, and  $r_k$  and  $H_k$  are the measurement error and observation operator for the  $k$ th datum. A pre-analysis increment,  $\Delta_p = \mathbf{x}_p - \bar{\mathbf{x}}$ , is then applied to each ensemble member. This places the ensemble mean at the original “location” of  $\mathbf{x}_p$  without affecting our estimates of  $\mathbf{P}$  and  $\mathbf{R}$  since the sample distribution relative to its mean is not changed (Figure 1).

The second step uses (2) and (3) to compute the usual EnKF analysis from the ensemble of  $\mathbf{x}'_i = \mathbf{x}_i + \Delta_p, i = 1, \dots, n$ , thus producing a set of analysis increments,

$$\Delta_i = \mathbf{P}\mathbf{H}^T (\mathbf{H}\mathbf{P}\mathbf{H}^T + \mathbf{R})^{-1} (\mathbf{y} - \mathbf{H}\mathbf{x}'_i + \mathbf{r}_i), \quad i = 1, \dots, n, \quad (5)$$

so that the total increment applied to ensemble member  $i$  is  $\Delta_p + \Delta_i$ .

Note that the observation vector used to identify  $\mathbf{x}_p$  could be different from that used in the Kalman filter update (5).

### 3. The GEOS-5 Coupled Model

The GEOS-5 CGCM couples the NASA GEOS-5 atmospheric model [Rienecker *et al.*, 2008] and GFDL’s Modular Ocean Model version 4 (MOM4) and the Los Alamos CICE sea ice model [Hunke *et al.*, 2013] using the Earth System Modeling Framework [ESMF: Hill *et al.*, 2006]. Ocean-atmosphere coupling is done through a physical interface layer of constant depth (2m). Atmospheric fluxes of heat and fresh water are applied at the top of this layer and oceanic turbulent fluxes are parameterized at the bottom of the layer. The mass of the layer and amount of mixing at its base are chosen so that the layer simulates the diurnal cycle of SST.

The GEOS-5 configuration used here runs MOM4 and CICE on a tripolar grid with uniform  $1/2^\circ$  zonal resolution, variable meridional resolution ranging from  $1/6^\circ$  in the Tropics to  $1/2^\circ$  in the high latitudes, and has 40 vertical levels. The dimension of the grid is  $720 \times 410 \times 40$  (zonal $\times$ meridional $\times$ vertical). The AGCM has uniform  $5/4^\circ$  zonal resolution,  $1^\circ$  meridional resolution and 72 vertical levels ( $288 \times 181 \times 72$  grid dimension). The system is run with a 15-minute time step. A faster, 15-second barotropic time step is also used to run MOM4.

The CGCM system has close to 100-million prognostic state variables in the configuration used herein. The ocean data assimilation modifies oceanic fields of temperature (T), salinity (S), zonal current (U), meridional current (V) and sea surface height (SSH) as well as sea ice concentration and thickness.

### 4. GEOS iODAS

The GEOS integrated ocean data assimilation system (iODAS) system has evolved from the first generation GMAO ocean data assimilation system introduced originally in Keppenne and Rienecker [2003] and discussed in more detail in Keppenne *et al.* [2005, 2008]. iODAS is implemented as an ESMF [Hill *et al.*, 2004] gridded component. The communications between the ocean and sea ice models and iODAS are managed by ESMF. The system is used routinely to produce the GMAO production ocean analysis [Vernieres *et al.*, 2012]. What follows is a summary of the main differences between the analysis algorithms used here and in Keppenne *et al.* [2008].

178

179 As is customary in EnKF applications to large-scale atmospheric or ocean models, the  
180 background-error covariances are localized to address the degree-of-freedom limitations  
181 encountered when the sample size is much less than that of the model state vector. The error-  
182 covariance localization is flow adaptive (following neutral density surfaces) and an iterative  
183 procedure is used to individually optimize the covariance localization scales involved in the  
184 processing of each observation. Incremental analysis updating [IAU: *Bloom et al.*, 1996] is used  
185 to apply them gradually over the assimilation interval.

186

187 The ocean assimilation is applied to the full CGCM, with the atmospheric component  
188 constrained during the integration by “replaying” the NASA Modern Era Retrospective-analysis  
189 for Research and Applications [MERRA: e.g., *Rienecker et al.*, 2008] into the GEOS-5 AGCM.  
190 The procedure involves integrating the AGCM to the next synoptic analysis time, reading the  
191 MERRA analysis fields and calculating analysis increments by taking their difference from the  
192 background atmospheric fields, rewinding the AGCM and, finally, integrating the full CGCM  
193 while incrementally applying both the atmospheric analysis increments thus computed and the  
194 ocean analysis increments produced by iODAS. The implementation is designed to facilitate  
195 consistent atmosphere-ocean states to initialize GEOS-5 seasonal climate forecasts. More details  
196 about the replay procedure are available in *Vernieres et al.* [2012].

197

## 198 5. Experiments and Results

199 To test the ERKF, four experiments were run spanning March-June 2006 and using 16 model  
200 trajectories. The initial condition comes from integrating a single instance of the CGCM while  
201 constraining its AGCM component by replaying the MERRA reanalysis. From this initial  
202 condition, the ensemble is first spun up during March 2006 by perturbing the analysis increments  
203 of the atmospheric analysis replay procedure and also applying daily perturbations to the OGCM  
204 T and salinity S fields. The ensemble configuration as of March 31 2006 is then used to  
205 initialize each ensemble run.

206

207 CE-16 is a control ensemble run without assimilation. ER-16 only applies the ensemble  
208 recentering step, EnKF-16 only applies the EnKF analysis step and ERKF-16 applies both the  
209 recentering and EnKF analysis steps. Starting on April 1, 2006, Argo T profiles are assimilated  
210 daily in EnKF-16 and in ERKF-16. These same Argo T profiles are also used to select  $\mathbf{x}_p$  in the  
211 recentering steps of ER-16 and ERKF-16. The average wallclock run times per simulation day  
212 on 960 2.8 GHz Intel Altix Sandy Bridge cores are 687 seconds in CE-16, 691 seconds in ER-16,  
213 744 seconds in EnKF-16 and 749 seconds in ERKF-16, illustrating the minimal cost of the  
214 recentering step. The performance of the data assimilation and of the ensemble recentering is  
215 assessed from how closely each run can reproduce the assimilated Argo T profiles and the  
216 unassimilated Argo S profiles.

217

218 The time mean April-June 2006 global RMS OMF differences of each run with the Argo T and S  
219 data are listed in table 1. They show that applying the recentering step without following it with  
220 an EnKF analysis has only a small impact on the RMS T OMF (7% reduction from CE-16 to ER-  
221 16) and that there is no noticeable benefit for T in applying the recentering step on top of the  
222 EnKF analysis step (1% RMS T OMF increase from EnKF-16 to ERKF-16). However, the  
223 results are markedly different in terms of the unassimilated S variable. The RMS S OMF from

ER-16 are 24% lower than those from CE-16 and those of ERKF-16 are 24% lower than those from EnKF-16. In contrast with ER-16 and ERKF-16, while EnKF-16 performs best for T, it does not significantly improve over CE-16 for S.

	CE-16	ER-16	EnKF-16	ERKF-16
RMS T OMF (°C)	1.62	1.51	0.91	0.92
RMS S OMF (PSU)	0.51	0.39	0.49	0.37

**Table 1.** April-June 2006 global mean RMS OMF differences with the assimilated Argo T data and unassimilated Argo S data in each of the CE-16, ER-16, EnKF-16 and ERKF-16 runs.

Figures 2 and 3 expand upon the information provided by Table 1. Figure 2 shows the vertical average of the RMS OMF for the assimilated Argo T data averaged over the 3-month data assimilation period and binned to  $3^\circ \times 3^\circ$  horizontal boxes. Figure 3 shows the corresponding binned RMS OMF horizontal distributions for the unassimilated Argo S data. For the T data, Figures 2c (EnKF-16) and 2d (ERKF-16) show very similar distributions, while the distribution of RMS T OMF in ER-16 (Figure 2b) resembles closely that of CE-16 (Figure 2a). For the unassimilated S data, there are strong similarities between the distributions of RMS OMF of CE-16 (Figure 3a) and EnKF-16 (Figure 3c) and of ER-16 (Figure 3b) and ERKF-16 (Figure 3d).

To better understand the effect of the recentering step on the unassimilated S variable, Figures 4 and 5 show the differences of the RMS S OMF of CE-16 from those of EnKF-16 and ER-16. Figure 4 corresponds to the upper 200 meters. Figure 5 is for the 200-2000 meter depth range. Warm (cold) colors indicate that the analysis is closer to (further away from) the unassimilated Argo S data than the control. Clearly, the better performance in terms of RMS S OMF of the recentering step from ER-16 over the EnKF analysis from EnKF-16 is mostly due the recentering's effectiveness in the upper 200 meters, where the RMS S OMF of EnKF-16 are not significantly different from those of CE-16 (Figure 4). Both the recentering and the EnKF analysis are effective below 200 meters, but while ER-16 is generally closer than both EnKF-16 and CE-16 to Argo S in the tropics where there are more Argo profiles than in the extratropics, it is less effective south of  $45^\circ\text{S}$  where the observations are sparse. However, when the recentering is applied first and then followed by the EnKF analysis, as is the case in ERKF-16, the poor performance of the recentering in the data-sparse high southern latitudes is compensated by the EnKF analysis ability to use its background error covariance estimates to propagate information from data rich regions to data poor regions.

## 6. Conclusions

The ensemble recentering Kalman filter (ERKF) first applies a recentering step to adjust the distribution of an ensemble of model states so that the ensemble mean is shifted to the prior position of the most likely sample without altering the ensemble statistics, after which an EnKF analysis step is applied. In our experiments, the ERKF is better able than the EnKF to improve model estimates of the unassimilated S variable when Argo T profiles are assimilated into the GEOS-5 CGCM. The effectiveness of the recentering step is attributed to its ability to gradually correct model biases over time and to the fact that it replaces the ensemble mean with a balanced state. Although the recentering is less effective in poorly observed regions, the EnKF analysis compensates for this fact by propagating information from data rich areas to data poor areas using its background error covariance estimates. In contrast to particle filters, which must



integrate a very large number of model trajectories to be competitive with the EnKF, the ERKF can improve upon the performance of the EnKF without requiring a larger ensemble.

## Acknowledgement

This work is supported by NASA's Modeling Analysis and Prediction Program under WBS 802678.02.17.01.25. The infrastructure for the runs is provided by the NASA Center for Climate Simulation (NCCS). Yuri Vikhliav, Max Suarez and Bin Zhao helped configure the GEOS-5 modeling system. Guillaume Vernieres helped configure the data assimilation system and Robin Kovach assisted with the observation preprocessing and with plotting the results. Michele Rienecker provided critical insights that helped refocus the initial draft manuscript.

## References

- Bloom, S., L. Takacs, A. DaSilva, and D. Ledvina (1996), Data assimilation using incremental analysis updates, *Mon. Wea. Rev.*, **124**, 1256-1271.
- Evensen, G. (1994), Sequential data assimilation with a nonlinear quasi-geostrophic model using Monte Carlo methods to forecast error statistics, *J. Geophys. Res.*, **C99**, 10,143-10,162.
- Gordon, N. J.; D.J. Salmond, and A.F.M. Smith (1993), Novel approach to nonlinear/non-Gaussian Bayesian state estimation, *IEEE Proceedings on Radar and Signal Processing*, **140** (2), 107-113.
- Hill, C., C. Deluca, V. Balaji, M. Suarez and A. Da Silva (2004), The architecture of the earth system modeling framework, *Computing in Science and Engineering*, **6** (1), 18-28.
- Hunke, E.C., W.H. Lipscomb, A.K. Turner, N. Jeffery, S. Elliott (2013), CICE: the Los Alamos Sea Ice Model Documentation and Software User's Manual Version 5.0, LA-CC-06-012, Los Alamos National Laboratory, Los Alamos NM.
- Kalman, R. (1960), A new approach to linear filtering and prediction problems, *J. Basic Eng.*, **D82**, 35-45.
- Keppenne, C.L., and M.M. Rienecker (2003), Assimilation of temperature into an isopycnal ocean general circulation model using a parallel ensemble Kalman filter, *J. Mar. Sys.*, **40-41**, 363-380.
- Keppenne, C.L., M.M. Rienecker, N.P. Kurkowski, and D.A. Adamec (2005), Ensemble Kalman filter assimilation of temperature and altimeter data with bias correction and application to seasonal prediction, *Nonlinear Processes in Geophysics*, **12**, 491-503.
- Keppenne, C.L., M.M. Rienecker, J.P. Jacob, and R.M. Kovach (2008), Error covariance modeling in the GMAO ocean ensemble Kalman filter, *Mon. Wea. Rev.*, **136**, 2964-2982.
- Keppenne, C.L., M.M. Rienecker, R.M. Kovach and G. Vernieres (2013), Background error covariance estimation using information from a single model trajectory with application to ocean data assimilation, *Ocean Modelling*, submitted
- Oke, P.R., G.B. Brassington, D.A. Griffin, and A. Schiller (2010), Ocean data assimilation: a case for ensemble optimal interpolation, *Australian Meteorological and Oceanographic Journal*, **59**, 67-76.
- Rienecker, M.M., M.J. Suarez, R. Todling, J. Bacmeister, L. Takacs, H.-C. Liu, W. Gu, M. Sienkiewicz, R.D. Koster, R. Gelaro, I. Stajner, and J.E. Nielsen (2008), The GEOS-5 data assimilation system. Documentation Versions 5.0.1, 5.1.0, and 5.2.0, *NASA Tech. Rep. Series on Global Modeling and Data Assimilation* **27**, NASA/TM-2008-104606, 1-118.



311 Vernieres, G., C.L. Keppenne, M.M. Rienecker, J.P. Jacob and R.N. Kovach (2012), The  
312 GEOS-ODAS description and evaluation. *NASA Tech. Rep. Series on Global Modeling and*  
313 *Data Assimilation* **30**, NASA/TM-2012-104606.  
314 Weerts, A.H., and G.Y.H. El Serafy (2006), Particle filtering and ensemble Kalman filtering for  
315 state updating with hydrological conceptual rainfall-runoff models, *Water Resources*  
316 *Research*, **42** (9), W09043, 17pp.  
317

318 **Figure Captions**

319

320 **Figure 1.** Schematic representation of the recentering step of the ERKF showing the sample  
321 mean,  $\bar{x}$  (doubly outlined small open circle), the most likely sample  $x_p$  (small open circle), the  
322 control observations,  $y$  (large open circle) and the recentering increment,  $\Delta_p$  (arrow connecting  
323  $\bar{x}$  and  $x_p$ ). The filled circles represent particles other than  $x_p$ . The cluster of circles shown to  
324 the left represents the prior sample forecast. The cluster to the right represents the recentered  
325 sample.

326

327 **Figure 2.** RMS OMF differences with respect to the assimilated Argo T data averaged vertically  
328 and binned to  $3^\circ \times 3^\circ$  boxes over April-June 2006 in (a) CE-16, (b) ER-16 (c) EnKF-16 and (d)  
329 ERKF-16.

330

331 **Figure 3.** Same as Figure 2 for the unassimilated Argo S data.

332

333 **Figure 4.** Difference of CE-16 RMS OMF for Argo S from those of (a) EnKF-16 and (b)  
334 ERKF-16 averaged over the upper 200 meters for the unassimilated Argo S data. Warm (cold)  
335 colors indicate areas where the analysis is closer to (further from) the observations than CE-16.

336

337 **Figure 5.** Same as Figure 4 for the 200-2000 meter depth range.

338