1966-3

# Fall Colloquium on the Physics of Weather and Climate: Regional Weather Predictability and Modelling

*29 September - 10 October, 2008*

**Atmospheric data assimilation**

Dusanka Zupanski

*Colorado State University
USA*

# Lecture 1: Atmospheric Data Assimilation

**Dusanka Zupanski**
**CIRA/Colorado State University**
**Fort Collins, Colorado**

*Fall Colloquium on the Physics of Weather and Climate: Regional Weather Predictability and Modelling*
*29 September - 7 October, 2008, Trieste, Italy*

Dusanka Zupanski, CIRA/CSU
Zupanski@CIRA.colostate.edu

# OUTLINE

**Lecture 1:**
➢ **What is data assimilation and why is it important?**

➢ **Ensemble Data Assimilation**

**Lecture 2:**
➢ **Maximum Likelihood Ensemble Filter (MLEF)**

➢ **Examples of data assimilation results**

Dusanka Zupanski, CIRA/CSU
Zupanski@CIRA.colostate.edu

# NOTE

If there is enough interest, we will perform simple data assimilation experiments, using Maximum Likelihood Ensemble Filter (MLEF) method and Burgers model, in one of the available Computer Labs.

Interested participants, please sign-in at the end of this class.

Suggested time for Lab exercises is today or Thursday after lunch break. Please write your suggestions next to your name when signing.

Dusanka Zupanski, CIRA/CSU
Zupanski@CIRA.colostate.edu

# Websites

If you would like to perform the exercise on your own, downloaded the MLEF algorithm, via anonymous ftp, from the following location:

ftp://ftp.cira.colostate.edu/Zupanski/ICTP_Lecture_2008/mlef-V1.2.tar

Then, do the following (on a Linux computer or laptop):
tar -xvf mlef-V1.2.tar
cd mlef-V1.2
and read the text file called "README" for further instructions

More information about the MLEF can be found at

http://www.cira.colostate.edu/projects/ensemble

Dusanka Zupanski,  CIRA/CSU
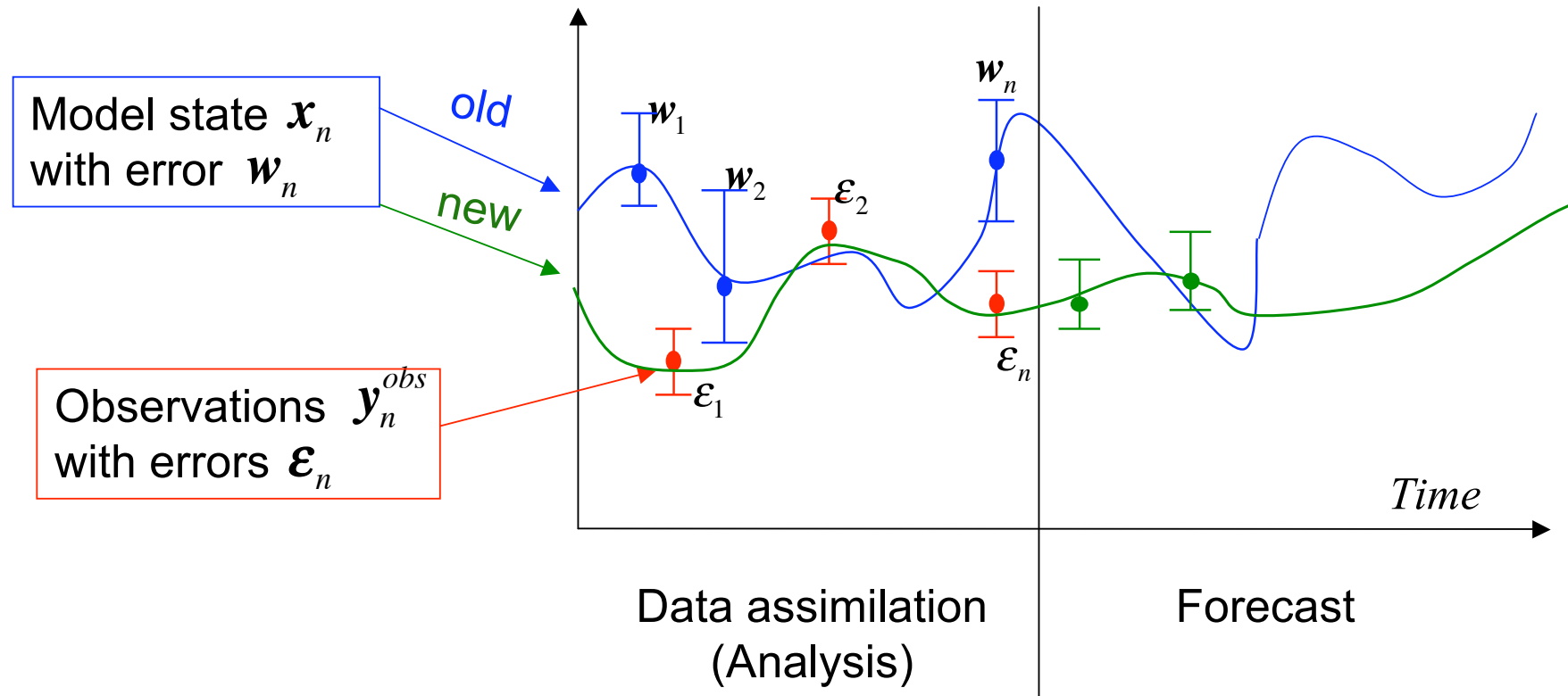Zupanski@CIRA.colostate.edu

# Questions

Please feel free to ask questions during and after my presentations (until Saturday, when I will be leaving).

You are also welcome to send me an e-mail.

Dusanka Zupanski, CIRA/CSU
Zupanski@CIRA.colostate.edu

# What is Data Assimilation?

➢ Data assimilation (DA) started as Atmospheric DA in 1960s-1970s (you might have heard of Gandin, who was a pioneer in the atmospheric DA).

➢ Initially, atmospheric DA was used to define "objective" **initial conditions** for Numerical Weather Predictions (NWP) models. In these early days, the term "**objective analysis**" was used.

➢ Today, DA is used to define "**optimal**" initial conditions for NWP, climate, hydrological, oceanic, land, canopy, pollution transport and other bio-geophysical models (it is not only atmospheric DA anymore).

➢ Advanced DA methods, available today, can also estimate and correct model errors (**model error estimation**) and define optimal empirical parameters (**parameter estimation**).

Dusanka Zupanski,  CIRA/CSU
Zupanski@CIRA.colostate.edu

# What is Data Assimilation?



New (posterior) model state fits the observations better than the old (prior) model state: Data Assimilation (DA) "**optimally**" combines **observations** and a **dynamical forecast model** in order to improve estimates of the current state (**analysis**) and the future state (**forecast**) of a dynamical system of interest (e.g., atmosphere, ocean, land).

# How is this done?

**There are several steps:**

**Step 1: Define forecast model time evolution equation:**

$$x_n = M_{n,n-1}(x_{n-1}) + w_n$$

$n$ - Time step index (denoting model time steps)

$M$ - Dynamical model for model state evolution (e.g., Eta model)

$x$ - Model state vector of dim *Nstate*

$w$ - Model error vector of dim *Nstate*

**Step 2: Define time evolution equation for the observations:**

$$y_k = H_k(x_k) + \varepsilon_k$$

$k$ - Time step index (denoting observation times)

$y$ - Observations vector of dim *Nobs*

$\varepsilon$ - Observation error

$H$ - Observation operator

Dusanka Zupanski, CIRA/CSU
Zupanski@CIRA.colostate.edu

**Step 3: Combine the two equations in the probability space, by employing the Bayesian conditional Probability Density Function (PDF):**

$$PDF = p(x|y) = \frac{p(y|x)p(x)}{p(y)}$$

In order to proceed with the next step, we assume that $p(x)$ and $p(y|x)$ are known. Typical assumption is that both PDFs are Gaussian and are mutually independent, thus we can write

$$p(x) \sim exp\left[-\frac{1}{2}(x - x_b)^T P_f^{-1}(x - x_b)\right]$$

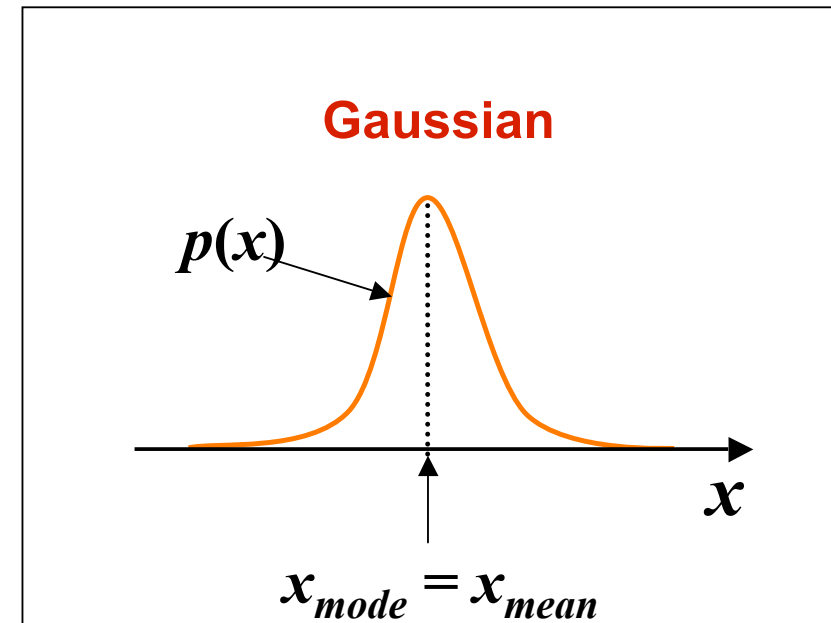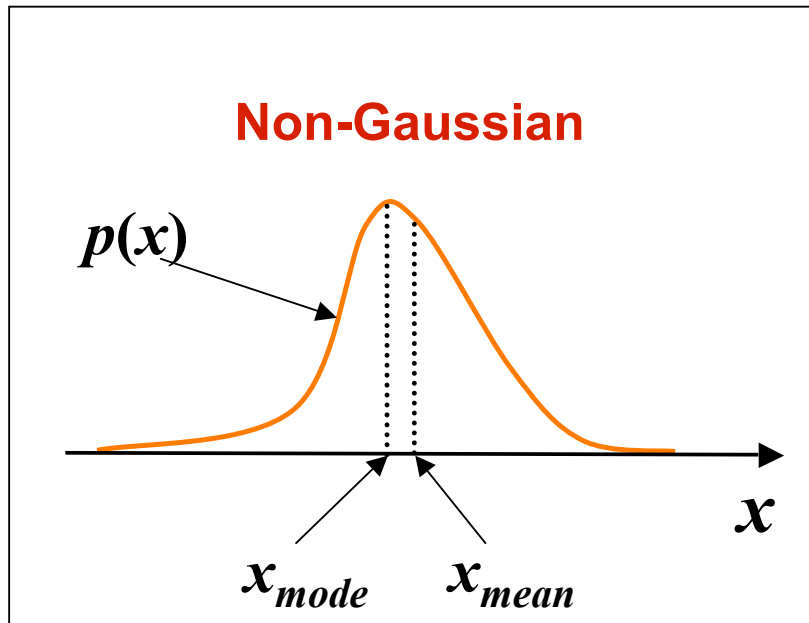$P_f$ is forecast error covariance matrix (prescribed or calculated)

$$p(y|x) \sim exp\left\{-\frac{1}{2}\left[y - H\left[M(x)\right]\right]^T R^{-1}\left[y - H\left[M(x)\right]\right]\right\}$$

$R$ is observation error covariance matrix (prescribed)

**Step 4: Finally, obtain the minimum variance (mean) or maximum likelihood (mode) solution by minimizing the following cost function J**

$$J = \frac{1}{2}[x - x_b]^T P_f^{-1}[x - x_b] + \frac{1}{2}[H[M(x)] - y]^T R^{-1}[H[M(x)] - y]$$

MEAN vs. MODE

**Non-Gaussian**

$p(x)$

$x_{mode}$  $x_{mean}$

$x$

**Gaussian**

$p(x)$

$x_{mode} = x_{mean}$

$x$

Minimum variance estimate= Maximum likelihood estimate!

For Gaussian PDFs mean=mode, thus the solution is easy to obtain. Or is it?

Dusanka Zupanski, CIRA/CSU
Zupanski@CIRA.colostate.edu

**Even under the Gaussian error assumption, there are different ways to obtain the optimal solution (analysis):**

1. Variational methods
   - 3d-var
   - 4d-var

2. Kalman Filter - like methods
   - Classical Kalman Filter (KF)
   - Extended Kalman Filter (EKF)
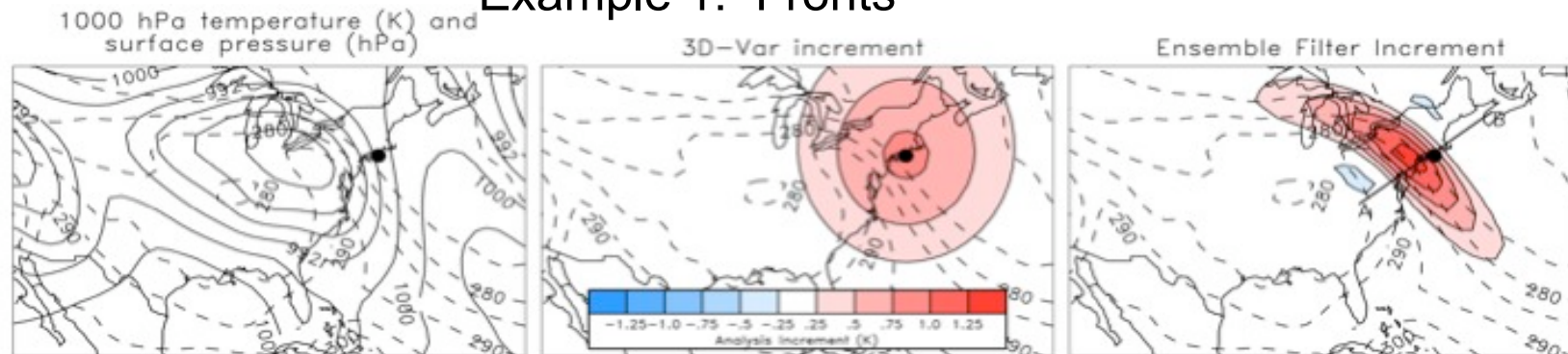   - Ensemble Kalman Filter (EnKF)

The most significant difference between the two groups of DA methods is in the way how $P_f$ is defined:

1. In variational methods $P_f$ is prescribed.

2. In KF-like methods, $P_f$ is evolving in time according to model dynamics.

Dusanka Zupanski, CIRA/CSU
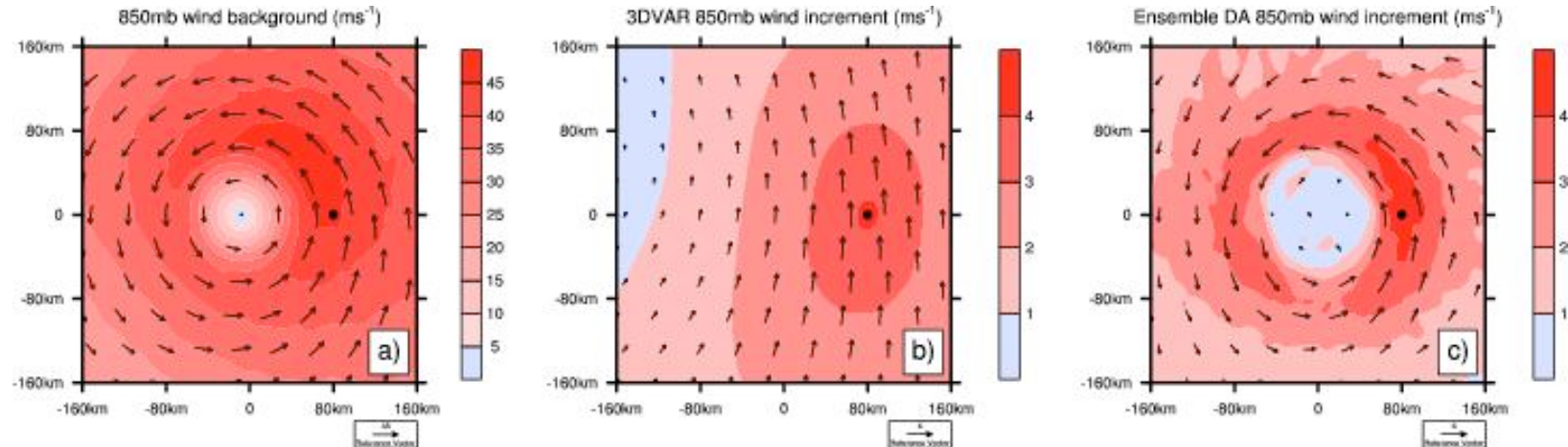Zupanski@CIRA.colostate.edu

# Why is it important to have time evolving $P_f$?

(From Whitaker et al., THORPEX web-page)

Example 1: Fronts



Example 2: Hurricanes



Ensemble methods produce realistic $P_f$. Nevertheless, covariance localization is sometimes needed to eliminate spurious long-distance correlations (covariance localization was applied in this example).

## More specifically

❑ **3d-var** method employs a prescribed forecast error covariance and never evolves it in time.

❑ **4d-var** method evolves the forecast error covariance in time, but only until the end of a data assimilation interval. In each new data assimilation interval it starts with the same, prescribed forecast error covariance.

❑ **Kalman Filter** (KF) Does evolve the forecast error covariance in time and from one data assimilation cycle to another, but it is computationally too expensive for applications to complex atmospheric models.
⇓

Ensemble data assimilation (e.g., EnKF) is a practical alternative to KF, applicable to most complex atmospheric models (Eta, WRF).
⇓

A bonus benefit: The ensemble data assimilation is easier to deal with, since we do not need to define $P_f$; the ensembles will do it for us.
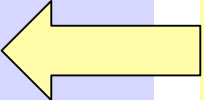
# More reasons for ensemble data assimilation

❑ Ensemble data assimilation provides initial ensemble perturbations for ensemble forecasting.

❑There is no need to develop adjoint models.

❑ As will be shown in Lecture 2, the ensemble data assimilation methods can deal with non-linear and discontinuous (on/off) functions better than variational methods

Dusanka Zupanski, CIRA/CSU
Zupanski@CIRA.colostate.edu
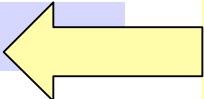
# There are many different versions of EnKF

➢ Monte Carlo EnKF (Evensen 1994; 2003)

➢ EnKF (Houtekamer et al. 1995; 2005; First operational version)

➢ Hybrid EnKF (Hamill and Snyder 2000)

➢ EAKF (Anderson 2001)

➢ ETKF (Bishop et al. 2001)    ⟵ Minimum variance solution

➢ EnSRF (Whitaker and Hamill 2002)

➢ LEKF (Ott et al. 2004)

➢ LETKF (Hunt et al. 2007)

➢ SEEK filter (Pham et al. 1998)

➢MLEF (Zupanski 2005; Zupanski and Zupanski 2006)    ⟵ Maximum likelihood solution

Why maximum likelihood solution? Because it is more adequate for non-linear and discontinuous models (Zupanski et al. 2008) and for problems involving non-Gaussian PDFs (e.g., Fletcher and Zupanski 2006).

# Kalman filter solution

**Analysis step:**

$$x_a = x_b + P_f H^T (HP_f H^T + R)^{-1} [y - H(x_b)]$$

$x_a$    - Optimal estimate of $x$ (analysis)

$x_b$    - Background (prior) estimate of $x$

$$P_a = [I - P_f H^T (HP_f H^T + R)^{-1} H] P_f = (I - KH) P_f$$

$P_a$    - Analysis (posterior) error covariance matrix (*Nstate* x *Nstate*)

$K$    - Kalman gain matrix (*Nstate* x *Nobs*)

**Forecast step:**

$$x_0 = x_a$$    ;    $$x_n = M_{n,n-1}(x_{n-1}) + w_n$$

Often neglected

$$P_f = MP_a M^T + Q$$    - Update of forecast error covariance

Dusanka Zupanski, CIRA/CSU
Zupanski@CIRA.colostate.edu

# Ensemble Kalman Filter (EnKF) solution

EnKF was first introduced by Evensen (1994) as a Monte Carlo filter.

## Equations following Evensen (2003)

**Analysis step:**

Analysis solution defined for each ensemble member $i$:

$$x_a^i = x_b^i + P_f^e H^T (HP_f^e H^T + R^e)^{-1}(y^i - H(x_b^i))$$

Analysis ensemble perturbations:

Mean analysis solution:

$$\overline{x_a} = \overline{x_b} + P_f^e H^T (HP_f^e H^T + R^e)^{-1}(\overline{y} - H(\overline{x_b}))$$

$$p_a^i = x_a^i - \overline{x_a}$$

Analysis error covariance in ensemble subspace:

$$\left(P_a\right)^{1/2} = \left[ p_a^1 \quad p_a^2 \quad . \quad p_a^{Nens} \right]$$

$$P_a = \frac{1}{N_{ens} - 1}\left(P_a\right)^{1/2}\left[\left(P_a\right)^{1/2}\right]^T$$

Dusanka Zupanski,  CIRA/CSU
Zupanski@CIRA.colostate.edu

# Ensemble Kalman Filter (EnKF) solution

**Forecast step:**

Ensemble forecasts employing a **non-linear** model $M$

$$x_n^j = M_{n,n-1}(x_{n-1}^j)$$

Non-linear forecast perturbations

$$p_f^i = M(x_a^i) - M(\overline{x_a})$$

Forecast error covariance calculated using ensemble perturbations:

$$\left(P_f\right)^{1/2} = \begin{bmatrix} p_f^1 & p_f^2 & . & p_f^{Nens} \end{bmatrix}$$

$$P_f = \frac{1}{N_{ens}-1}\left(P_f^e\right)^{1/2}\left[\left(P_f^e\right)^{1/2}\right]^T$$

Dusanka Zupanski, CIRA/CSU
Zupanski@CIRA.colostate.edu

# MLEF solution

**Analysis step:**

Analysis solution $x_a$ obtained by minimizing the cost function

$$J = \frac{1}{2}[x - x_b]^T P_f^{-1}[x - x_b] + \frac{1}{2}[H(x) - y]^T R^{-1}[H(x) - y]$$

Analysis error covariance in ensemble subspace:

$$\left(P_a\right)^{1/2} = P_f^{1/2}\left[I + \left(Z(x_a)\right)^T Z(x_a)\right]^{-1/2}$$

$$Z(x) = \left[z_1(x) \quad z_2(x) \quad \cdot \quad \cdot \quad z_{Nens}(x)\right] \; ; \; z_i(x) = R^{-1/2}\left[H(x + p_i^f) - H(x)\right]$$

$$P_a = \left(P_a\right)^{1/2}\left[\left(P_a\right)^{1/2}\right]^T$$

Dusanka Zupanski, CIRA/CSU
Zupanski@CIRA.colostate.edu

# MLEF solution

**Forecast step:**

Ensemble forecasts employing a **non-linear** model $M$

$$x_n^j = M_{n,n-1}(x_{n-1}^j)$$

Non-linear ensemble forecast perturbations

$$p_f^i = M(x_a + p_a^i) - M(x_a)$$

Forecast error covariance calculated using ensemble perturbations:

$$\left(P_f\right)^{1/2} = \left[\begin{array}{cccc} p_f^1 & p_f^2 & \cdot & p_f^{Nens} \end{array}\right]$$

$$P_f = \left(P_f\right)^{1/2}\left[\left(P_f\right)^{1/2}\right]^T$$

Dusanka Zupanski, CIRA/CSU
Zupanski@CIRA.colostate.edu

## Connections between MLEF, KF and 3d-var
*(Zupanski et al. 2007)*

➢ **MLEF=KF**, **if PDFs are Gaussian, observation operator H is linear and differentiable and the MLEF uses a full-rank $P_f$ (Full-rank means *Nens=Nstate*).**

$$x = x_b + \alpha P_f H^T (HP_f^T H^T + R)^{-1}[y - H(x_b)]$$ ; for $\alpha=1$

*MLEF is also applicable to Non-Gaussian PDFs (e.g., Fletcher and Zupanski 2006).*

➢ **MLEF=3d-var**, **if observation operator is differentiable and $P_f$ is full-rank, but not updated from one data assimilation cycle to another. In such case the same cost function is minimized:**

$$J(x) = \frac{1}{2}[x - x_b]^T P_f^{-1}[x - x_b] + \frac{1}{2}[y - H(x)]^T R^{-1}[y - H(x)]$$

*MLEF is also applicable to non-differentiable observation operators (e.g., Zupanski et al. 2008).*

Dusanka Zupanski, CIRA/CSU
Zupanski@CIRA.colostate.edu

# Current status of EnKF applications

➢ EnKF is operational in Canada, since January 2005 (Houtekamer et al. 2005).

➢ Pseudo operational EnKF at University of Washington (Torn and Hakim 2008).

➢ EnKF is better than 3d-var (Meng and Zhang 2008; Whitaker et al., 2008).

➢ Superior performance in application to non-hydrostatic, cloud resolving models (Caya et al. 2005; Xue et al. 2006; Carrio et al. 2008).

➢ Superior performance in applications to non-linear and discontinuous (with on/off switches) models (Zupanski et al. 2008).

➢ Superior performance for ocean (Keppenne et al. 2008), climate (Karspeck and Anderson 2007), and soil hydrology models (Reichle et al. 2007).

Theoretical advantages of ensemble-based DA methods are getting confirmed in an increasing number of practical applications.

There is still many possibilities for further improvements of the ensemble based DA methods.. For example, covariance localization is one of the critical areas that need further improvements.

Dusanka Zupanski, CIRA/CSU
Zupanski@CIRA.colostate.edu

# A hint on how to run the MLEF algorithm

Dusanka Zupanski,  CIRA/CSU
Zupanski@CIRA.colostate.edu