

# EnKF overview

Jeff Whitaker\*

NOAA Earth System Research Lab

[<jeffrey.s.whitaker@noaa.gov>](mailto:jeffrey.s.whitaker@noaa.gov)

\* GSI EnVar development team: Daryl Kleist, Dave Parrish, John Derber (NCEP) and Xuguang Wang (OU)

# GSI Var cost function

$$J_{3DVAR}(\mathbf{x}') = \frac{1}{2}(\mathbf{x}')^T \mathbf{B}_f^{-1}(\mathbf{x}') + \frac{1}{2}(\mathbf{H}\mathbf{x}' - \mathbf{y}')^T \mathbf{R}^{-1}(\mathbf{H}\mathbf{x}' - \mathbf{y}')$$

$J$  : Penalty (Fit to background + Fit to observations)

$\mathbf{x}'$  : Analysis increment ( $\mathbf{x}^a - \mathbf{x}^b$ ) ; where  $\mathbf{x}^b$  is a background

$\mathbf{B}_f$  : Background error covariance

$\mathbf{H}$  : Observations (forward) operator

$\mathbf{R}$  : Observation error covariance (Instrument + representativeness)

$\mathbf{y}' = \mathbf{y}^o - \mathbf{H}\mathbf{x}^b$ , where  $\mathbf{y}^o$  are the observations

Cost function ( $J$ ) is minimized to find solution,  $\mathbf{x}'$  [ $\mathbf{x}^a = \mathbf{x}^b + \mathbf{x}'$ ]

# GSI hybrid ensemble Var cost function

$$\mathbf{J}_{\text{hybrid}}(\mathbf{x}') = \frac{\beta}{2}(\mathbf{x}')^T \mathbf{B}_f^{-1}(\mathbf{x}') + \frac{1-\beta}{2}(\mathbf{x}')^T \mathbf{B}_{\text{ens}}^{-1}(\mathbf{x}') + \frac{1}{2}(\mathbf{H}\mathbf{x}' - \mathbf{y}')^T \mathbf{R}^{-1}(\mathbf{H}\mathbf{x}' - \mathbf{y}')$$

$\mathbf{B}_f$  : (Fixed) background-error covariance (estimated offline)

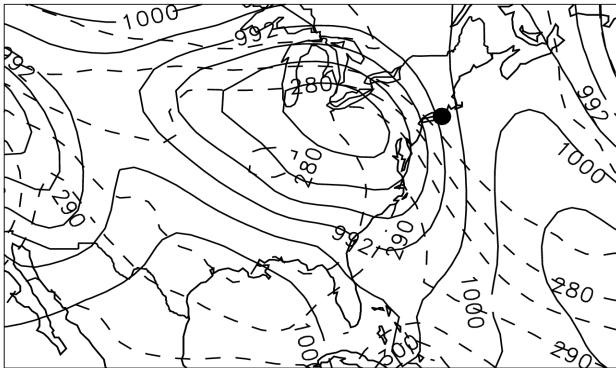
$\mathbf{B}_{\text{ens}}$  : (Flow-dependent) background-error covariance (estimated from ensemble)

$\beta$ : Weighting factor (0.25 means total  $\mathbf{B}$  is  $\frac{3}{4}$  ensemble).

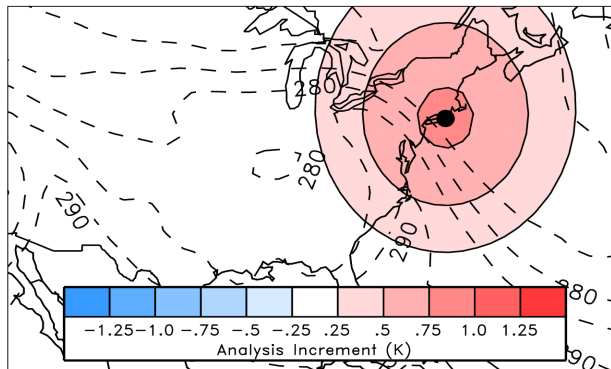
# What does $B_{ens}$ do?

Temperature observation near a warm front

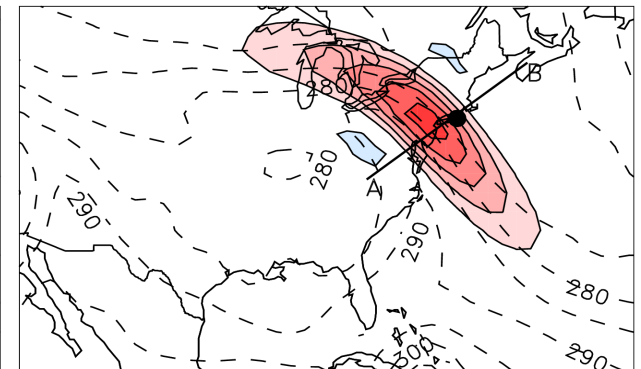
1000 hPa temperature (K) and surface pressure (hPa)



Increment (all static)

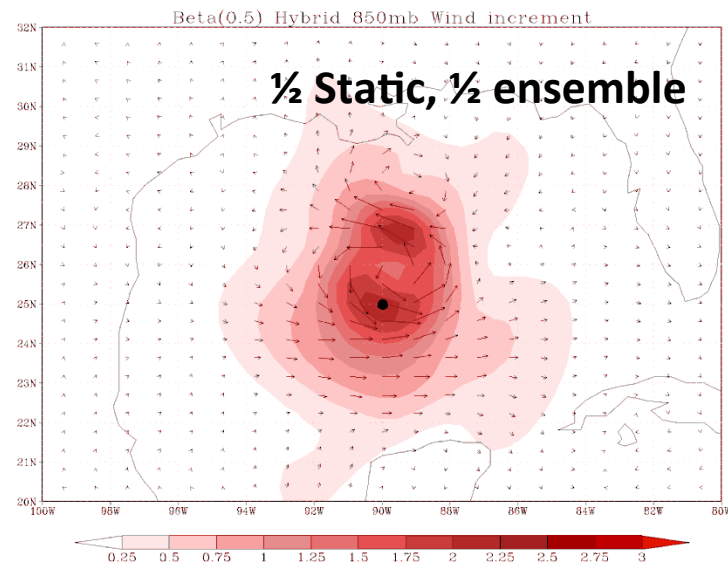
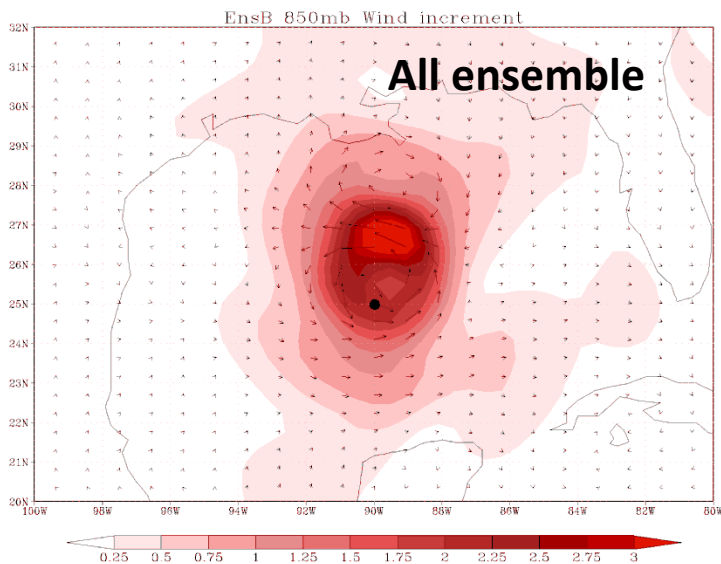
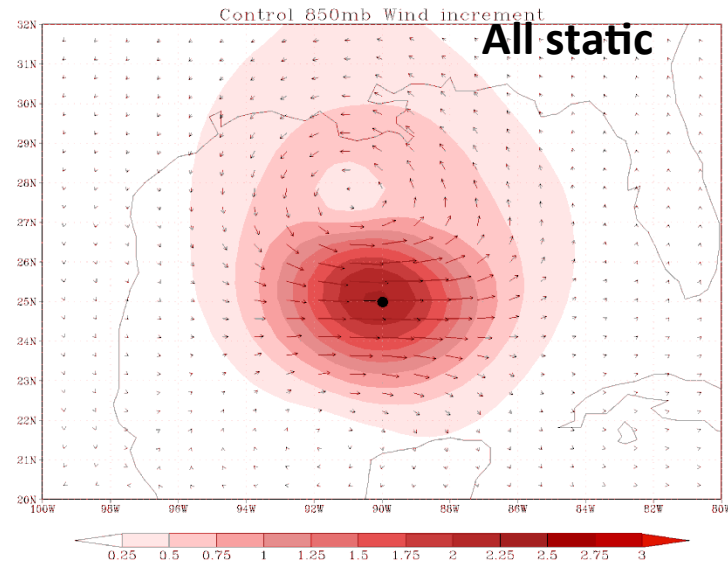
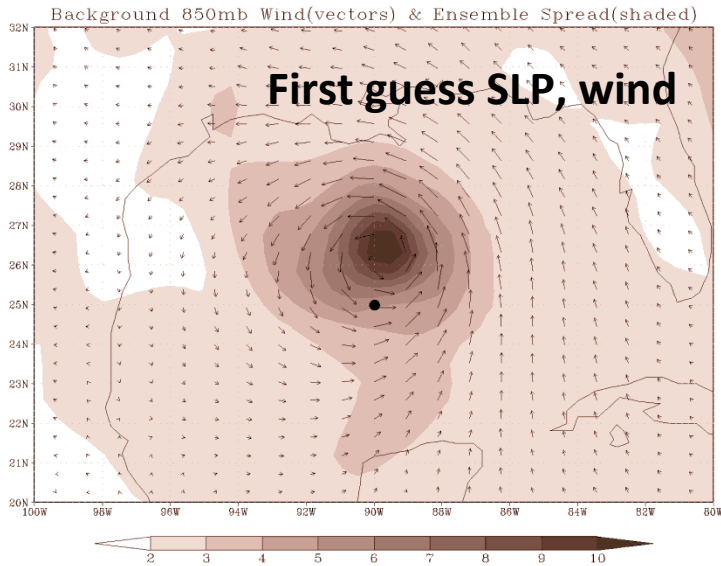


Increment (all ensemble)



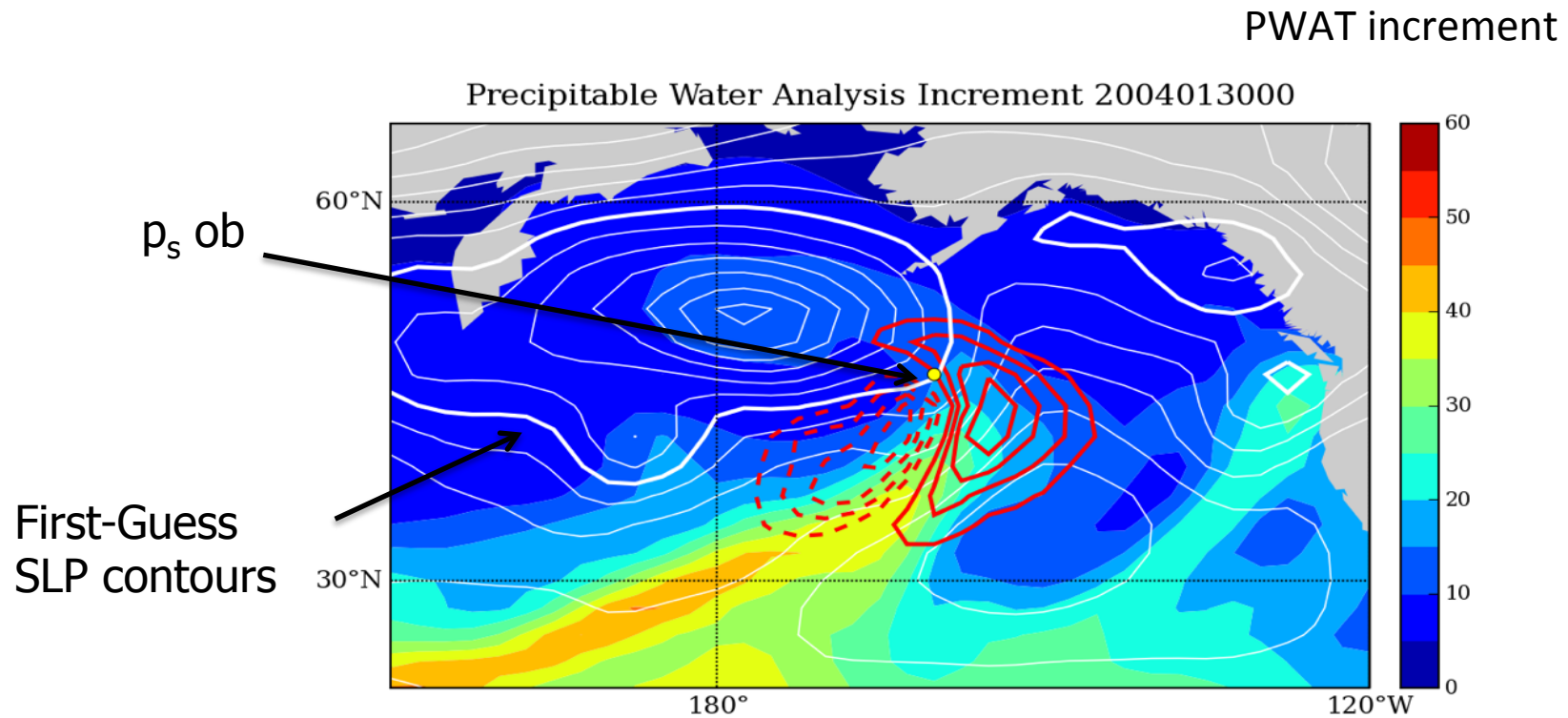
# What does $B_{ens}$ do?

Zonal wind observation near a hurricane (Ike)



# What does $B_{ens}$ do?

Surface pressure observation near an “atmospheric river”



***3Dvar increment would be zero!***  
***(cross-variable covariances hard to model with static  $B_f$ )***

## What does $\mathbf{B}_{\text{ens}}$ do?

- Adds flow-dependence to analysis increments.
- Sparse observations near coherent dynamical features used more effectively.
- Changes in the observing network can be captured in background-error variance.
- ***More information extracted from observations => More skillful forecasts***

## So what's the catch?

- Need an ensemble (fairly large) that accurately represents the uncertainty in the first-guess forecast.
- “Fairly large” means  $O(50-100)$  -- smaller ensembles will have large sampling errors (and more weight will have to be given to  $\mathbf{B}_f$ ). Expensive to run.
- The GSI variational system does not provide the ensemble – it provides an analysis that can be interpreted as the ensemble mean, given an ensemble that represents forecast uncertainty.
- In NCEP operations, an “Ensemble Kalman Filter” (EnKF) is used to generate the background ensemble.



# Data assimilation terminology

- $\mathbf{y}$  : Observation vector (weather balloons, satellite radiances, etc.)
- $\mathbf{x}$  : the state of the atmosphere as represented by the model
- $\mathbf{x}^b$  : Background state vector (“prior”)
- $\mathbf{x}^a$  : Analysis state vector (“posterior”)
- $\mathbf{H}$  : (hopefully linear) operator to convert model state  $\rightarrow$  observation location & type
- $\mathbf{R}$  : Observation - error covariance matrix
- $\mathbf{P}^b$  : Background - error covariance matrix
- $\mathbf{P}^a$  : Analysis - error covariance matrix

## From Bayes theorem to 4DVar and the (Ensemble) Kalman Filter

$$p(\mathbf{x}|\mathbf{y}) \propto \exp \left( -(\mathbf{x} - \mathbf{x}_b)^T \mathbf{P}^b{}^{-1} (\mathbf{x} - \mathbf{x}_b) - (\mathbf{y} - \mathbf{H}\mathbf{x})^T \mathbf{R}^{-1} (\mathbf{y} - \mathbf{H}\mathbf{x}) \right)$$

Variational methods maximize the posterior PDF to find the state trajectory  $\mathbf{x}$  that best fits the obs  $\mathbf{y}$  in a least-squares sense. In practice, this is done by minimizing a cost function, which is what's inside the *exp*:

$$J(\mathbf{x}) \propto (\mathbf{x} - \mathbf{x}_b)^T \mathbf{P}^b{}^{-1} (\mathbf{x} - \mathbf{x}_b) + (\mathbf{y} - \mathbf{H}\mathbf{x})^T \mathbf{R}^{-1} (\mathbf{y} - \mathbf{H}\mathbf{x})$$

The minimum can be found analytically if  $\mathbf{H}$  is linear (see Lorenc 1986 *QJRM*S for the algebra). This gives the equations for the Kalman Filter

$$\begin{aligned} \mathbf{x}_a &= \mathbf{x}_b + \mathbf{K} (\mathbf{y} - \mathbf{H}\mathbf{x}_b), \quad \mathbf{P}^a = (\mathbf{I} - \mathbf{K}\mathbf{H}) \mathbf{P}^b \\ \mathbf{K} &= \mathbf{P}^b \mathbf{H}^T (\mathbf{H}\mathbf{P}^b \mathbf{H}^T + \mathbf{R})^{-1} \end{aligned}$$

- *Matrix  $\mathbf{P}^b$  is too big for any computer, covariance update step impractical.*
- *Instead, represent PDFs of  $\mathbf{x}$  and  $\mathbf{y}$  by an ensemble, compute sample estimate of  $\mathbf{P}^b$  and  $\mathbf{x}_b$ . Evolve the sample, not the full covariance. **EnKF** gives same result as full KF if ensemble size becomes infinite.*

## Computational shortcuts in EnKF:

### (1) Simplifying Kalman gain calculation

$$\mathbf{K} = \mathbf{P}^b H^T \left( H \mathbf{P}^b H^T + \mathbf{R} \right)^{-1}$$

$$\text{define } \overline{H\mathbf{x}^b} = \frac{1}{m} \sum_{i=1}^m H\mathbf{x}_i^b$$

$$\mathbf{P}^b H^T = \frac{1}{m-1} \sum_{i=1}^m \left( \mathbf{x}_i^b - \overline{\mathbf{x}^b} \right) \left( H\mathbf{x}_i^b - \overline{H\mathbf{x}^b} \right)^T$$

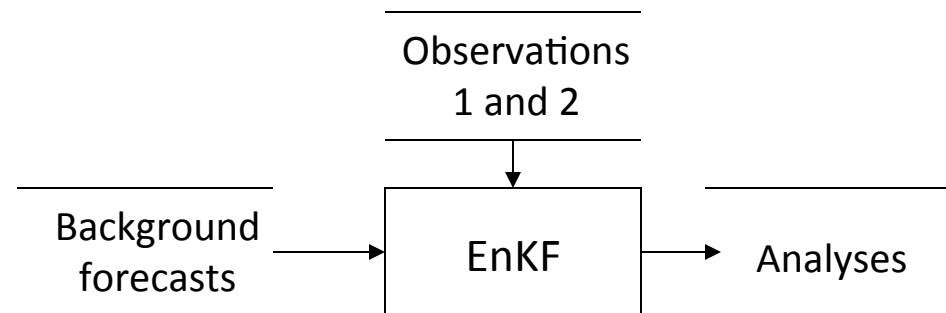
$$H \mathbf{P}^b H^T = \frac{1}{m-1} \sum_{i=1}^m \left( H\mathbf{x}_i^b - \overline{H\mathbf{x}^b} \right) \left( H\mathbf{x}_i^b - \overline{H\mathbf{x}^b} \right)^T$$

The key here is that the huge matrix  $\mathbf{P}^b$  is never explicitly formed

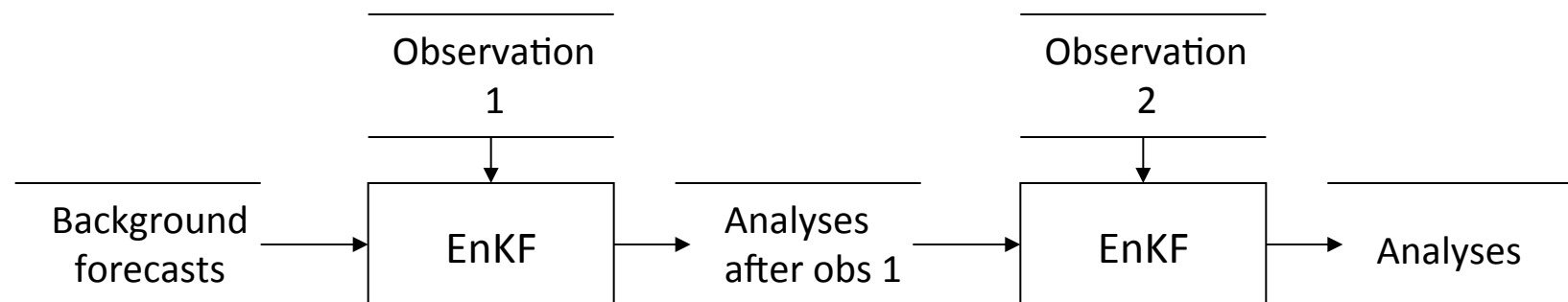
# Computational shortcuts in EnKF:

(2) serial processing of observations (requires observation error covariance  $\mathbf{R}$  to be diagonal)

Method 1



Method 2



# The serial EnKF – a recipe

Given a single ob  $y^o$  with expected error variance  $R$ , an ensemble of model forecasts  $\mathbf{x}^b$  (model priors), and an ensemble of predicted observations  $\mathbf{y}^b = \mathbf{H}\mathbf{x}^b$  (observation priors):

**Step 1:** Update observation priors.

$$(1a) \quad \bar{\mathbf{y}}_a = (1 - K)\bar{\mathbf{y}}_b + Ky^o \quad \text{update for ob prior means}$$

$$(1b) \quad \mathbf{y}'_a = \sqrt{(1 - K)}\mathbf{y}'_b \quad \text{rescaling of ob prior perturbations}$$

where the scalar  $K = \text{var}(\mathbf{y}^b)/(\text{var}(\mathbf{y}^b) + R)$ , overbar denotes means, prime denotes perturbations, superscript  $b$  denotes prior,  $a$  denotes analysis.

**Linear interpolation between observation and observation prior mean with weight  $K$  ( $0 \leq K \leq 1$ ), rescaling of observation prior ensemble so posterior variance is consistent with Kalman filter, i.e.  $\text{var}(\mathbf{y}^a) = (1 - K) \text{var}(\mathbf{y}^b) = \text{var}(\mathbf{y}^b)R/(\text{var}(\mathbf{y}^b) + R)$ .**

when  $\text{var}(\mathbf{y}^b) \ll R$ , all weight given to prior.

when  $\text{var}(\mathbf{y}^b) \gg R$ , all weight given to observation.

## The serial EnKF – a recipe (2)

**Step 2:** Update model priors.

Let  $\Delta\mathbf{x} = \mathbf{x}^a - \mathbf{x}^b$  be analysis increment for model priors,  $\Delta\mathbf{y} = \mathbf{y}^a - \mathbf{y}^b$  is analysis increment for observation priors.

$$(2) \quad \Delta\mathbf{x} = \mathbf{G}\Delta\mathbf{y} \quad \textit{computation of increments to model prior}$$

where  $\mathbf{G} = \text{cov}(\mathbf{x}^b, \mathbf{y}^{bT}) / \text{var}(\mathbf{y}^b)$

***Linear regression of model priors on observation priors.***

Only changes model priors when  $\mathbf{x}^b$  and  $\mathbf{y}^b$  are correlated within the ensemble.

If there is more than one ob to be assimilated, the observation priors for other (not yet assimilated) obs ( $\mathbf{Y}^b$ ) should be also be updated using (2) with  $\Delta\mathbf{x}$  replaced by  $\Delta\mathbf{Y}$ . Next iteration, replace  $\mathbf{y}^b$  with next column of  $\mathbf{Y}^b$ , removing that column from  $\mathbf{Y}^b$ . After each iteration the model priors and observation priors are set to the latest analysis values ( $\mathbf{x}^a$  replaces  $\mathbf{x}^b$ ,  $\mathbf{Y}^a$  replaces  $\mathbf{Y}^b$ ). Continue iterating until  $\mathbf{Y}^b$  is empty.

# Factors limiting EnKF performance

## *1) Treatment of model error*

Must account for the background error covariance associated with “model error” (any difference between simulated and true environment). Methods used so far:

- 1) multiplicative inflation (mult. ens perts by a factor  $> 1$ ).
- 2) additive inflation (random perts added to each member – e.g. differences between 24 and 48-h forecasts valid at the same time).
- 3) model-based schemes (e.g. stochastic kinetic energy backscatter for representing unresolved processes, stochastically perturbed physics tendencies for representing parameterization uncertainty).

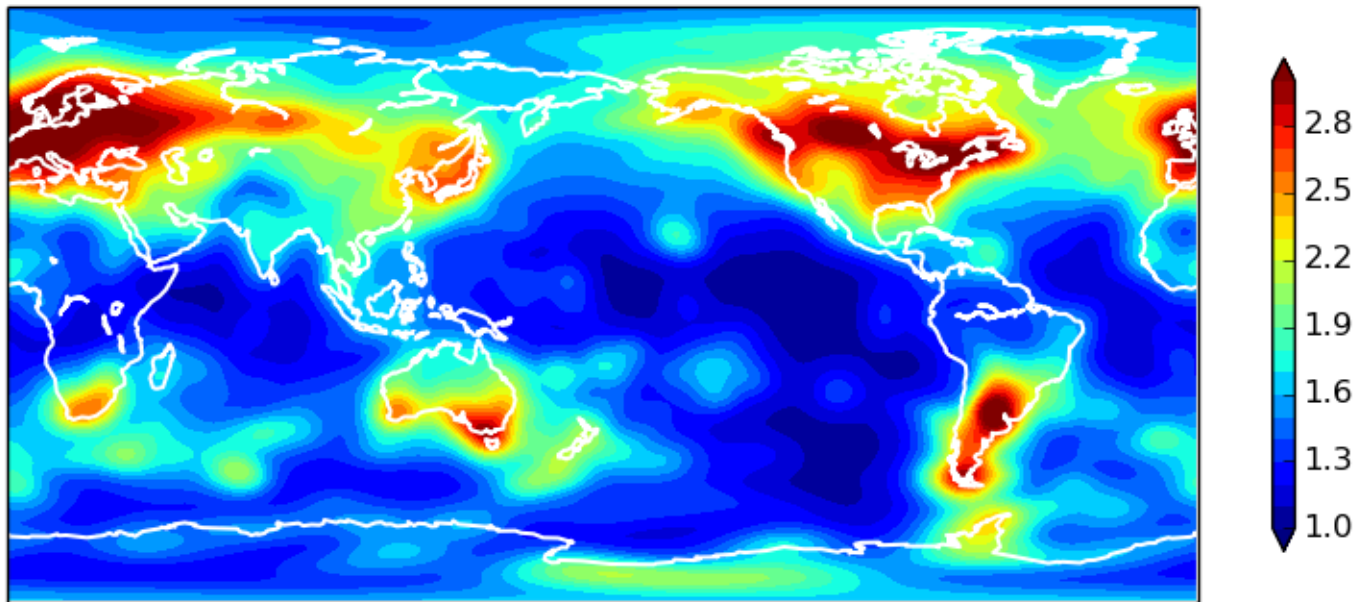
*Opnl NCEP system uses a combination all three.*

# Relaxation To Prior Spread (RTPS) Inflation

Described in [DOI: 10.1175/MWR-D-11-00276.1](https://doi.org/10.1175/MWR-D-11-00276.1)

Inflate posterior spread (std. dev)  $\sigma^a$  back toward prior spread  $\sigma^b$ :  $\sigma^a \leftarrow \alpha\sigma^b + (1-\alpha)\sigma^a$

Equivalent to:  $x'^a \leftarrow x'^a [1 + \alpha(\sigma^b - \sigma^a)/\sigma^a]$



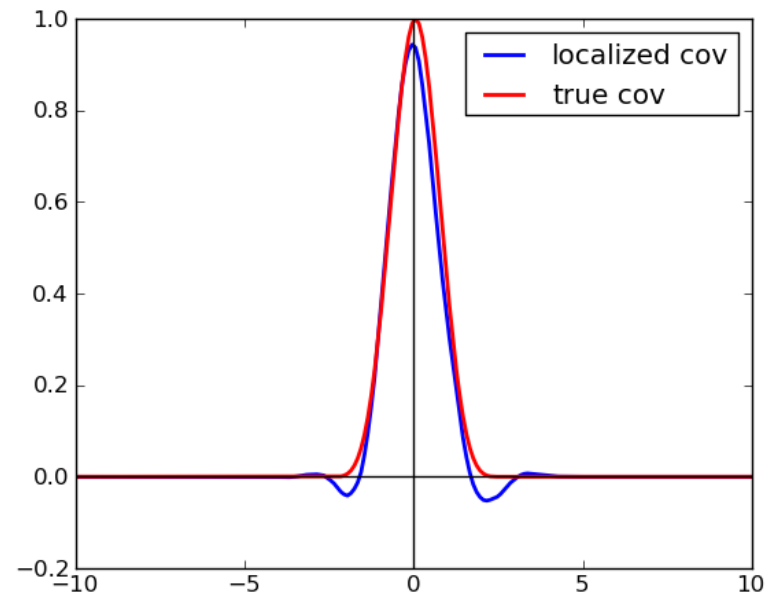
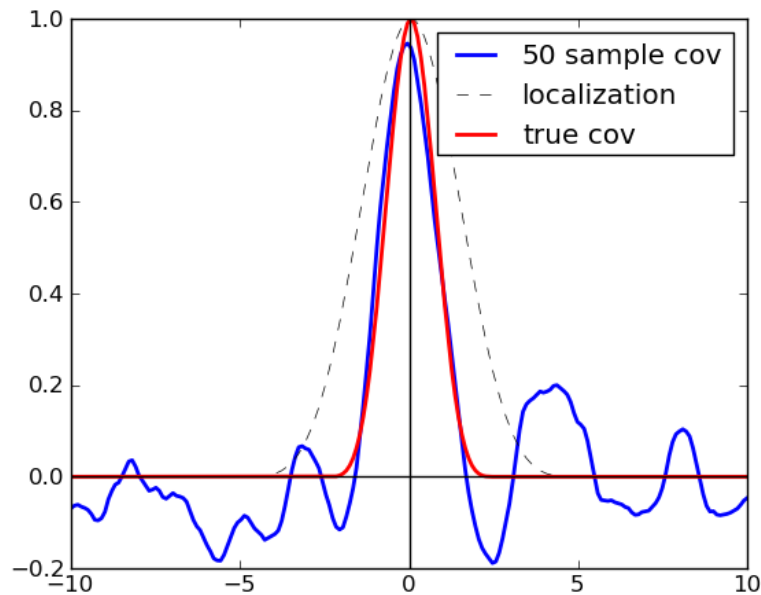


# Factors limiting EnKF performance

## *2) Treatment of sampling error (localization)*

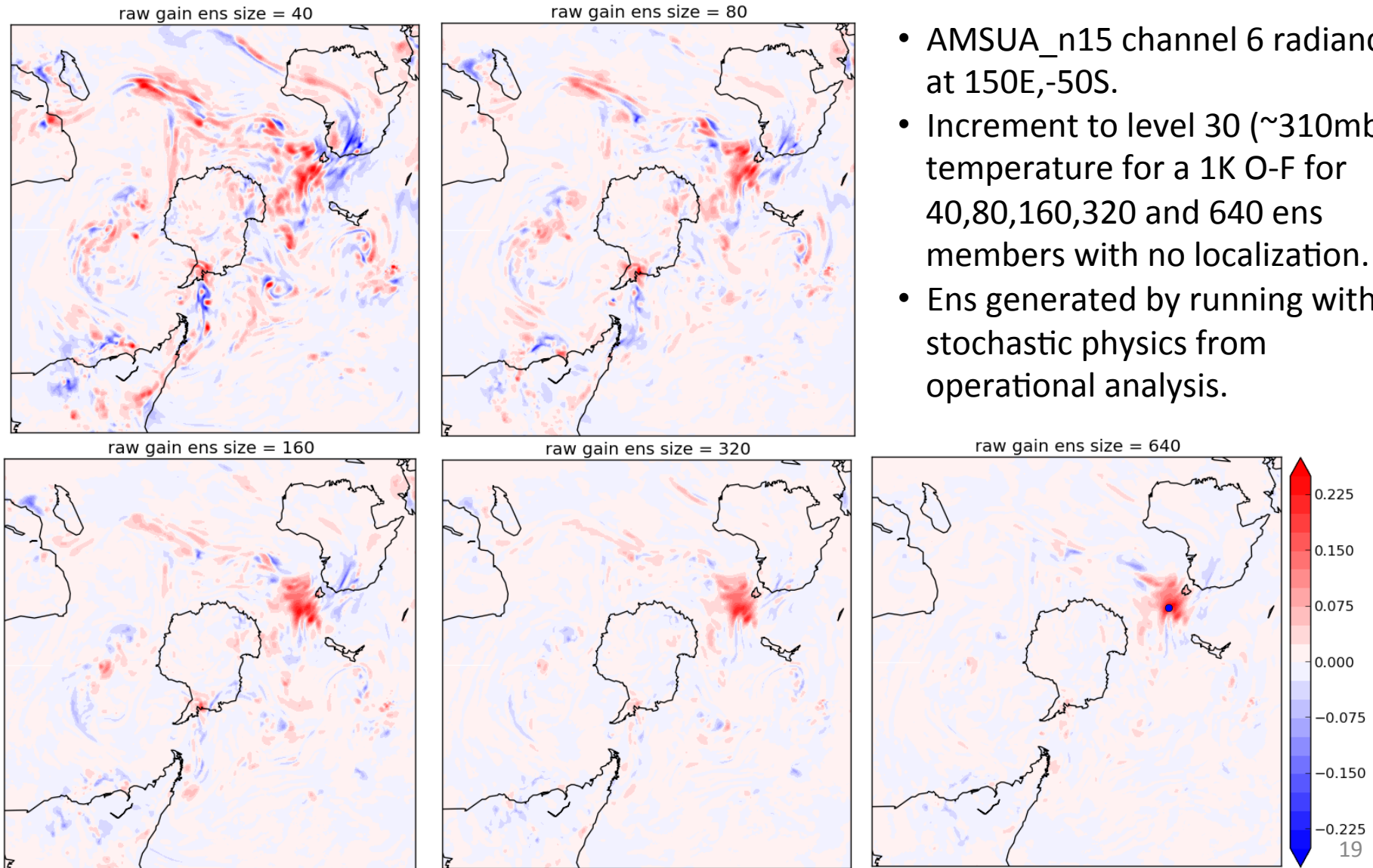
- All EnKF implementations localize the spatial impact of observations on the model state.
- Done by spatially modulating covariance between obs. prior and model state, or by only using observations 'close' to a model state variable to update that variable.
- Needed to account for low rank of ensemble (compared to model state).
- Methods used currently are not flow dependent, and assume there is no sampling error at ob location.

# A simple example of covariance localization

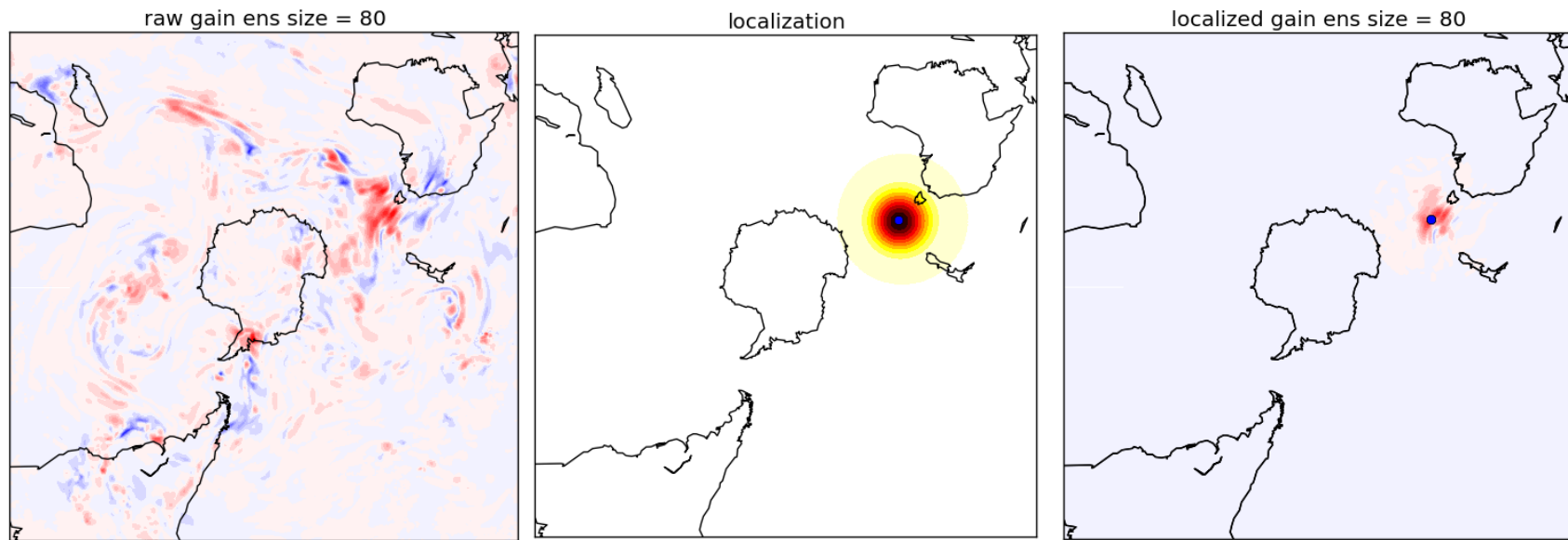


Estimates of covariances from a small ensemble will be noisy, with signal-to-noise small especially when covariance is small

# Covariance localization



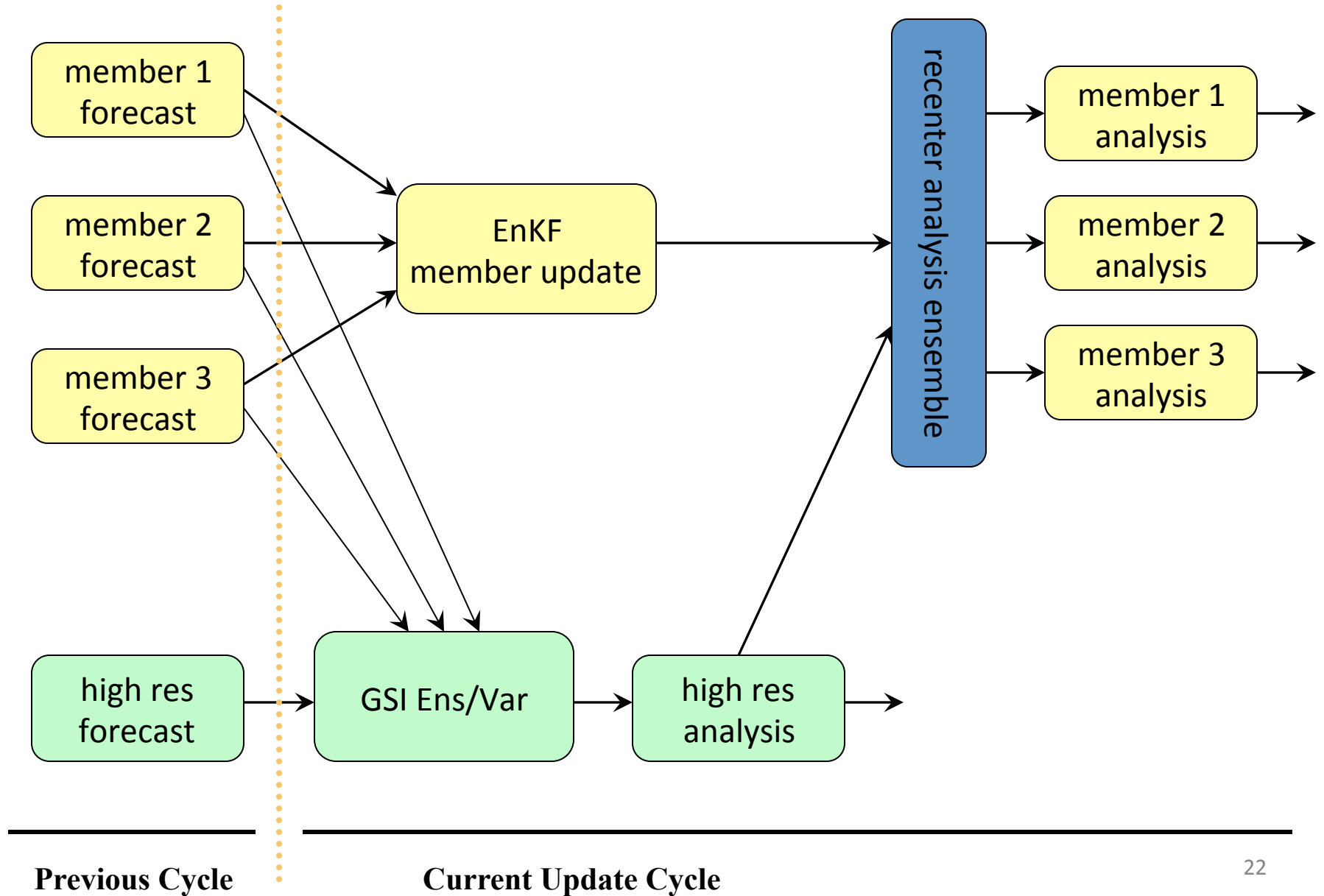
# Covariance Localization



# Why combine EnKF and Var?

Features from EnKF	Features from Var
Can propagate $\mathbf{P}^b$ from across assimilation windows	Treatment of sampling error in ensemble $\mathbf{P}^b$ estimate does not depend on $\mathbf{H}$ .
More flexible treatment of model error (can be treated in ensemble)	Dual-resolution capability – can produce a high-res “control” analysis.
Automatic initialization of ensemble forecasts.	Ease of adding extra constraints to cost function, including a static $\mathbf{P}^b$ component.

# Ensemble-Var workflow



# Summary

- The EnKF uses an ensemble of first-guess forecasts to estimate the background-error covariance. Every ensemble member is updated at each analysis time.
  - Parallel code, scalable out to  $O(1000's)$  of processors as long as number of obs  $\ll$  number of state vars.
  - Requires state vector in model and ob space, plus obs, as input.
  - GSI used to compute forward observation operator (separate step run before EnKF).
- Need to carefully tune localization length scales (depends on model resolution, observing network).
- Ensemble (co)variances must be representative of control forecast error. Treatment of model and sampling error is crucial.

# GSI ensemble 3DVar cost function (with localization)

$$\mathbf{J}_{\text{hybrid}}(\mathbf{x}') = \frac{\beta}{2}(\mathbf{x}')^T \mathbf{B}_f^{-1}(\mathbf{x}') + \frac{1-\beta}{2}(\mathbf{x}')^T (\mathbf{B} \circ \mathbf{S})_{\text{ens}}^{-1}(\mathbf{x}') + \frac{1}{2}(\mathbf{H}\mathbf{x}' - \mathbf{y}')^T \mathbf{R}^{-1}(\mathbf{H}\mathbf{x}' - \mathbf{y}')$$

$\mathbf{B}_f$  : (Fixed) background-error covariance (estimated offline)

$\mathbf{B}_{\text{ens}}$  : (Flow-dependent) background-error covariance (estimated from ensemble). **Schur product with correlation matrix  $\mathbf{S}$  implies localization.**

$\beta$ : Weighting factor (0.25 means total  $\mathbf{B}$  is  $\frac{3}{4}$  ensemble).

Extra parameters control horizontal and vertical scales in  $\mathbf{S}$ .