

The Ensemble Kalman Filter: Theoretical formulation and practical implementation

Geir Evensen

Norsk Hydro Research Centre, Bergen, Norway

Based on Evensen 2003, Ocean Dynamics, Vol 53, No 4

The Ensemble Kalman Filter (EnKF)

- Represents error statistics using an ensemble of model states.
- Evolves error statistics by ensemble integrations.
- “Variance minimizing” analysis scheme operating on the ensemble.



- Monte Carlo, low rank, error subspace method.
- Converges to the Kalman Filter with increasing ensemble size.
- Fully nonlinear error evolution, contrary to EKF.
- Assumption of Gaussian statistics in analysis scheme.

The error covariance matrix

Define ensemble covariances around the ensemble mean

$$\mathbf{P}^f \simeq \mathbf{P}_e^f = \overline{(\boldsymbol{\psi}^f - \overline{\boldsymbol{\psi}}^f)(\boldsymbol{\psi}^f - \overline{\boldsymbol{\psi}}^f)^T}$$

$$\mathbf{P}^a \simeq \mathbf{P}_e^a = \overline{(\boldsymbol{\psi}^a - \overline{\boldsymbol{\psi}}^a)(\boldsymbol{\psi}^a - \overline{\boldsymbol{\psi}}^a)^T}$$

- The ensemble mean $\overline{\boldsymbol{\psi}}$ is the best-guess.
- The ensemble spread defines the error variance.
- The covariance is determined by the smoothness of the ensemble members.
- A covariance matrix can be represented by an ensemble of model states (not unique).

Dynamical evolution of error statistics

- Each ensemble member evolve according to the model dynamics which is expressed by a stochastic differential equation

$$d\psi = \mathbf{f}(\psi)dt + \mathbf{g}(\psi)dq.$$

- The probability density then evolve according to Kolmogorov's equation

$$\frac{\partial \phi}{\partial t} + \sum_i \frac{\partial(f_i \phi)}{\partial \psi_i} = \frac{1}{2} \sum_{i,j} \frac{\partial^2 \phi (\mathbf{g} \mathbf{Q} \mathbf{g}^T)_{ij}}{\partial \psi_i \partial \psi_j}.$$

- This is the fundamental equation for evolution of error statistics and can be solved using Monte Carlo methods.

Analysis scheme (1)

- Given an ensemble of model forecasts, ψ_j^f , defining forecast error covariance

$$P^f \simeq \mathbf{P}_e^f = \overline{(\psi^f - \bar{\psi}^f)(\psi^f - \bar{\psi}^f)^T}.$$

- Create an ensemble of observations

$$\mathbf{d}_j = \mathbf{d} + \boldsymbol{\epsilon}_j,$$

with

- \mathbf{d} , the real observations,
- $\boldsymbol{\epsilon}_j$, a vector of observation noise,
- $\overline{\boldsymbol{\epsilon}\boldsymbol{\epsilon}^T} = \mathbf{R}_e \simeq \mathbf{R}$.

Analysis scheme (2)

- Update each ensemble member according to

$$\psi_j^a = \psi_j^f + K_e(d_j - H\psi_j^f),$$

where

$$K_e = P_e^f H^T (H P_e^f H^T + R_e)^{-1}.$$

- Thus the update of the mean becomes

$$\bar{\psi}^a = \bar{\psi}^f + K_e(d - H\bar{\psi}^f).$$

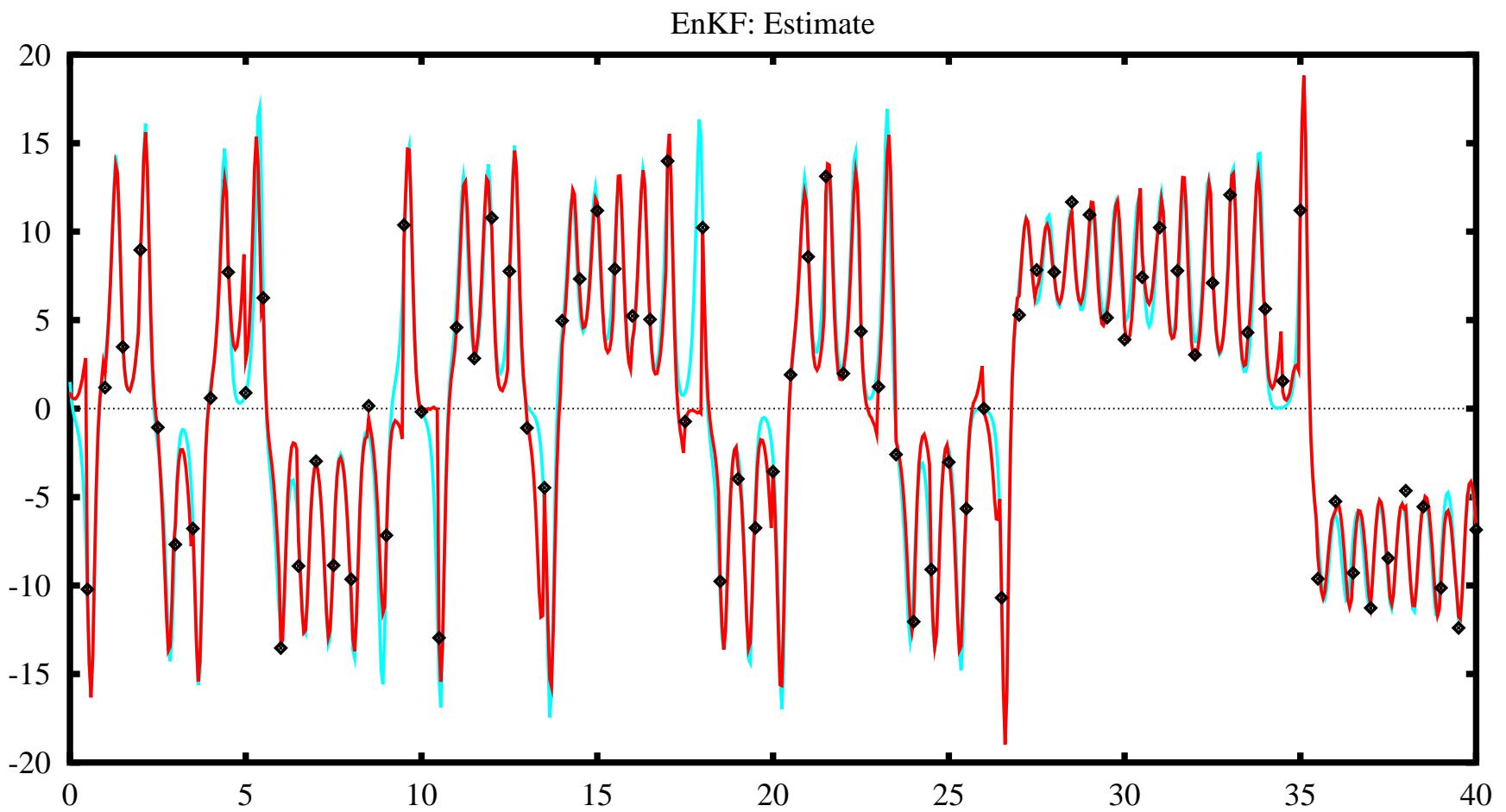
- The posterior error covariance becomes

$$P_e^a = (I - K_e H) P_e^f.$$

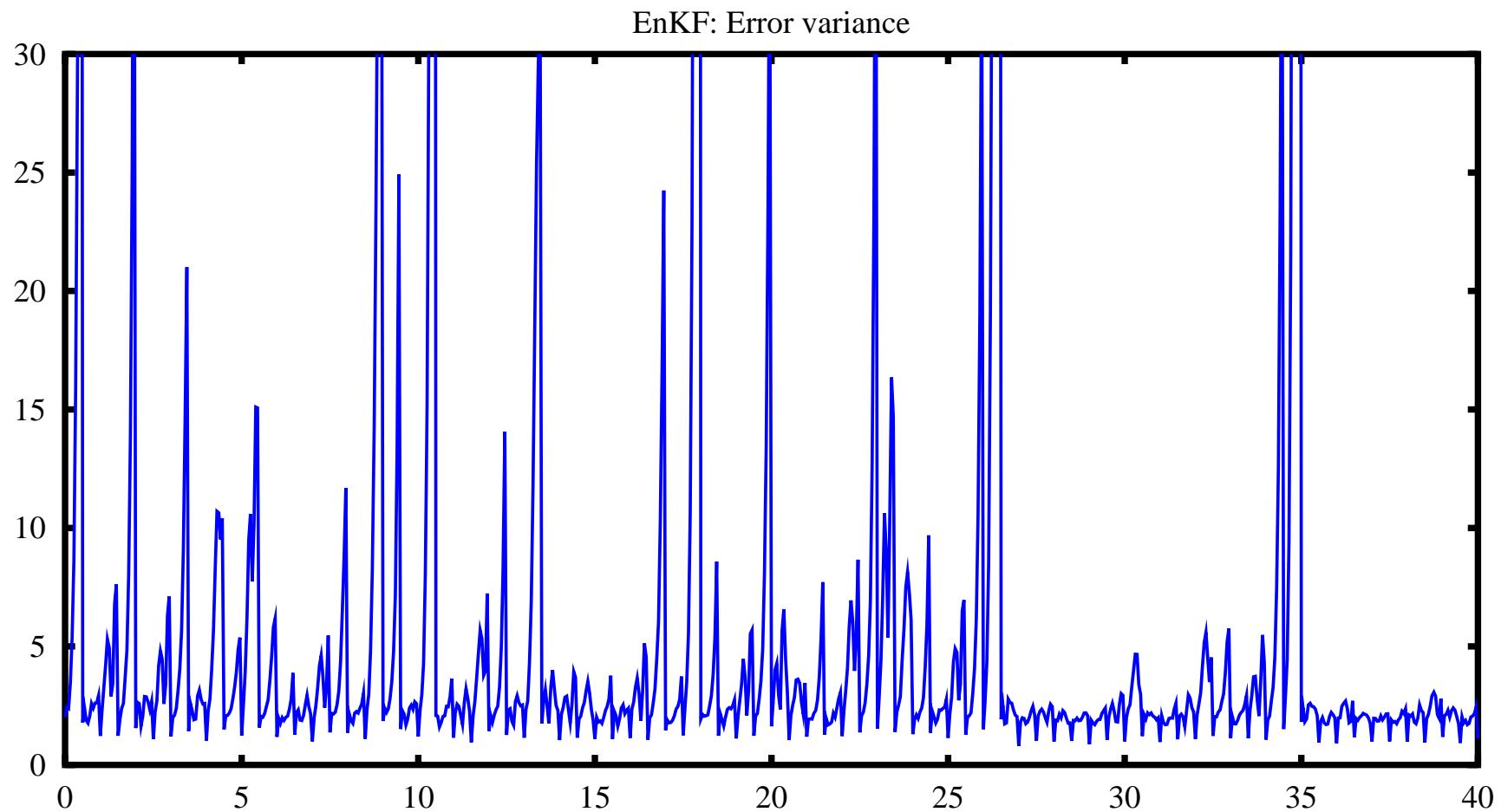
Example: Lorenz model

- Application with the chaotic Lorenz model
- Illustrates properties with highly nonlinear dynamical models.
- From Evensen (1997), MWR.

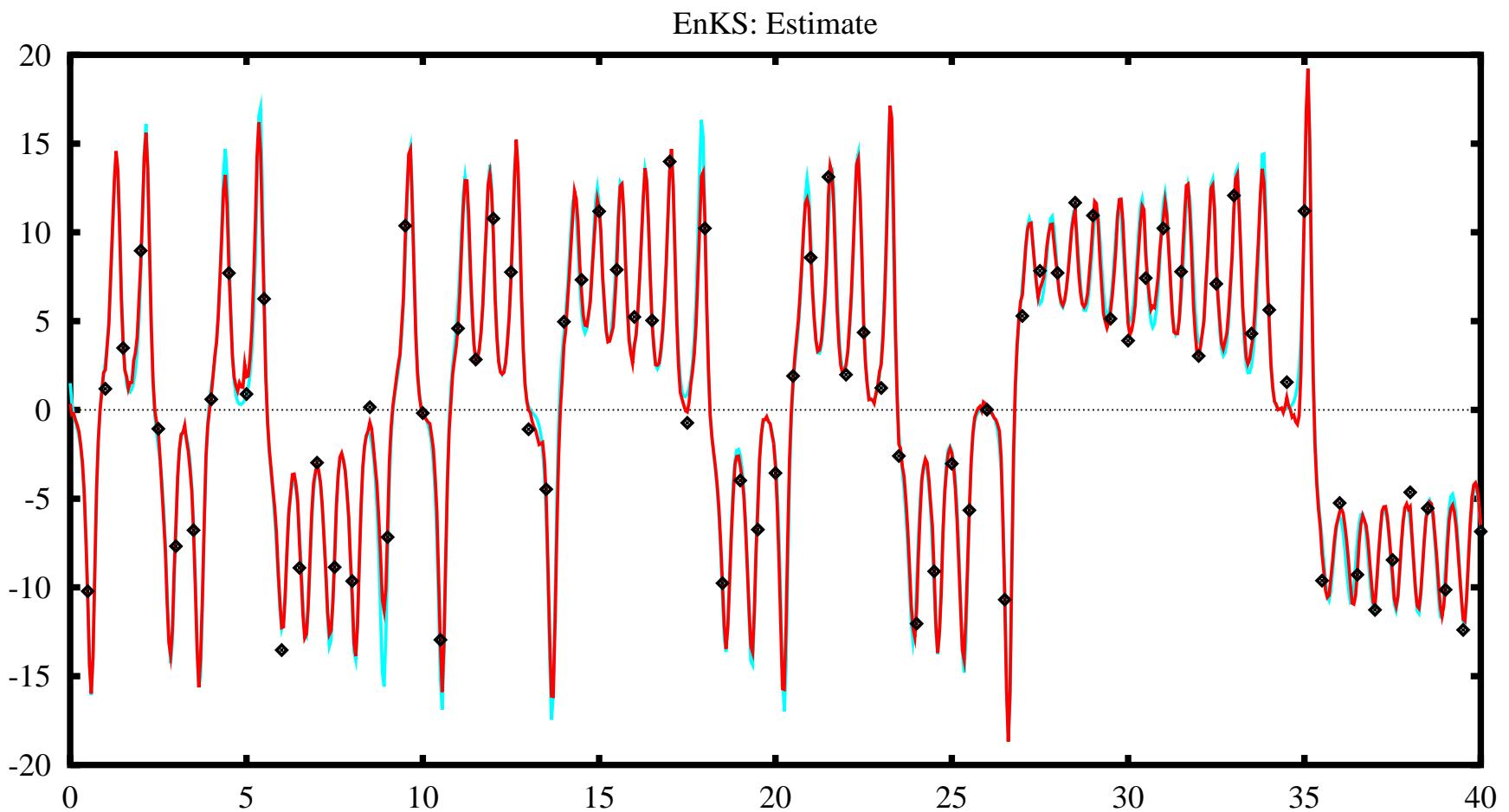
EnKF solution



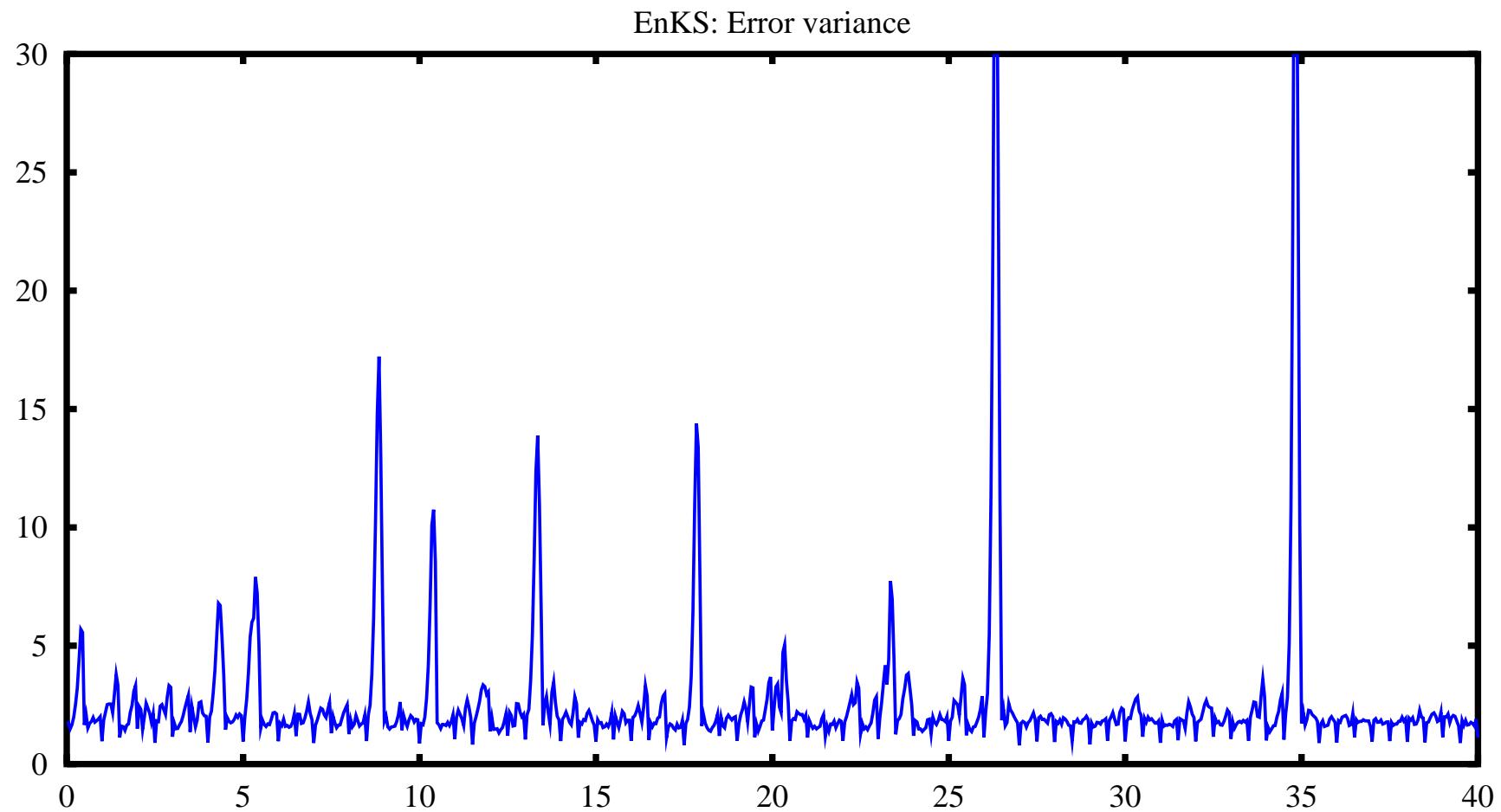
EnKF error variance



EnKS solution



EnKS error variance



Summary: Lorenz model

- The EnKF and EnKS works well with highly nonlinear dynamical models.
- There is no linearization in the evolution of error statistics.
- Methods using tangent linear or adjoint operators have problems with the Lorenz equations:
 - limited by the predictability time,
 - limited by the validity time of tangent linear operator.
- Can we expect the same to be true for high resolution ocean and atmosphere models?

Analysis equation (1)

- Define the ensemble matrix

$$\mathbf{A} = (\psi_1, \psi_2, \dots, \psi_N) \in \mathbb{R}^{n \times N}.$$

- The ensemble mean is (defining $\mathbf{1}_N \in \mathbb{R}^{N \times 1} \equiv \mathbf{1}/N$)

$$\overline{\mathbf{A}} = \mathbf{A}\mathbf{1}_N.$$

- The ensemble perturbations become

$$\mathbf{A}' = \mathbf{A} - \overline{\mathbf{A}} = \mathbf{A}(\mathbf{I} - \mathbf{1}_N).$$

- The ensemble covariance matrix $\mathbf{P}_e \in \mathbb{R}^{n \times n}$ becomes

$$\mathbf{P}_e = \frac{\mathbf{A}'(\mathbf{A}')^T}{N - 1}.$$

Analysis equation (2)

- Given a vector of measurements $\mathbf{d} \in \mathbb{R}^m$, define

$$\mathbf{d}_j = \mathbf{d} + \boldsymbol{\epsilon}_j, \quad j = 1, \dots, N,$$

stored in

$$\mathbf{D} = (\mathbf{d}_1, \mathbf{d}_2, \dots, \mathbf{d}_N) \in \mathbb{R}^{m \times N}.$$

- The ensemble perturbations are stored in

$$\mathbf{E} = (\boldsymbol{\epsilon}_1, \boldsymbol{\epsilon}_2, \dots, \boldsymbol{\epsilon}_N) \in \mathbb{R}^{m \times N},$$

thus, the measurement error covariance matrix becomes

$$\mathbf{R}_e = \frac{\mathbf{E} \mathbf{E}^T}{N - 1}.$$

Analysis equation (3)

- The analysis equation can now be written

$$\mathbf{A}^a = \mathbf{A} + \mathbf{P}_e \mathbf{H}^T (\mathbf{H} \mathbf{P}_e \mathbf{H}^T + \mathbf{R}_e)^{-1} (\mathbf{D} - \mathbf{H} \mathbf{A}).$$

- Defining the innovations $\mathbf{D}' = \mathbf{D} - \mathbf{H} \mathbf{A}$ and using previous definitions:

$$\mathbf{A}^a = \mathbf{A} + \mathbf{A}' (\mathbf{H} \mathbf{A}')^T \left((\mathbf{H} \mathbf{A}') (\mathbf{H} \mathbf{A}')^T + \mathbf{E} \mathbf{E}^T \right)^{-1} \mathbf{D}'.$$

i.e., analysis expressed entirely in terms of the ensemble

Analysis equation (4)

- Define $S = HA'$ and $C = SS^T + EE^T$.
- Use $A' = A(I - 1_N)$.
- Use $1_N S^T \equiv 0$.

$$\begin{aligned} A^a &= A + A'S^T \left(SS^T + EE^T \right)^{-1} D' \\ &= A + A(I - 1_N)S^T C^{-1} D' \\ &= A \left(I + (I - 1_N)S^T C^{-1} D' \right) \\ &= A \left(I + S^T C^{-1} D' \right) \\ &= AX \end{aligned} \tag{1}$$

Remarks on the analysis equation (1)

- Covariances only needed between observed variables at measurement locations ($HP_e = SA'^T$).
- P_e never computed but indirectly used to determine $HP_eH^T = SS^T$.
- Analysis may be interpreted as:
 - combination of forecast ensemble members, or,
 - forecast plus combination of covariance functions.
- Accuracy of analysis is determined by:
 - the accuracy of X ,
 - the properties of the ensemble error space.

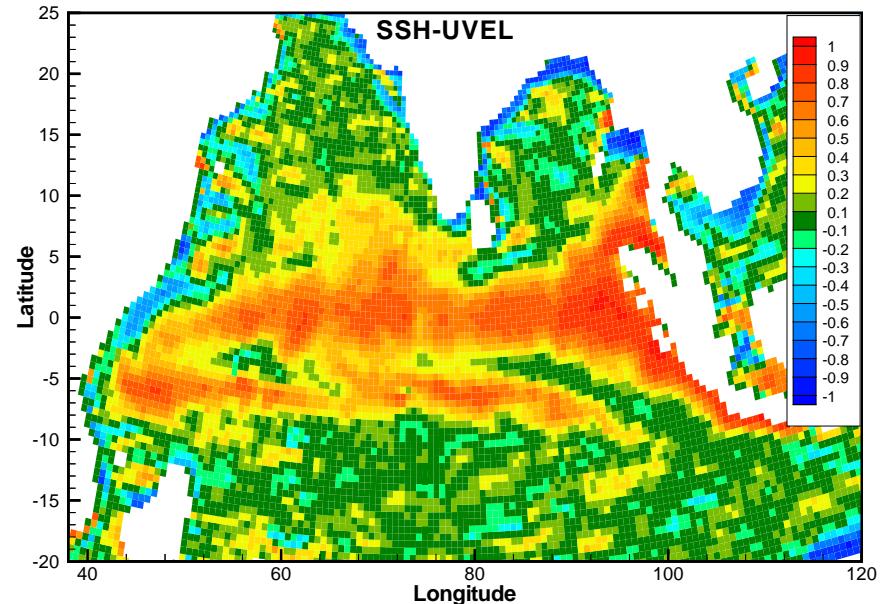
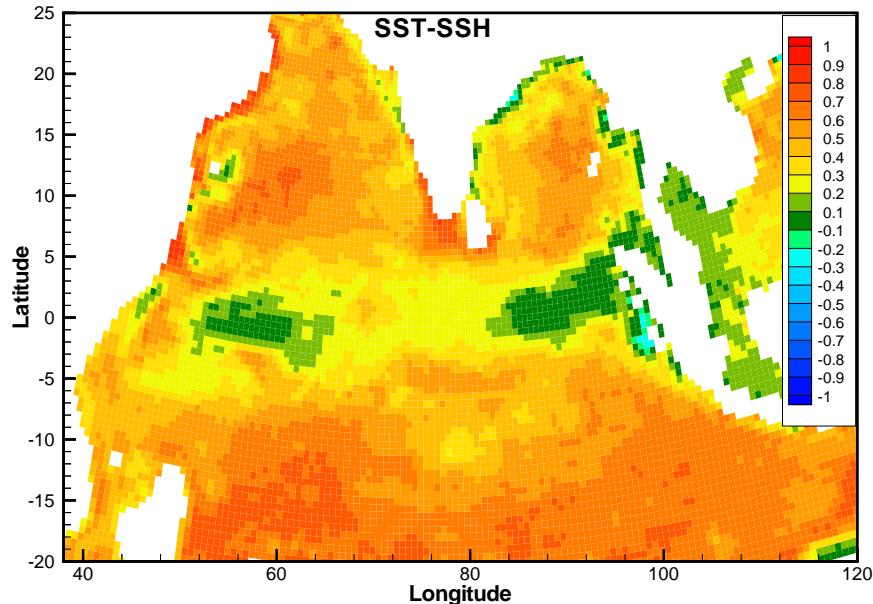
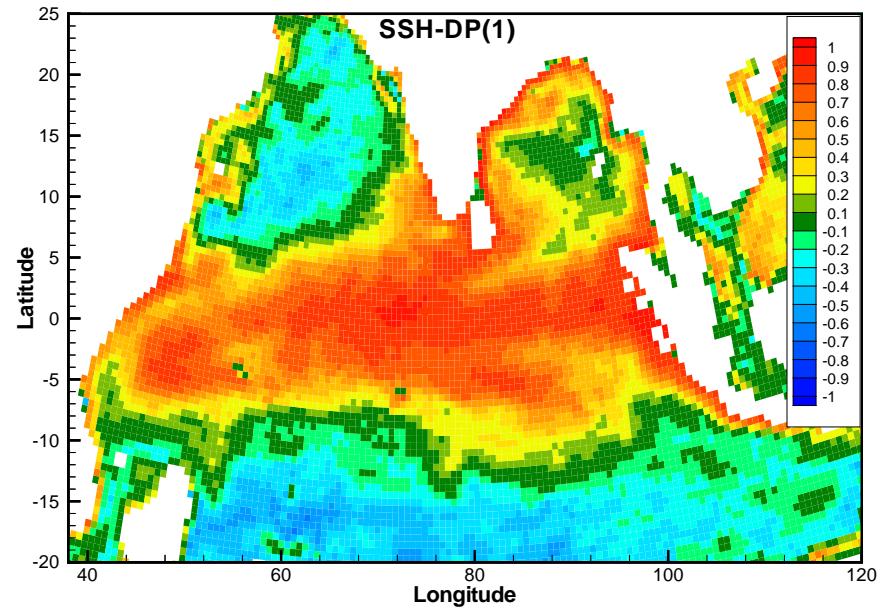
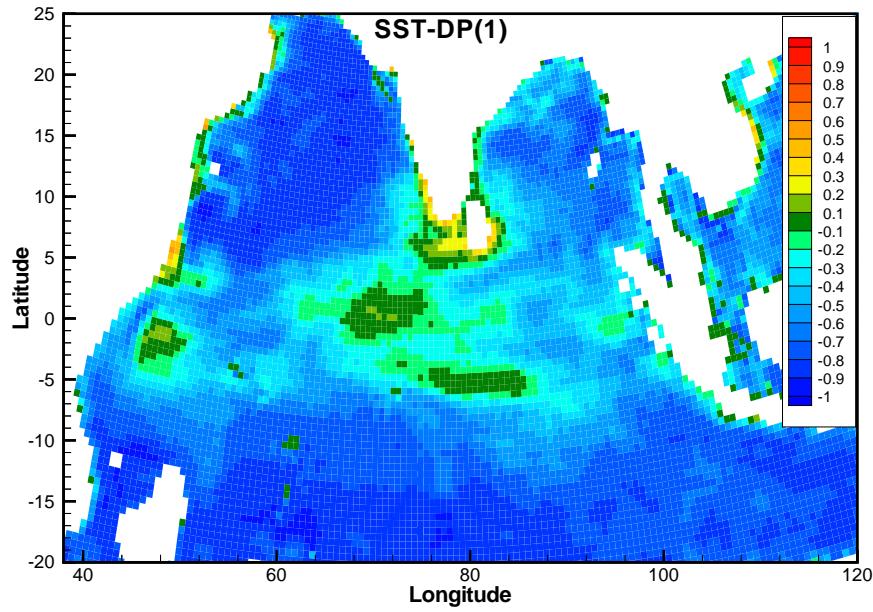
Remarks on the analysis equation (2)

- For a linear model, any choice of X will result in an analysis which is also a solution of the model.
- Filtering of covariance functions introduces nondynamical modes in the analysis.

Examples of ensemble statistics

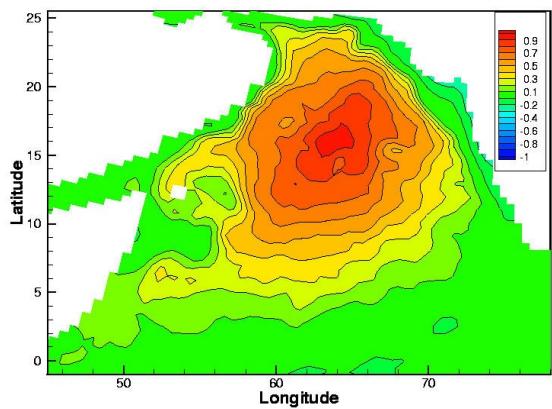
- Taken from Haugen and Evensen (2002), Ocean Dynamics.
- OGCM (Micom) for the Indian Ocean.
- Assimilation of SST and SLA data.

Spatial correlations

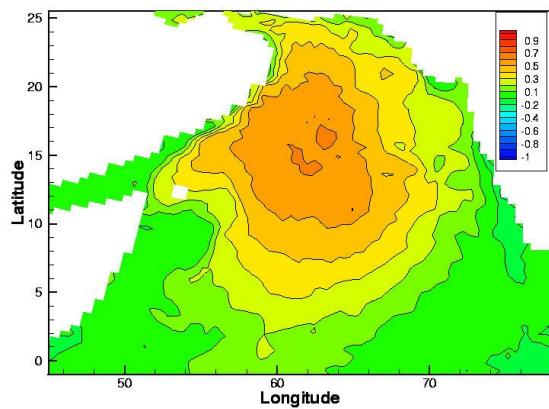


Correlation functions

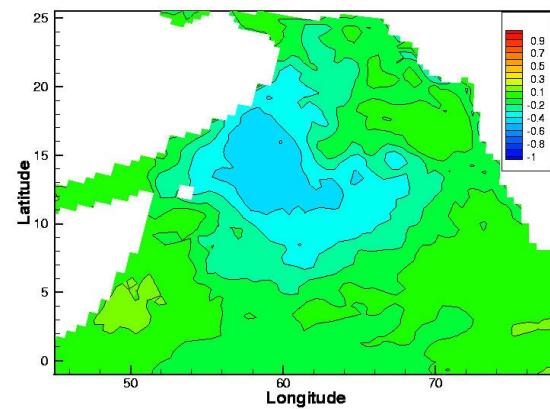
SSH–SSH



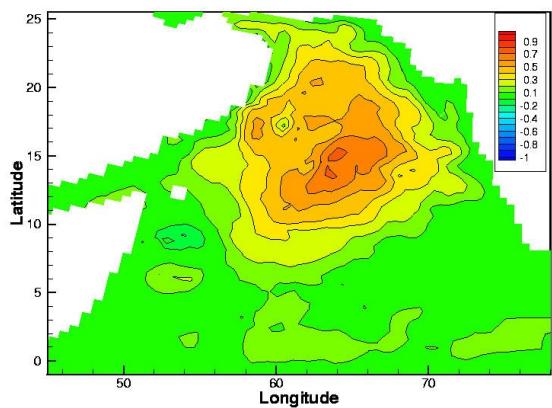
SSH–SST



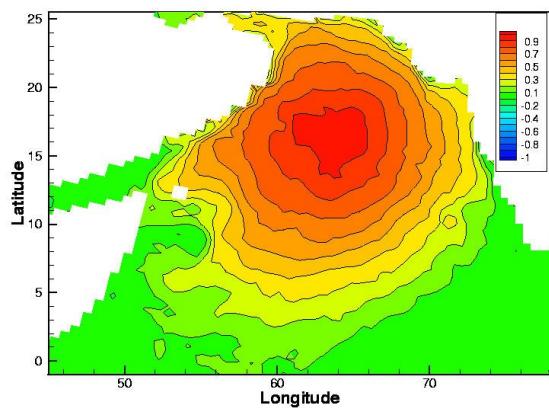
SSH–DP



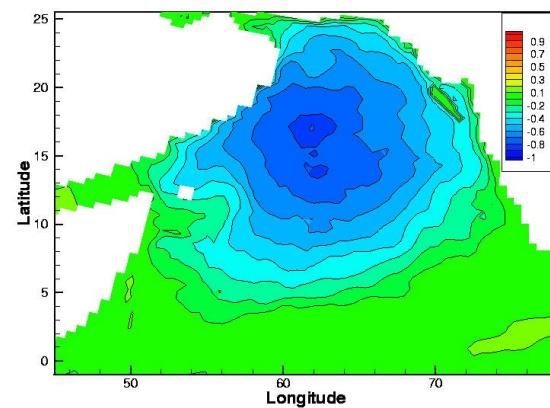
SST–SSH



SST–SST

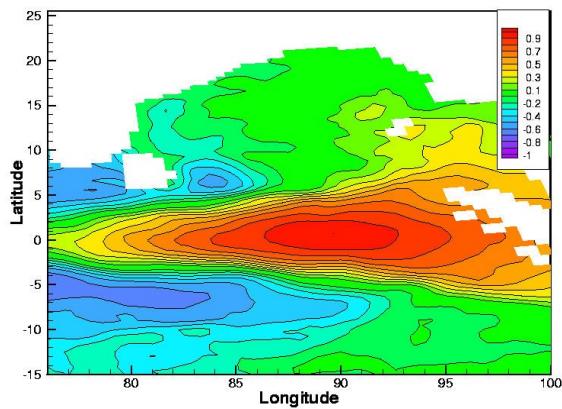


SST–DP

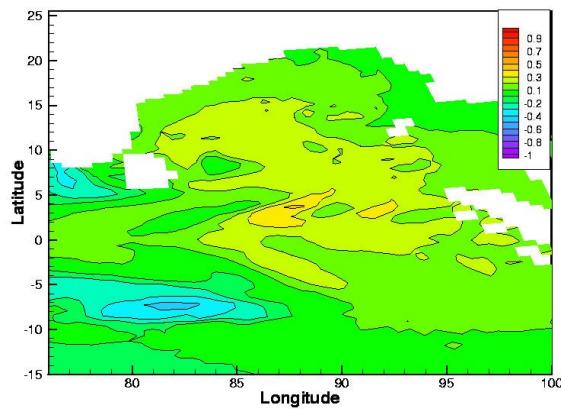


Correlation functions

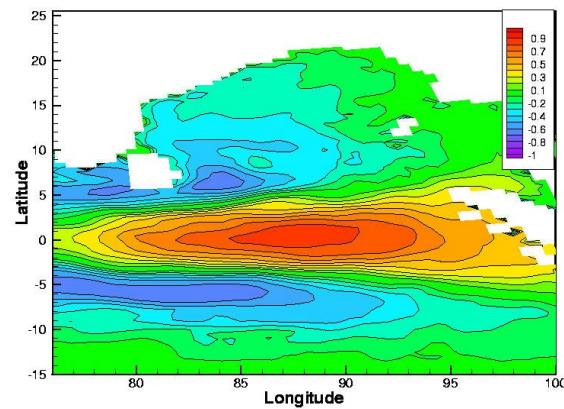
SSH–SSH



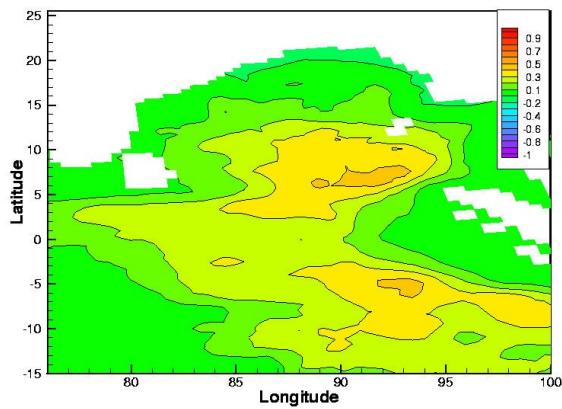
SSH–SST



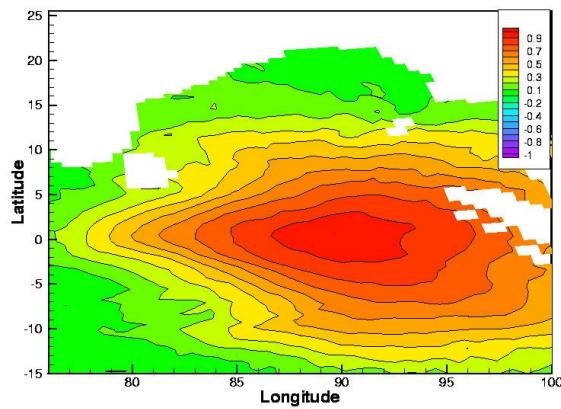
SSH–DP



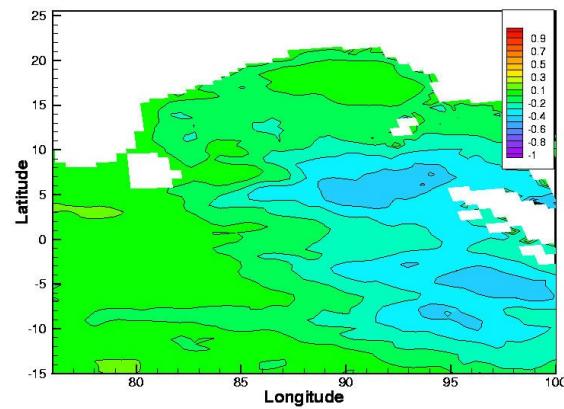
SST–SSH



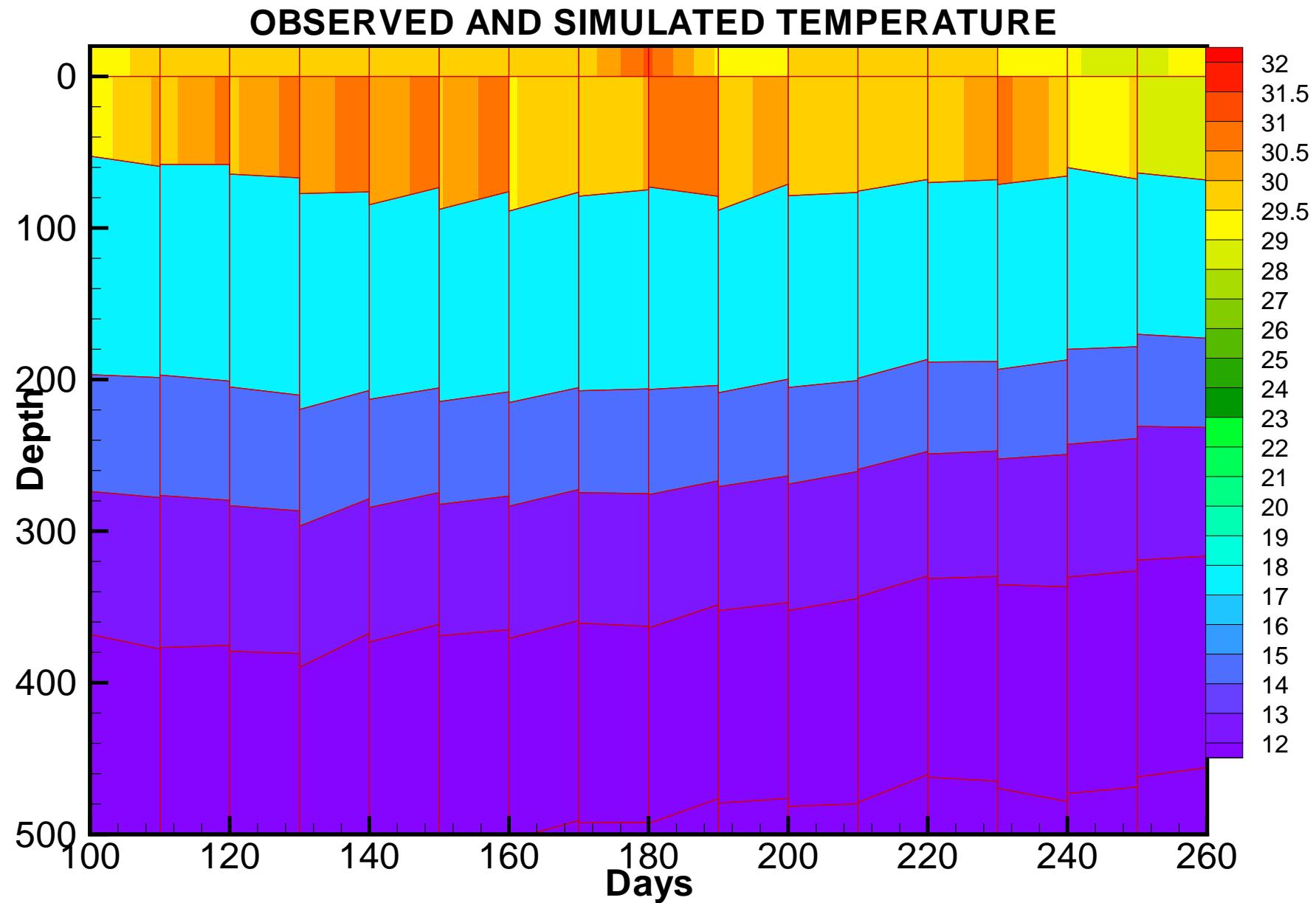
SST–SST



SST–DP



Time–Depth: Temperature



Ensemble Kalman Smoother (EnKS)

- Derived in Evensen and van Leeuwen (2000), MWR.
- Starts with EnKF solution.
- Computes updates backward in time;
 - sequentially for each measurement time,
 - using covariances in time,
 - no backward integrations.
- The analysis becomes for $t_{i-1} \leq t' < t_i \leq t_k$:

$$A_{\text{EnKS}}^a(t') = A_{\text{EnKF}}(t') \prod_{j=i}^k X(t_j)$$

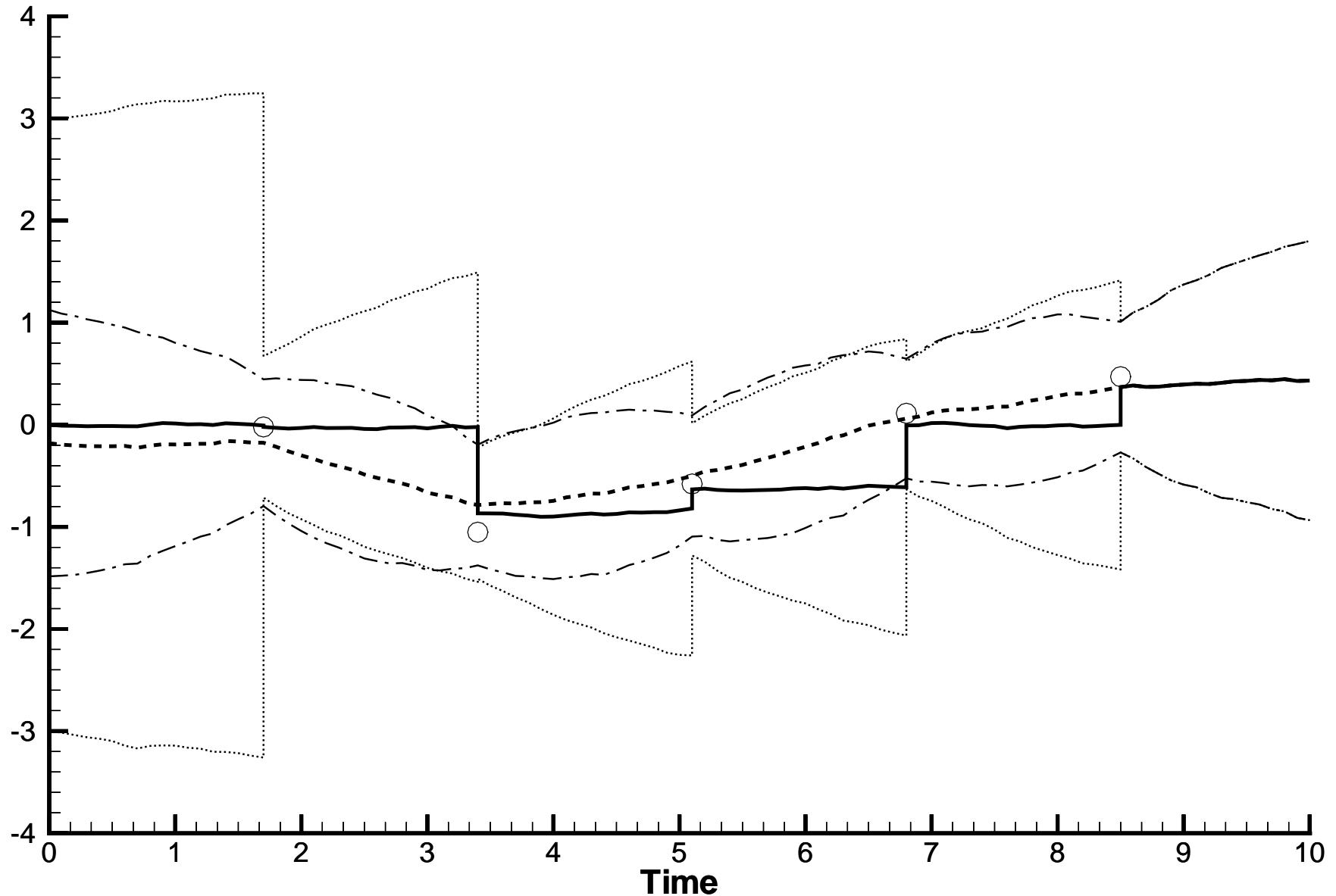
Time correlated model noise

- Most schemes assume white model noise in time.
- Augment model state to include time correlations

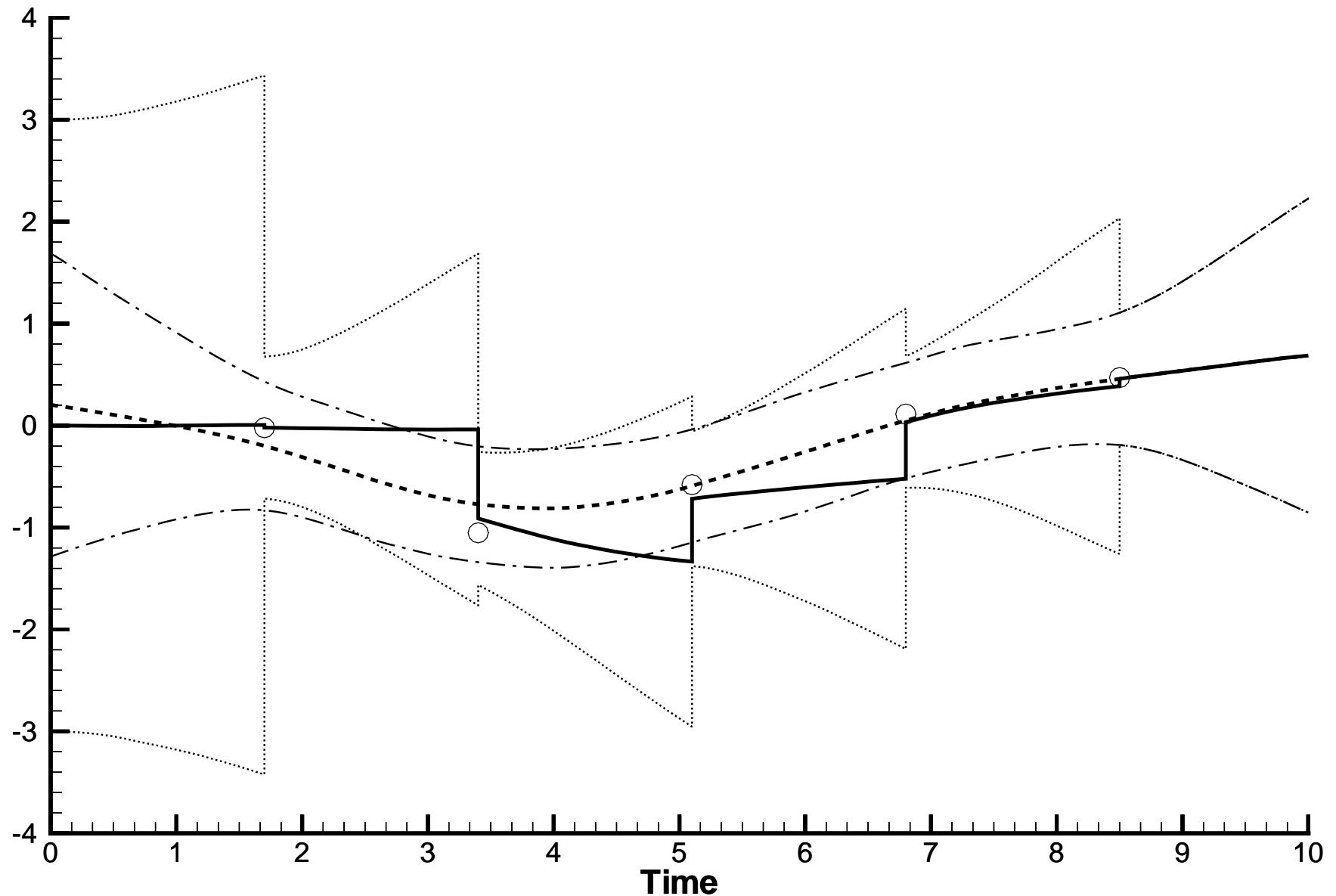
$$\begin{pmatrix} q_k \\ \psi_k \end{pmatrix} = \begin{pmatrix} \alpha q_{k-1} \\ \psi_{k-1} + \sqrt{\Delta t} \sigma \rho q_k \end{pmatrix} + \begin{pmatrix} \sqrt{1 - \alpha^2} w_{k-1} \\ 0 \end{pmatrix}.$$

- White noise when $\alpha = 0$ and $\rho = 1$.

Results ($\alpha = 0$)



Results ($\alpha = 0.95$)



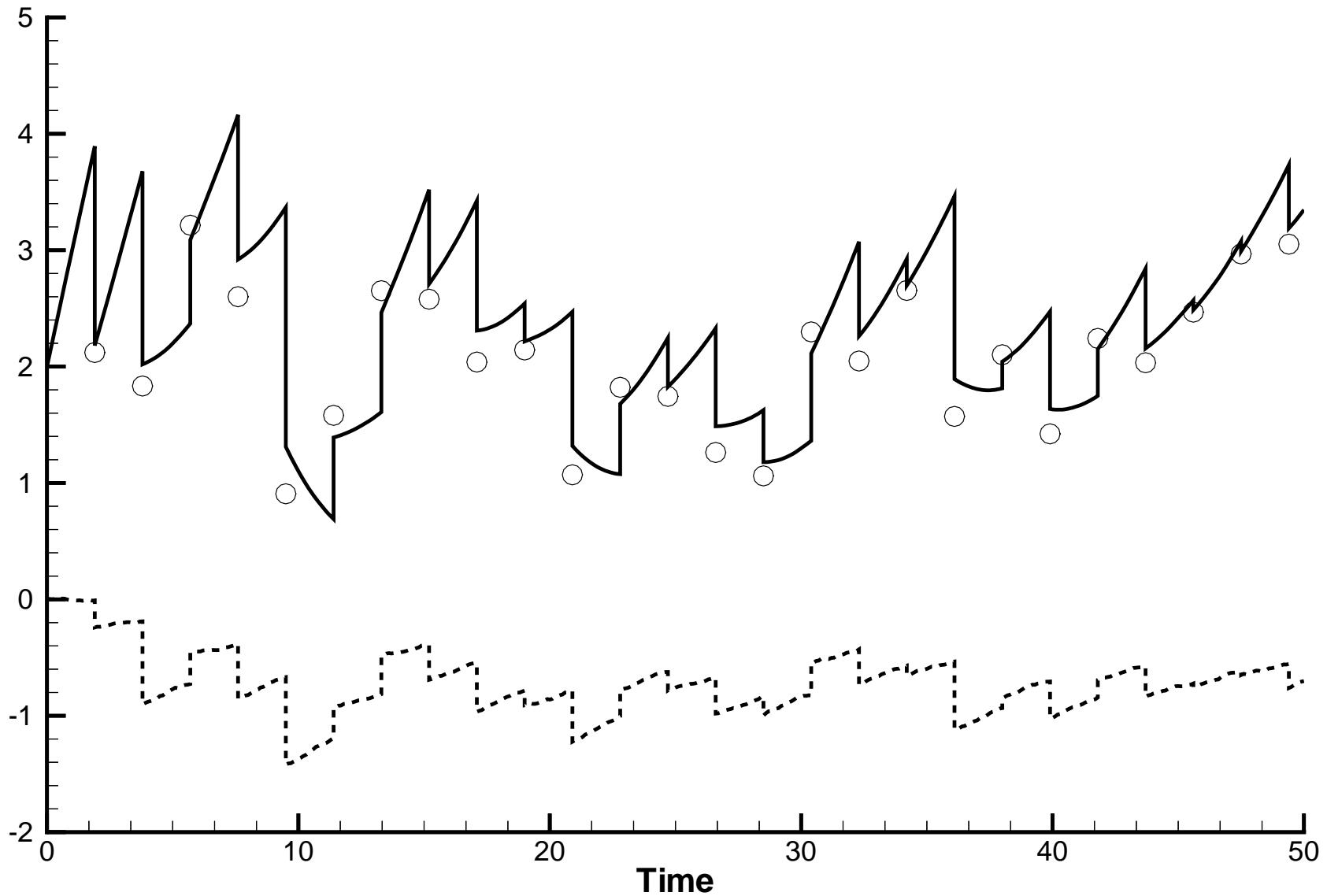
Parameter and bias estimation

- Introduces poorly known parameter β_k in model

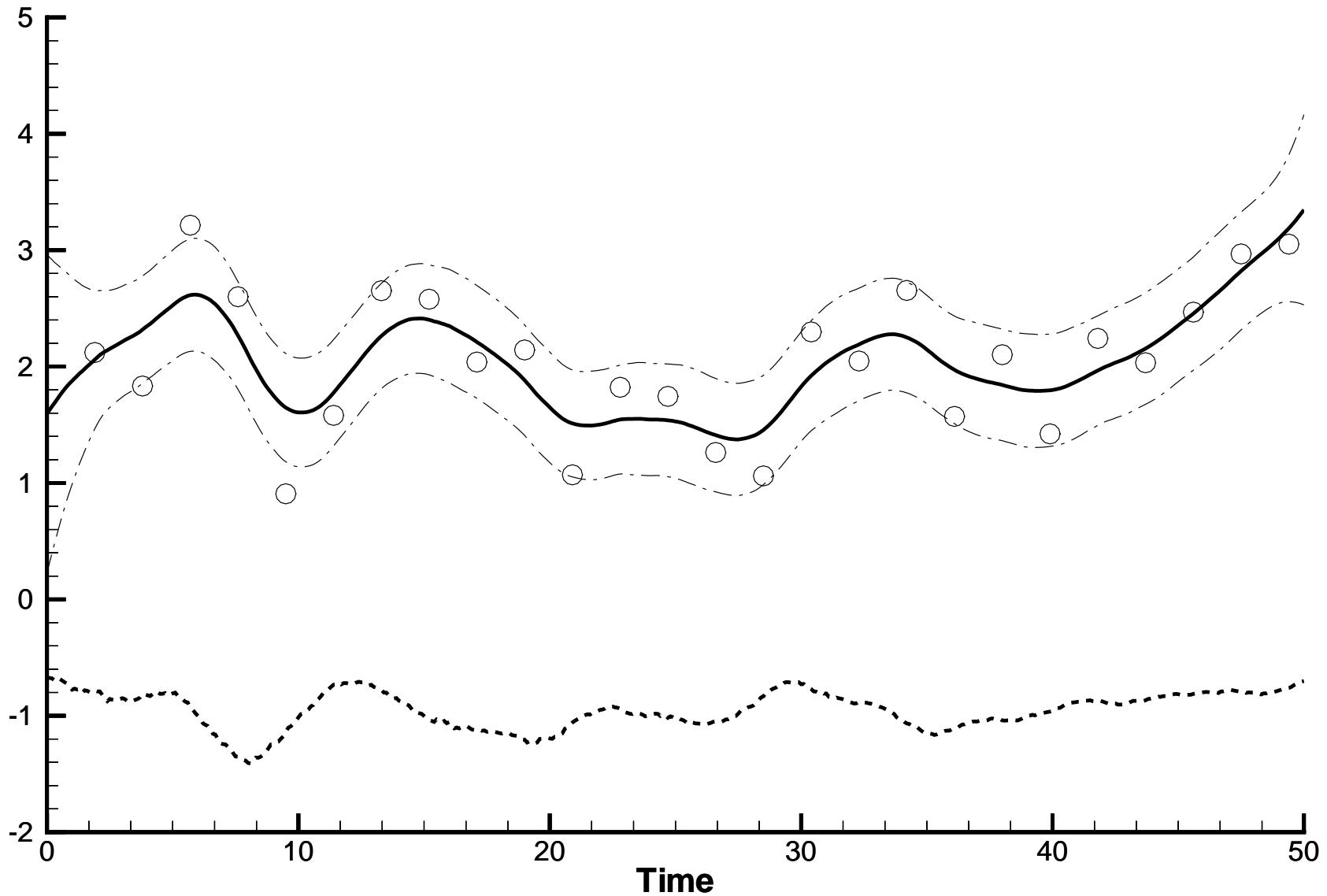
$$\begin{pmatrix} q_k \\ \beta_k \\ \psi_k \end{pmatrix} = \begin{pmatrix} \alpha q_{k-1} \\ \beta_{k-1} \\ \psi_{k-1} + (1 + \beta_k) \Delta t + \sqrt{\Delta t} \sigma \rho q_k \end{pmatrix} + \begin{pmatrix} \sqrt{1 - \alpha^2} w_{k-1} \\ 0 \\ 0 \end{pmatrix}.$$

- Two cases
 - $\beta \equiv 0$.
 - β a poorly known parameter with $\beta_0 = 0$.

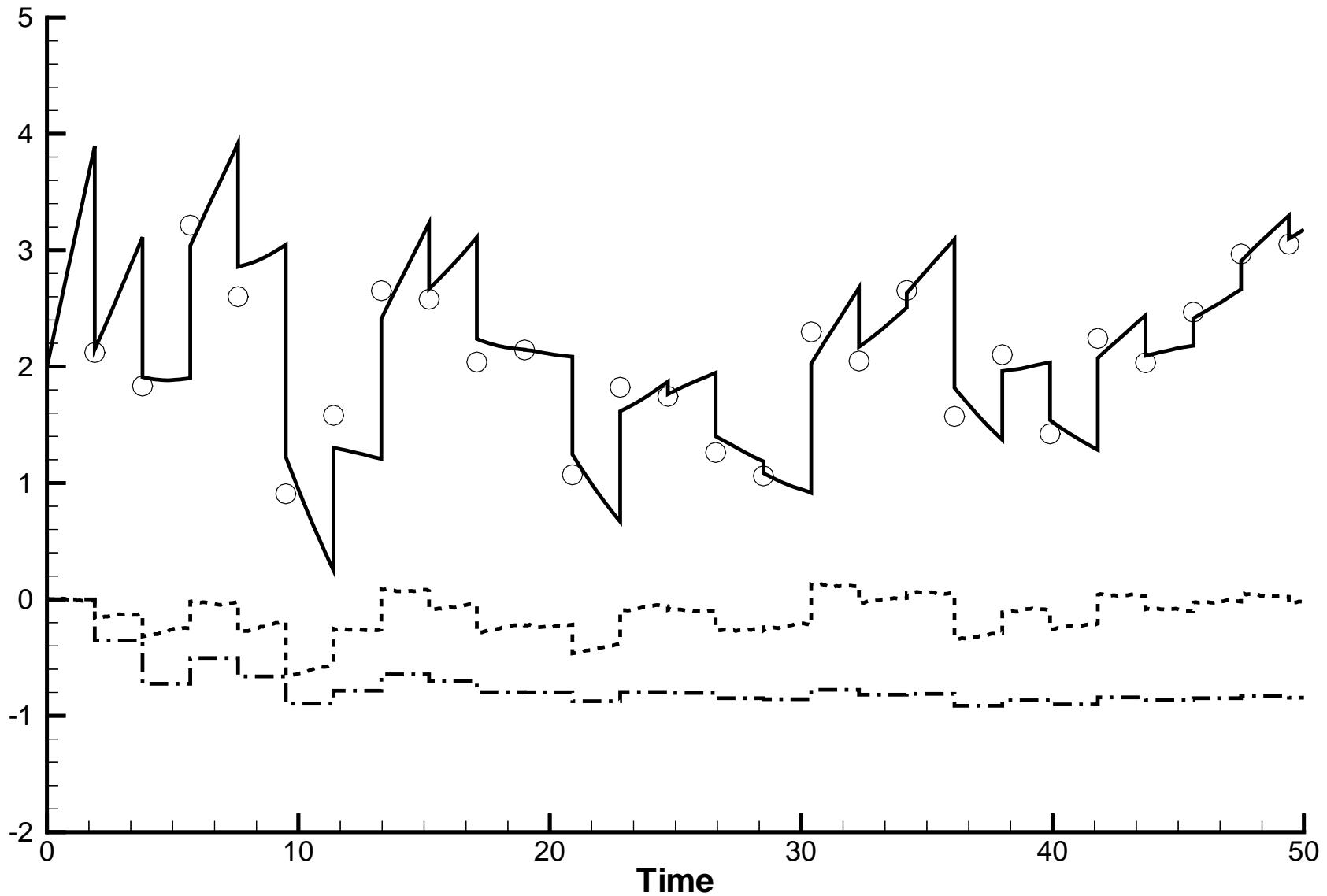
Estimate and model error, EnKF



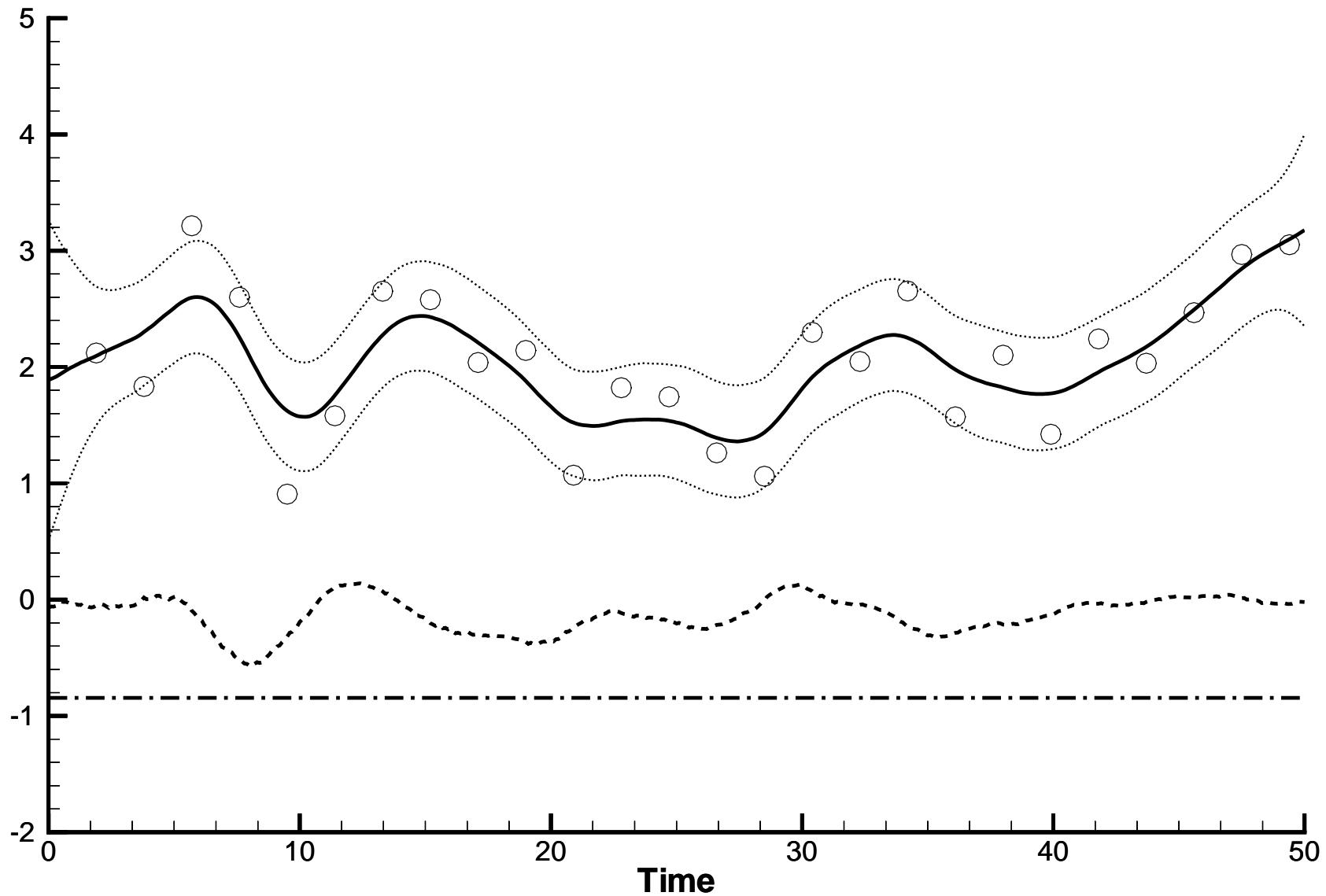
Estimate and model error, EnKS



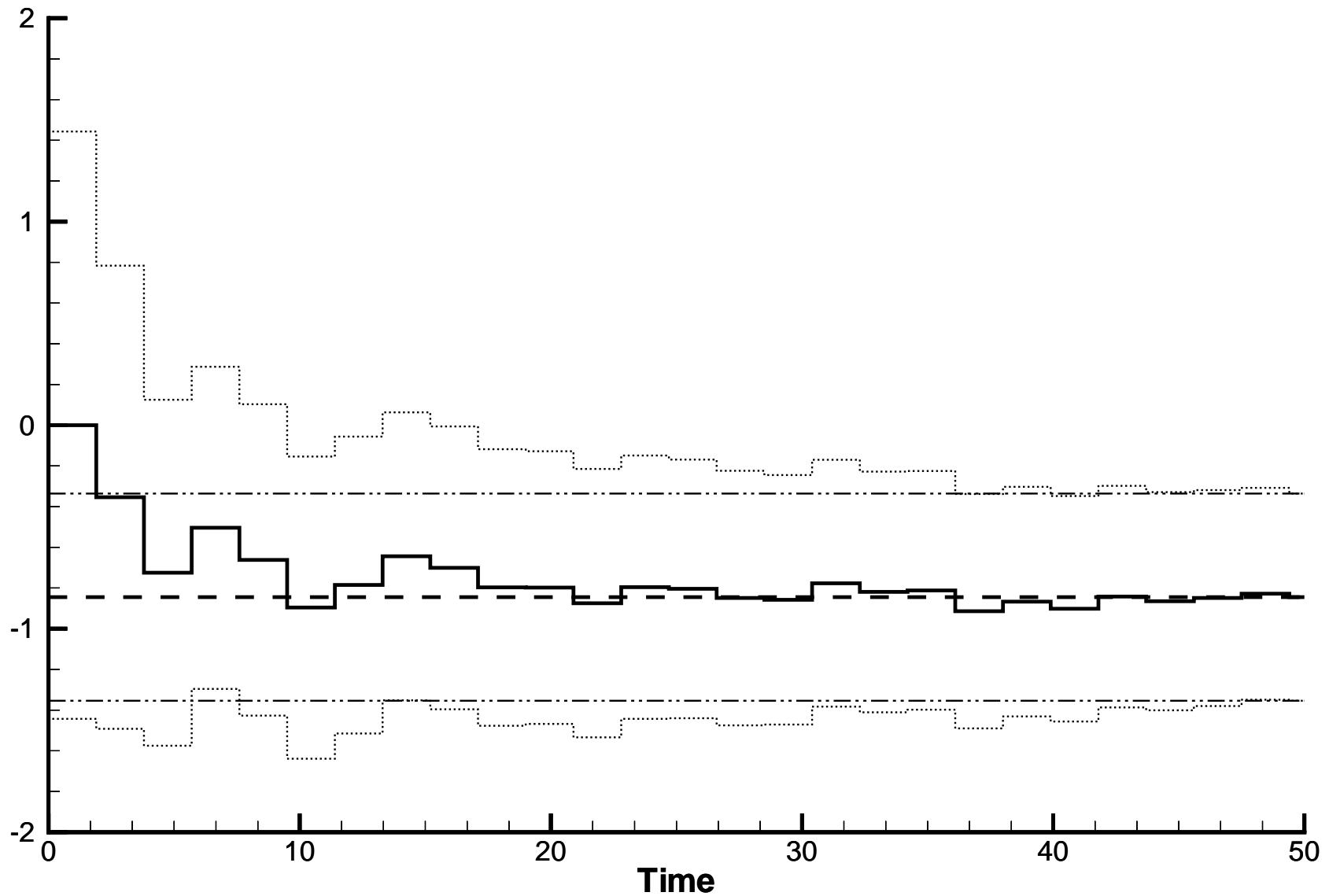
Estimate, model error and β , EnKF



Estimate, model error and β , EnKS



Estimated β and std dev



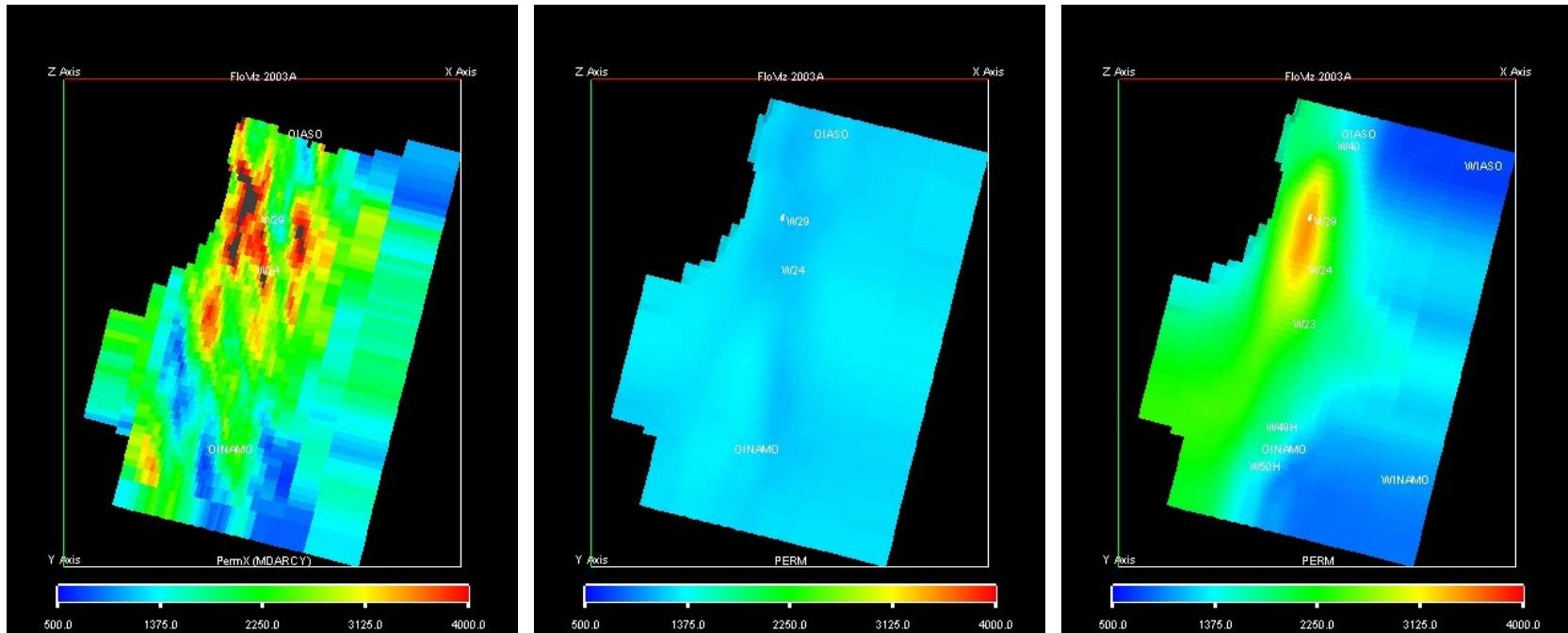
Summary

- The EnKF and EnKS can handle time correlated model errors.
- The EnKF and EnKS can be used for parameter estimation by augmenting the model state with the unknown parameters.

Oil reservoir application

- The EnKF has been used to estimate permeability in a reservoir simulation model using production well data.
- The greatest uncertainty is the reservoir permeability and porosity.
- The well data consists of pressures and oil, gas and water production rates.

Estimated bias and std dev



Local analysis

- May be useful when m is large or $n \gg N$.
- Analysis is computed grid point by grid point using only nearby measurements.
- Inverts many small matrices instead of one large.
- Different X for each grid point, thus, allows us to reach a larger class of solutions.
- Suboptimal solution which introduces nondynamical modes.

Nonlinear measurements

- Measurement equation

$$\mathbf{d} = \mathbf{h}(\boldsymbol{\psi}) + \boldsymbol{\epsilon}.$$

- Define ensemble of model prediction of the measurements

$$\widehat{\mathbf{A}} = (\mathbf{h}(\boldsymbol{\psi}_1), \dots, \mathbf{h}(\boldsymbol{\psi}_N)), \in \mathbb{R}^{\hat{m} \times N}.$$

- The analysis then becomes

$$\mathbf{A}^a = \mathbf{A} + \mathbf{A}' \widehat{\mathbf{A}}'^T \left(\widehat{\mathbf{A}}' \widehat{\mathbf{A}}'^T + \mathbf{E} \mathbf{E}^T \right)^{-1} (\mathbf{D} - \widehat{\mathbf{A}}),$$

- Analysis based on covariances between $\mathbf{h}(\boldsymbol{\psi})$ and $\boldsymbol{\psi}$.

TOPAZ

- Operational ocean prediction system for the Atlantic and Arctic oceans.
- topaz.nersc.no
- Based on HYCOM
- State space is 26 500 000 unknowns.
- Ensemble size is 100.
- Assimilates SSH, SST, Ice concentration, and ARGO T&S profiles.
- Total of more than 10 000 measurements in each assimilation step.
- Uses local analysis as well as nonlinear measurements.

TOPAZ

