# Sequential methods for ocean data assimilation

## From theory to practical implementations (I)

*P. BRASSEUR*

*CNRS/LEGI, Grenoble, France*

*Pierre.Brasseur@hmg.inpg.fr*

*J. Ballabrera, L. Berline, F. Birol, J.M. Brankart, G. Broquet, V. Carmillet, F. Castruccio, F. Debost, D. Rozier, Y. Ourmières, T.Penduff, J. Verron*

*E. Blayo, S. Carme, Pham D.T., C. Robert*

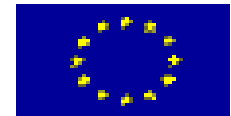*C.E. Testut, B. Tranchant, N. Ferry, E. Remy, M. Benkiran*

*F. Durand, L. Gourdeau, Ch. Maes, P. De Mey*

*A. Barth, C. Raick, M. Grégoire*

*L. Parent*

# Ocean data assimilation

❑ **Data assimilation involves the <u>optimal combination of measurements</u> with the underlying <u>dynamical principles</u> governing the system under observation.**

❑ **Data assimilation can serve several oceanographic objectives:**
- ➢ **Ocean state estimation in space & time (4D) ;**
- ➢ **Detection of model errors ;**
- ➢ **Estimation of budgets & model parameters ;**
- ➢ **Initialisation, prediction, monitoring ;**
- ➢ **Optimal design of complex observation systems ;**
- ➢ **…**

❑ **Theories: optimal _control_ (VAR) and optimal _estimation_ (Kalman)**

# MERCATOR Assimilation Systems

❑ **Incremental implementation strategy**

|  | OI | Kalman filters | 3D/4D-VAR |
|---|---|---|---|
| **Research** | 1993 (SOFA) | 1998 (SEEK) | 1999 (OPAVAR) |
| **R&D** | 1997 | 2002 | 2004 |
| **DEV** | 1999 | 2005 | 2008 ? |
| **OP** | 2001 | 2007 ? | ? |
|  | **SAM-1** | **SAM-2** | **SAM-3** |

❑ **State-of-the-art**

1. **Introduction**
2. **Kalman filter: fundamentals**
3. *Applied* **ocean data assimilation: specific issues**
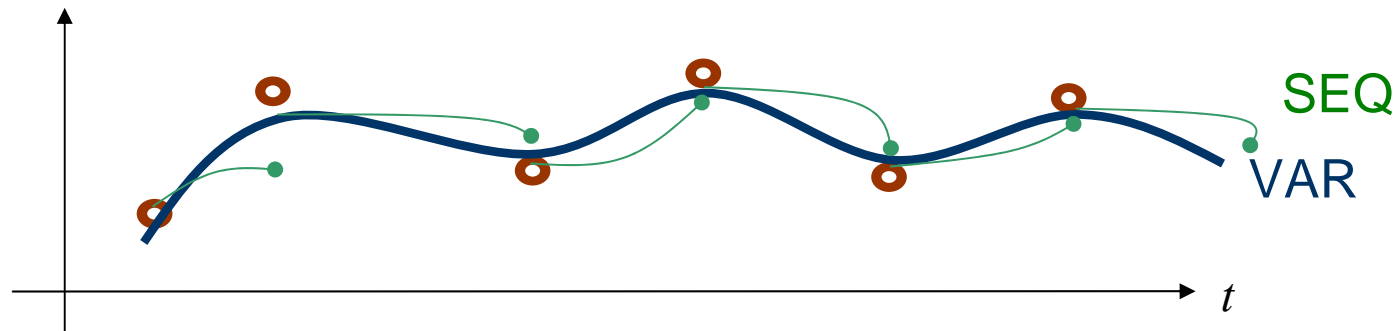4. **Simplifications of the KF – Optimal Interpolation**

❑ **Advanced issues**

5. **Space reduction: state and error sub-spaces**
6. **Low rank filters: SEEK and EnKF**
7. **Objective validation and adaptive schemes**
8. **Improved temporal strategies : FGAT and IAU**
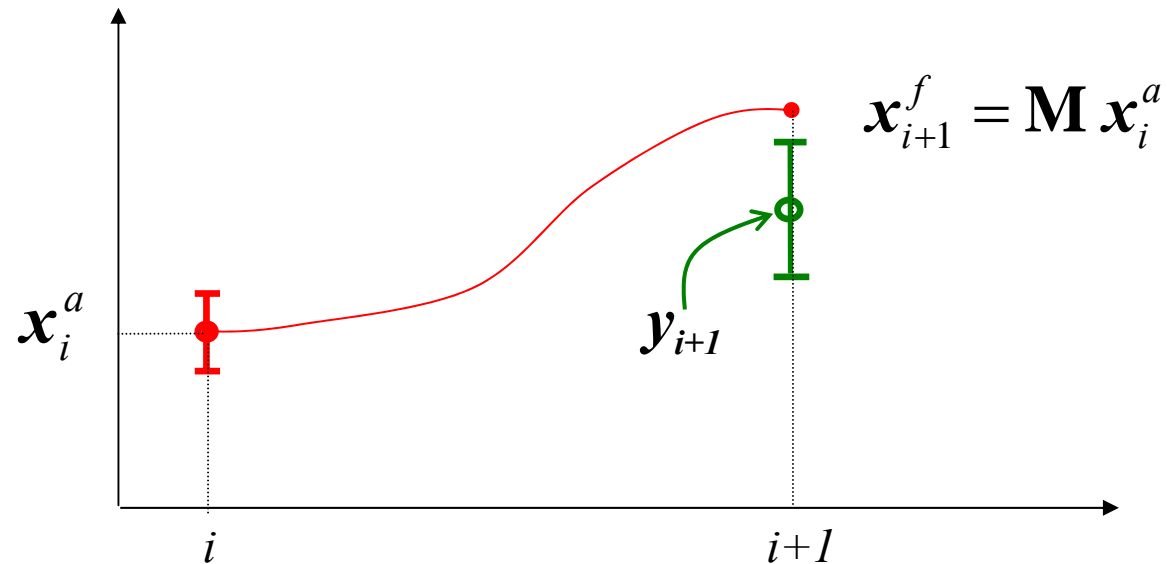
# « Sequential » data assimilation ?

❑   **Methods/algorithms**

➢     **« Variational »        vs.      « Sequential »** (Talagrand 1997, Ide et al., 1997)

➢     **« Smoothers »         vs.        « Filters »**

➢   **« Global problem »  vs.    « sub-problems »**   (Schröter 2004)

**Notations: Ide *et al.* (1997)**



$$x_{i+1}^{f} = \mathbf{M}\,x_{i}^{a}$$

$x_i^a$ : estimation of the « true » state vector $x_i^t$ at time $i$ , dimension $n$

$x_{i+1}^{f}$ : forecast of the state vector at time $i{+}1$, using the linear model $\mathbf{M}$

$y_{i+1}$ : observations available at time $i{+}1$ , dimension $p$

**How can the true state be best estimated
from this prior information ?**

Sea Surface Temperature on Gulf Stream
avhrr SST (August 26, 1993)

$$y_{i+1}$$

**incomplete data**

Gulf Stream temperature
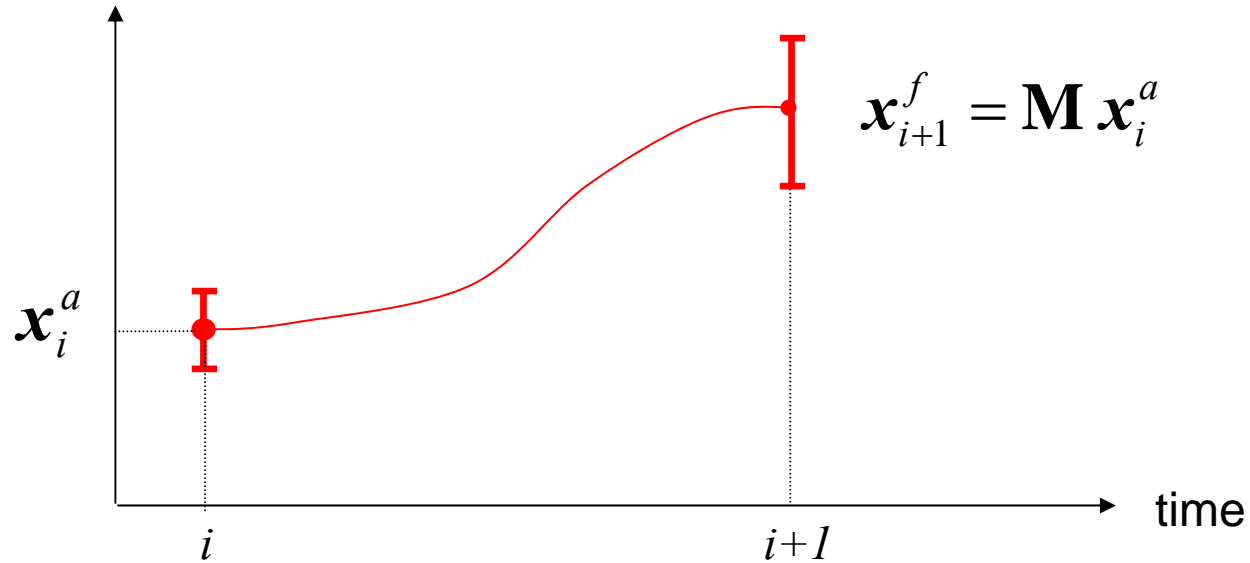
Isosurfaces:
V=0.7 m/s
T=10 degC

degC

LEGI/MEOM

$$x^{f}_{i+1} = \mathbf{M}\, x^{a}_{i}$$

**imperfect model**

$$\boldsymbol{\varepsilon}_i^a = \boldsymbol{x}_i^a - \boldsymbol{x}_i^t$$ : error on state estimate at time $i$ ; unknown quantity, but assume
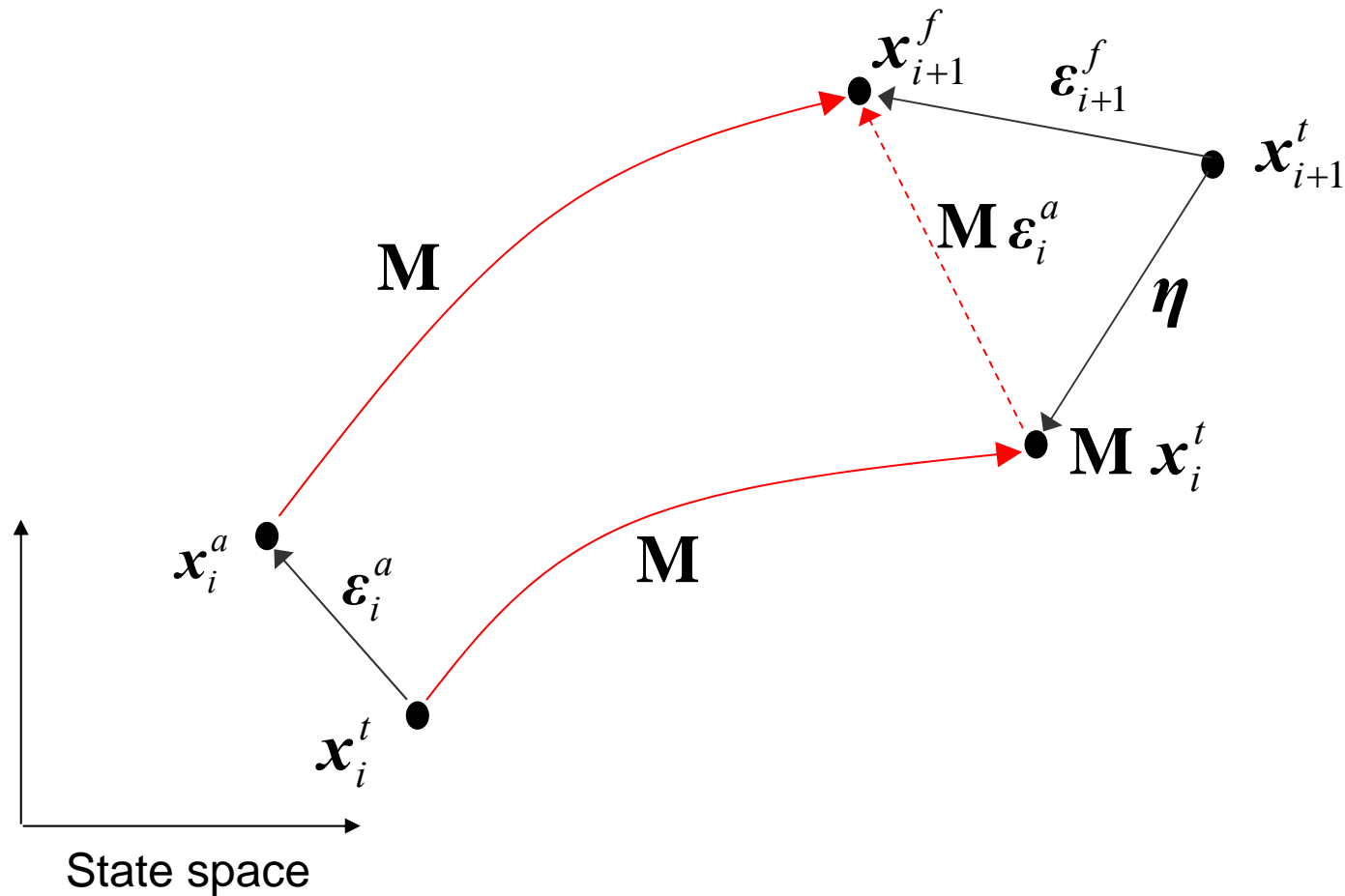
$$\boldsymbol{\varepsilon}_i^a \rightarrow N(0, \mathbf{P}_i^a) \sim \exp\left[-\frac{1}{2} \boldsymbol{\varepsilon}_i^{a^T} \mathbf{P}_i^{a^{-1}} \boldsymbol{\varepsilon}_i^a\right] \quad (1)$$

$$\boldsymbol{\eta} = \mathbf{M} \boldsymbol{x}_i^t - \boldsymbol{x}_{i+1}^t$$ : error on model forecast between time $i$ and time $i+1$ ; assume:

$$\boldsymbol{\eta} \rightarrow N(0, \mathbf{Q}) \sim \exp\left[-\frac{1}{2} \boldsymbol{\eta}^T \mathbf{Q}^{-1} \boldsymbol{\eta}\right] \quad (2)$$

$$\mathbf{M}\,\varepsilon_i^a = \mathbf{M}\,x_i^a - \mathbf{M}\,x_i^t = x_{i+1}^f - (x_{i+1}^t + \eta) = \varepsilon_{i+1}^f - \eta$$

Assuming pdf (1) and (2), model linearity and uncorrelated initial and modelling errors, the forecast error is distributed as :

$$\varepsilon_{i+1}^f \rightarrow N(0, \mathbf{P}_{i+1}^f) \sim \exp\left[ -\frac{1}{2}\varepsilon_{i+1}^{f\,T} \mathbf{P}_{i+1}^{f\,-1} \varepsilon_{i+1}^f \right] \qquad (3)$$

with

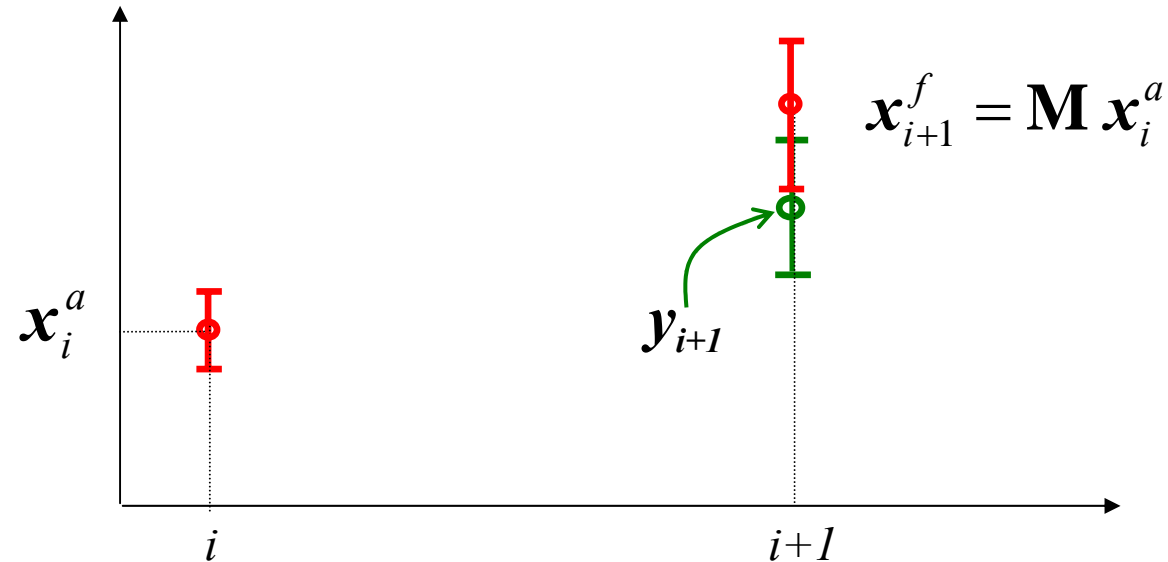$$\varepsilon_{i+1}^f = \mathbf{M}\,\varepsilon_i^a + \eta \quad \Longrightarrow \quad \overline{\varepsilon_{i+1}^f \varepsilon_{i+1}^{f\,T}} = \mathbf{M}\,\overline{\varepsilon_i^a \varepsilon_i^{a\,T}}\,\mathbf{M}^T + \overline{\eta\eta^T}$$

$$\Longrightarrow \quad \mathbf{P}_{i+1}^f \quad = \quad \mathbf{M}\mathbf{P}_i^a\mathbf{M}^T \quad + \mathbf{Q} \qquad (4)$$

**The estimation error is amplified by:**
**- unstable model dynamics (M) ;**
**- modelling errors Q .**

$$y_{i+1} = \mathbf{H}\, x_{i+1}^{t} + \varepsilon_{i+1}^{o} \quad : \text{observation available at time } i + 1$$

A probability distribution for the observation error is assumed:

$$\varepsilon_{i+1}^{o} \rightarrow N(0, \mathbf{R}) \sim \exp\left[ -\frac{1}{2}\, \varepsilon_{i+1}^{o}{}^{T}\, \mathbf{R}^{-1}\, \varepsilon_{i+1}^{o} \right] \qquad (5)$$

Using Bayes rule at time $i+1$ :

given by (5)    given by (3)

$$P\left(x_{i+1}^t \middle| y_{i+1}\right) = \frac{\overbrace{P\left(y_{i+1} \middle| x_{i+1}^t\right)} \cdot \overbrace{P\left(x_{i+1}^t\right)}}{\underbrace{P(y_{i+1})}} \quad (6)$$

a scaling factor

$$P\left(x_{i+1}^t\right) \cdot P\left(y_{i+1} \middle| x_{i+1}^t\right) \sim$$

$$\exp\left[-\frac{1}{2}(x_{i+1}^f - x_{i+1}^t)^T \mathbf{P}_{i+1}^{f\ -1}(x_{i+1}^f - x_{i+1}^t)\right] \cdot \exp\left[-\frac{1}{2}(y_{i+1} - \mathbf{H}\, x_{i+1}^t)^T \mathbf{R}^{-1}(y_{i+1} - \mathbf{H}\, x_{i+1}^t)\right]$$

$$= \exp\left[-\frac{1}{2}\left\{(x_{i+1}^f - x_{i+1}^t)^T \mathbf{P}_{i+1}^{f\ -1}(x_{i+1}^f - x_{i+1}^t) + (y_{i+1} - \mathbf{H}\, x_{i+1}^t)^T \mathbf{R}^{-1}(y_{i+1} - \mathbf{H}\, x_{i+1}^t)\right\}\right] \quad (7)$$

The best estimate of $x_{i+1}^t$ is the value of $x$ which maximize (7), i.e. the minimum of :

$$J(x) = (x_{i+1}^f - x)^T \mathbf{P}_{i+1}^{f\ -1}(x_{i+1}^f - x) + (y_{i+1} - \mathbf{H}\, x)^T \mathbf{R}^{-1}(y_{i+1} - \mathbf{H}\, x) \quad (8)$$

$$\delta_x J(x) = 0 \quad \Rightarrow \quad x = x_{i+1}^f + \mathbf{P}_{i+1}^f \mathbf{H}^T \mathbf{R}^{-1} (y_{i+1} - \mathbf{H} x) \quad (9)$$
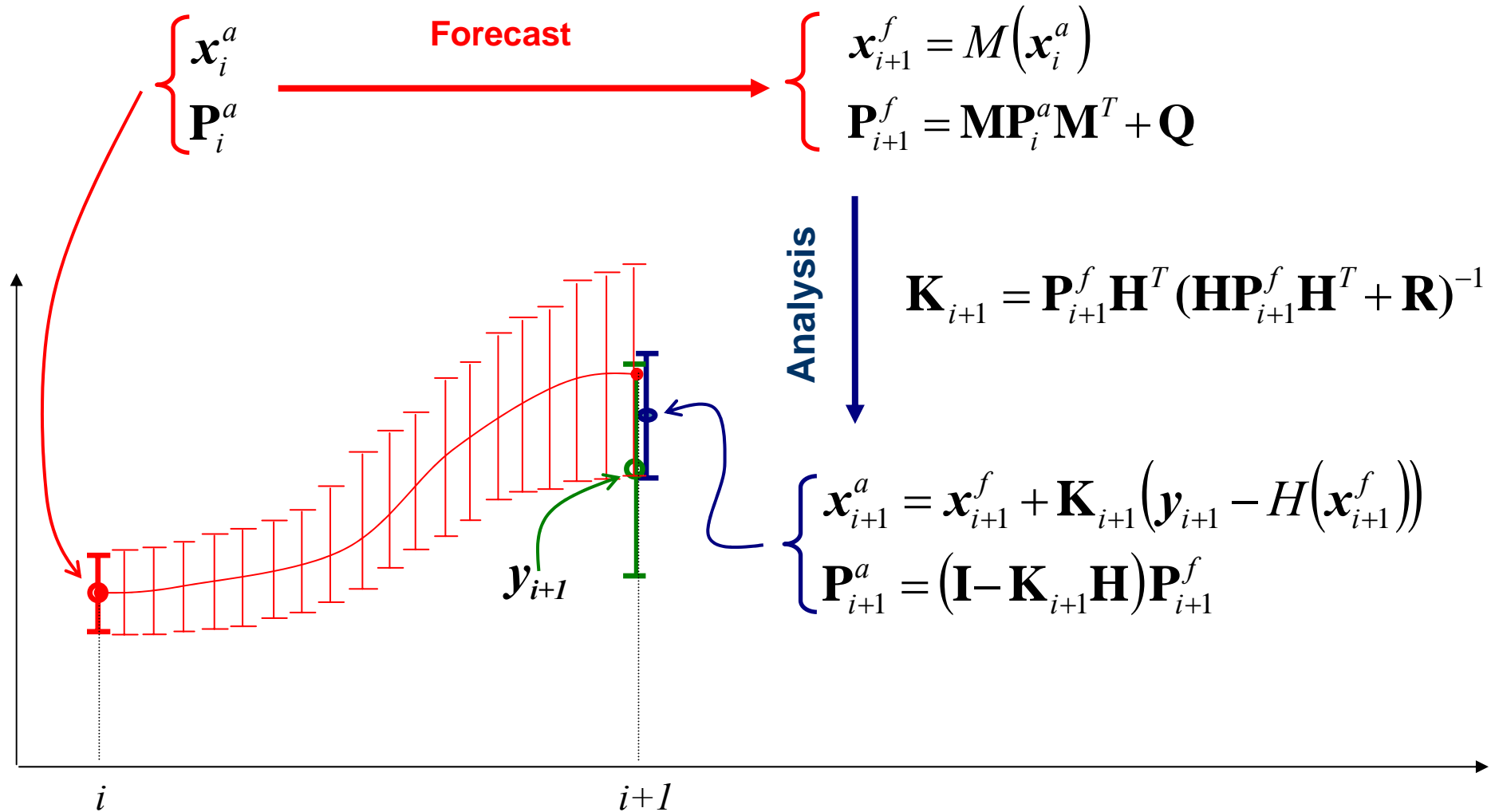
Equation (9) can be solved for $x$ using simple algebra (*), leading to:

$$x = x_{i+1}^f + \underbrace{\mathbf{P}_{i+1}^f \mathbf{H}^T (\mathbf{H}\mathbf{P}_{i+1}^f \mathbf{H}^T + \mathbf{R})^{-1}} (y_{i+1} - \mathbf{H} x_{i+1}^f)$$

Kalman gain = $\mathbf{K}_{i+1}$

Note: the forecast and analysis equations can be extended to *weakly* non-linear models $M$ and observation operator $H$.
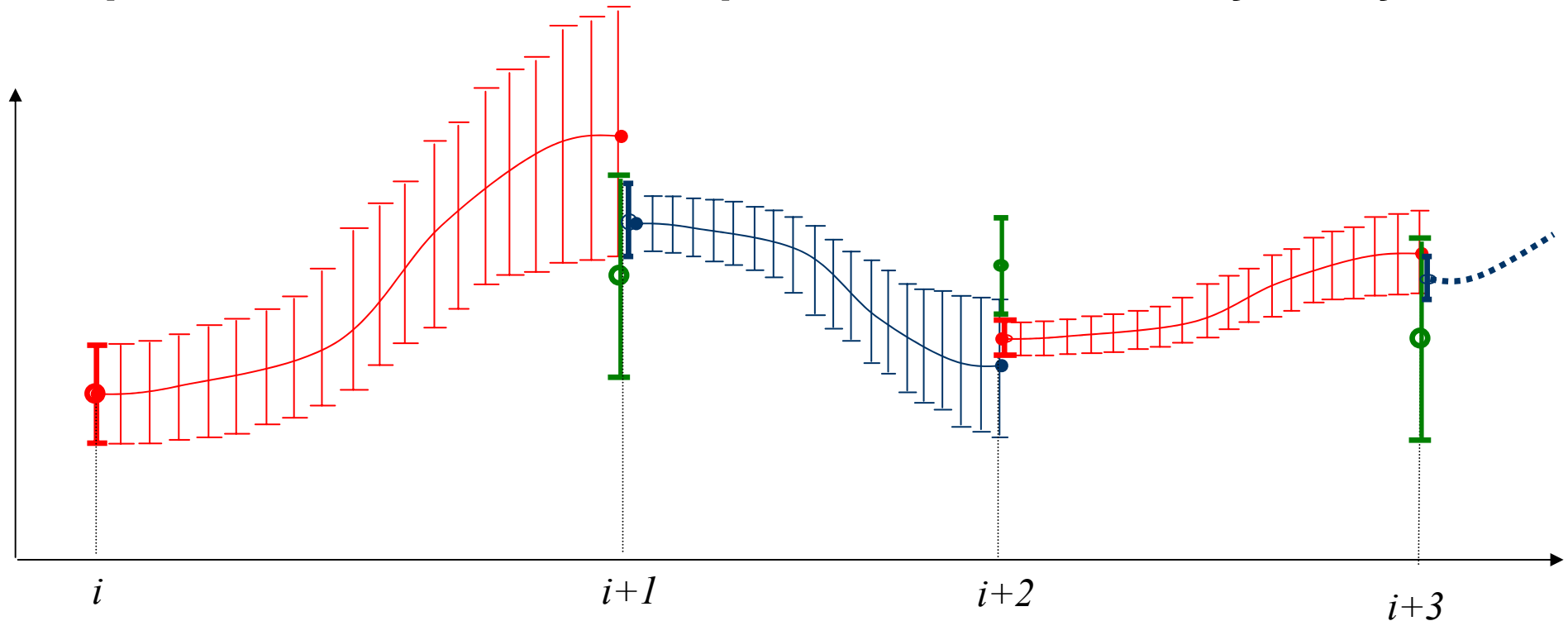
(*) Hint : use matrix equality $\left[ \mathbf{X}_1 + \mathbf{X}_{12} \, \mathbf{X}_2^{-1} \, \mathbf{X}_{21} \right]^{-1} = \mathbf{X}_1^{-1} - \mathbf{X}_1^{-1} \mathbf{X}_{12} \left[ \mathbf{X}_2 + \mathbf{X}_{21} \mathbf{X}_1^{-1} \mathbf{X}_{12} \right]^{-1} \mathbf{X}_{21} \mathbf{X}_1^{-1}$

**Forecast**

$$x_i^a$$
$$P_i^a$$

$$x_{i+1}^f = M\left(x_i^a\right)$$

$$P_{i+1}^f = MP_i^a M^T + Q$$

**Analysis**

$$K_{i+1} = P_{i+1}^f H^T (HP_{i+1}^f H^T + R)^{-1}$$

$$x_{i+1}^a = x_{i+1}^f + K_{i+1}\left(y_{i+1} - H\left(x_{i+1}^f\right)\right)$$

$$P_{i+1}^a = \left(I - K_{i+1}H\right)P_{i+1}^f$$

$$y_{i+1}$$

$i$

$i+1$

## Sequential assimilation = repeated forecast/analysis cycles



The best estimate at a given time is influenced by all previous observations (Kalman « filter »), and the analysis error covariance reflects the competition between this accumulation of past information and the error growth due to model imperfections .

❑ « **The scientific difficulty of data assimilation is to find algorithms which simplify the BLUE (Best Linear Unbiased Estimation) to an affordable amount of computer resources, while preserving some of the essential characteristics.** »

*Courtier*

*J. Meteor. Soc. Japan, 1997*

□ **Model variables in HYCOM(*)**

temperature, salinity, velocity, layer thicknesses, sea-surface height (SSH)



(*) ocean circulation model developed at Univ. Miami (RSMAS, E. Chassignet)

❑ **HYCOM state vector $x$** : 3D grid of the 5 scalar model variables
+ 2D grid for SSH

❑ **R&D prototype**: 1/3° horizontal resolution , 19 hybrid layers

$$n \sim 5 \times 350 \times 350 \times 19 \sim 1.1 \times 10^7$$

**Operational prototype**: 1/12° horizontal resolution , 26 hybrid layers

$$n \sim 5 \times 1400 \times 1400 \times 26 \sim 2.5 \times 10^8$$

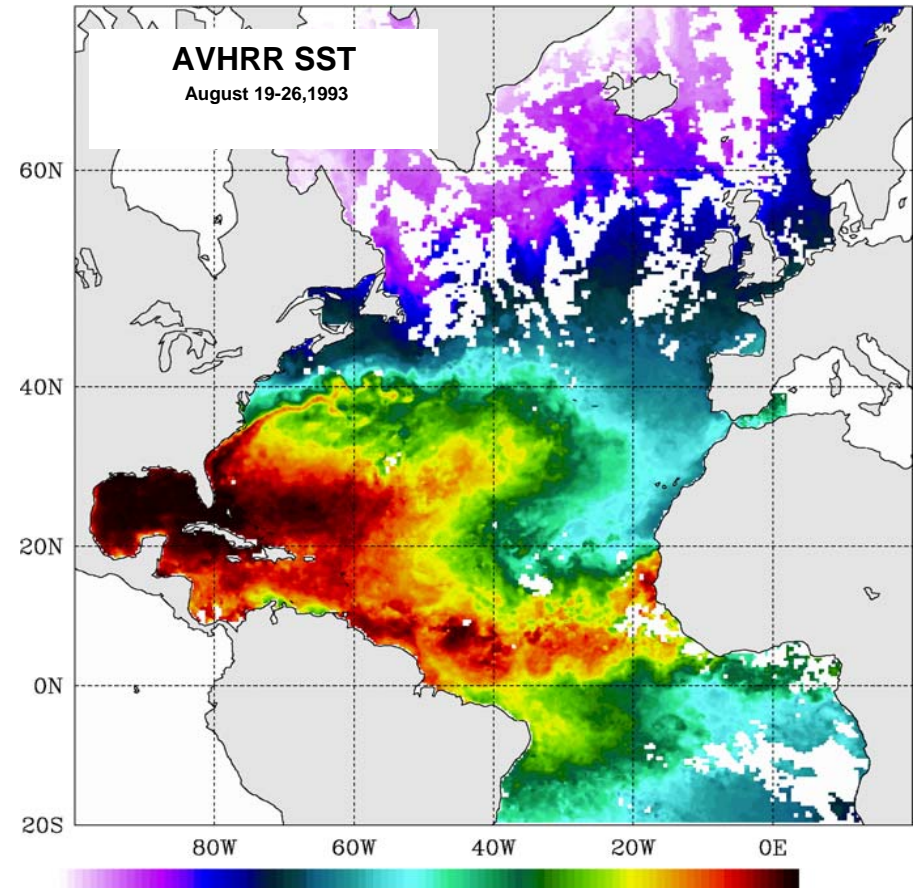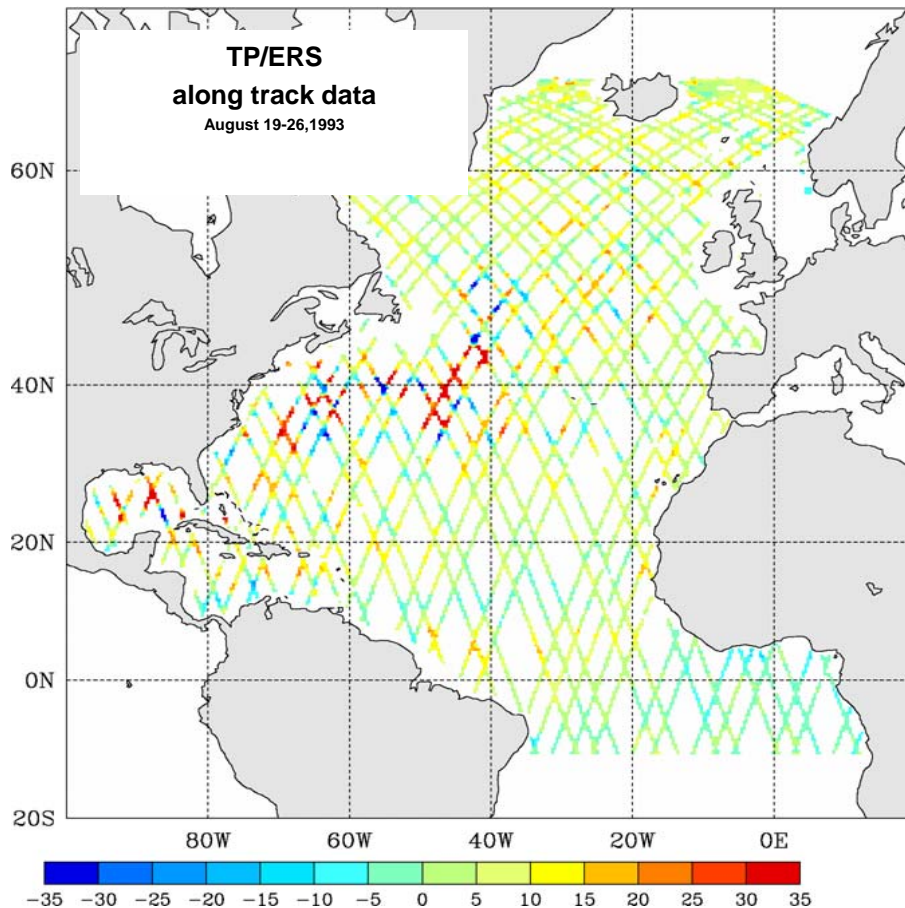❑ **M operator** : dim $n \times n \sim 6 \times 10^{16}$ real    (i.e. ~ 6000 Earth Simulators !)

➢The state vector dimension can be huge
➢The model transition matrix « **M** » cannot be represented explicitely.
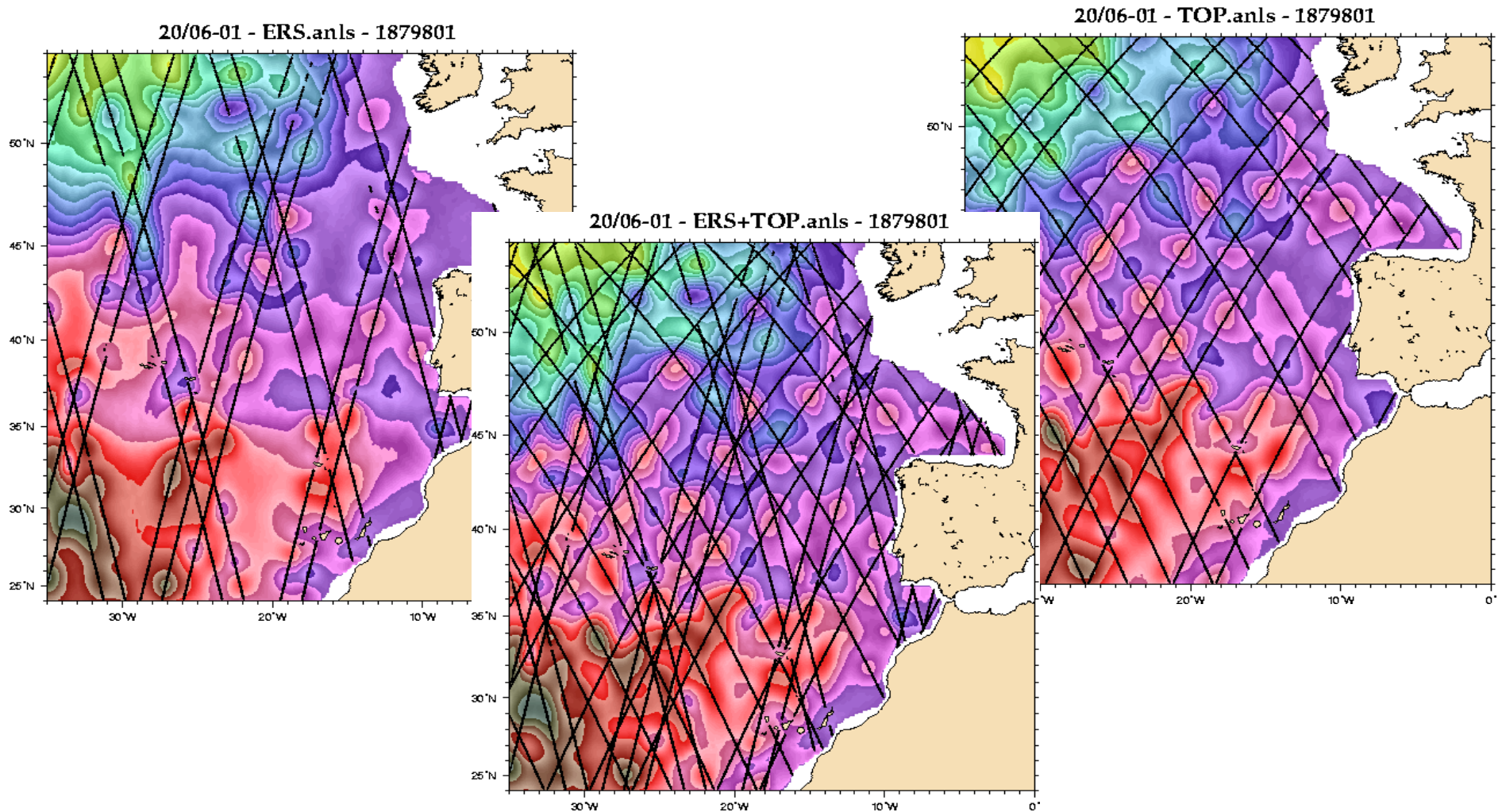➢ Instead, a computer code is used to transition  « $x$ » from time $i$ to $i+1$

❑ **Observed variables**:

- *from space*: sea-surface height (SSH), sea-surface temperature (SST)

❑ One week of altimeter data

❑ **Observed variables**:

- *in situ*: T/S profiles (drifting floats, field campaigns, …)



| | | |
|---|---|---|
| ● AUSTRALIA | ● GERMANY | ● NEW ZEALAND |
| ● CANADA | ● INDIA | ● NORWAY |
| ● CHINA | ● IRELAND | ● RUSSIAN FEDERATION |
| ● DENMARK | ● JAPAN | ● SPAIN |
| ● EUROPEAN UNION | ● KOREA (Rep. of) | ● UNITED KINGDOM |
| ● FRANCE | ● MAURITIUS | ● UNITED STATES |

**ARGO, July 2004**

❑ **<u>Observation vector *y*</u>** : data from various sources, at different time and space resolution

❑ **Radar Altimetry** : along-track measurements of SSH anomalies
(JASON: 1 obs. / 7 km, ~ 300 km equatorial tracks separation, repeated every 10 days ;
ERS/ENVISAT:  ~ 80 km equatorial tracks separation, repeated every 35 days)

❑ **AVHRR SST** : weekly composite images at 4 km resolution (if no clouds)

❑ **ARGO flots** : 3° x 3° horizontal resolution (targetted), profiles (between 2000 m depth to surface) every 10 days with 1 obs / m along vertical

➢ $p = \dim y$ is much smaller than $n = \dim x$ :  too few observations !

➢ The ocean surface relatively well observed by satellites: vertical extrapolation of data assimilated at the surface into the ocean's interior has to be consistent with vertical data profiles

➢ The observed variables are closely related to model variables:
$\mathbf{H}$ is mainly an interpolation operator  ( ~ simple)

❑ **Specification of error covariance matrix** $\mathbf{P}_0^a$ ?

❑ Assume a background state $\pmb{x}_0$ and associated error covariance $\mathbf{P}_0$
Consider the analysis step with only one data $\eta$ at a model grid point and the associated observation error $\varepsilon$.

➤ $p = 1$ , $\pmb{y}$ is a scalar and $\mathbf{H}$ is a vector of the form $\mathbf{H} = \begin{bmatrix} 0,...,0,1,0,...,0 \end{bmatrix}$

➤ The Kalman gain is then a ($n$ x $1$) vector:

$$\mathbf{K} = \mathbf{P}_0\mathbf{H}^T(\mathbf{H}\mathbf{P}_0\mathbf{H}^T + \mathbf{R})^{-1} = \frac{1}{(p_{\eta\eta} + \varepsilon^2)}\{\mathbf{P}_0\}_\eta \quad \text{with} \quad p_{\eta\eta} = \{\mathbf{P}_0\}_{\eta\eta}$$

➤ The posterior estimate is a correction of the background using the $\eta$–column of $\mathbf{P}_0$

$$\hat{\pmb{x}} = \pmb{x}_0 + \frac{1}{(p_{\eta\eta} + \varepsilon^2)}\{\mathbf{P}_0\}_\eta(\eta - \eta_0) \quad \text{with} \quad \eta_0 = \{\pmb{x}_0\}_\eta$$

Example: Horizontal covariance relative to a SSH ($\eta$) point at (32$^o$N,70$^o$W)
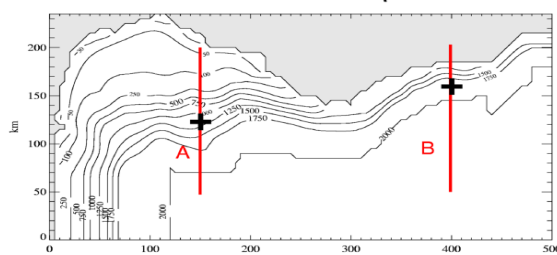MERCATOR Assimilation System - Testut *et al.*(2004)

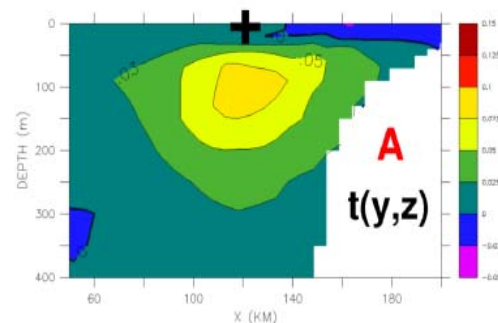

*Influence function of a single altimeter measurement in the sub-tropical gyre*

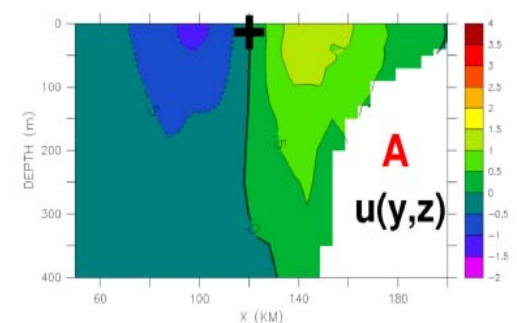- **Multivariate representers in 3D space, showing covariance structures consistent with the model dynamics**



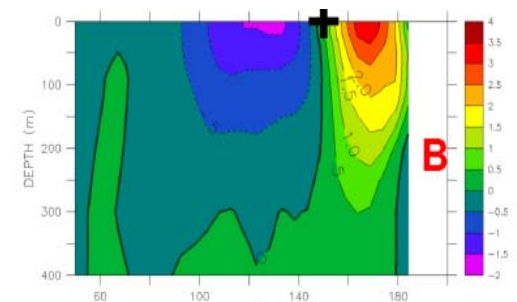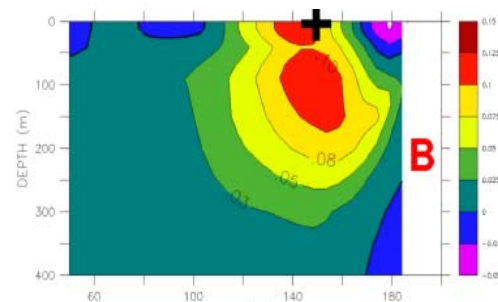Temperature and velocity corrections to a free surface elevation misfit δη at two locations

(+): location of sea surface measurement, δη=5 cm

**Representer functions of SSH in a free-surface coastal model**
(*Echevin et al., JPO, 2000*)

<u>Example</u>: covariance relative to a SSH ($\eta$) point at (0°,144°W)  -  Weaver *et al.*(2003)
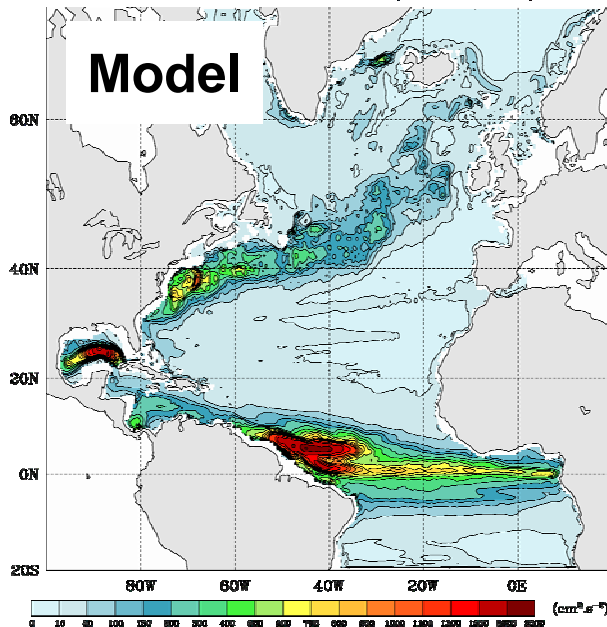


> The rows/colums of **P** should be « balanced » dynamically.

> This requires multivariate covariances

> A full-rank representation of **P** (dim $n$ x $n$) is still impossible !
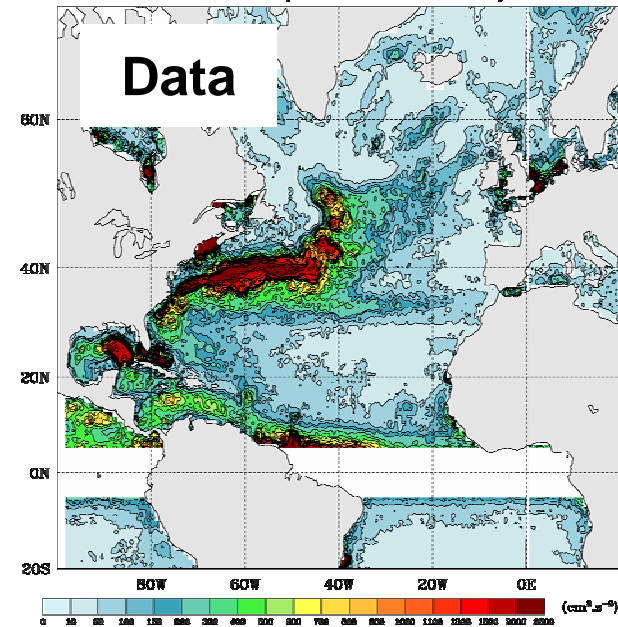
❑ **Model variability differs from observed variability**

**EKE example :**



Eddy Kinetic Energy at 50m depth on North Atlantic
NATL3 free simulation (1990-1999)

Model

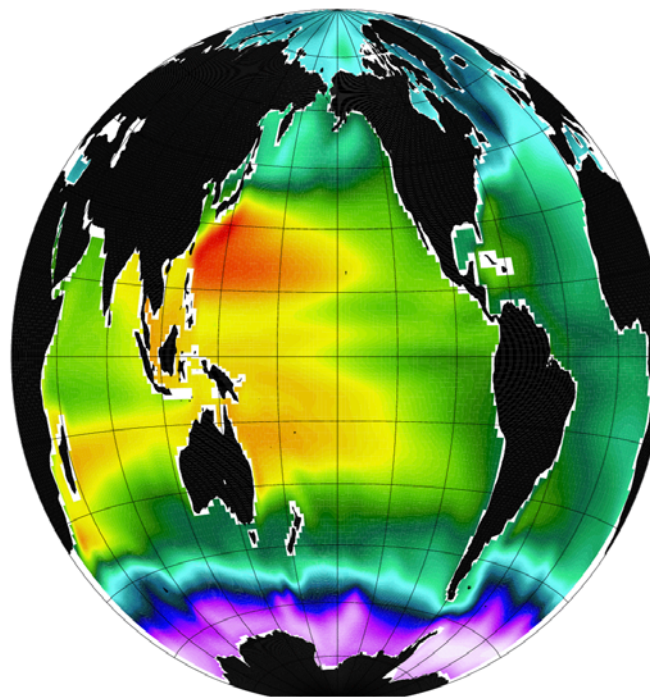Eddy Kinetic Energy on North Atlantic
TP/ERS (Oct 92 to Oct 97)

Data

➢ Many different model error sources (finite discretizations, representation of bottom topography, atmospheric forcings, etc … ) which cannot be easily quantified in terms of a $\mathbf{Q}$ matrix
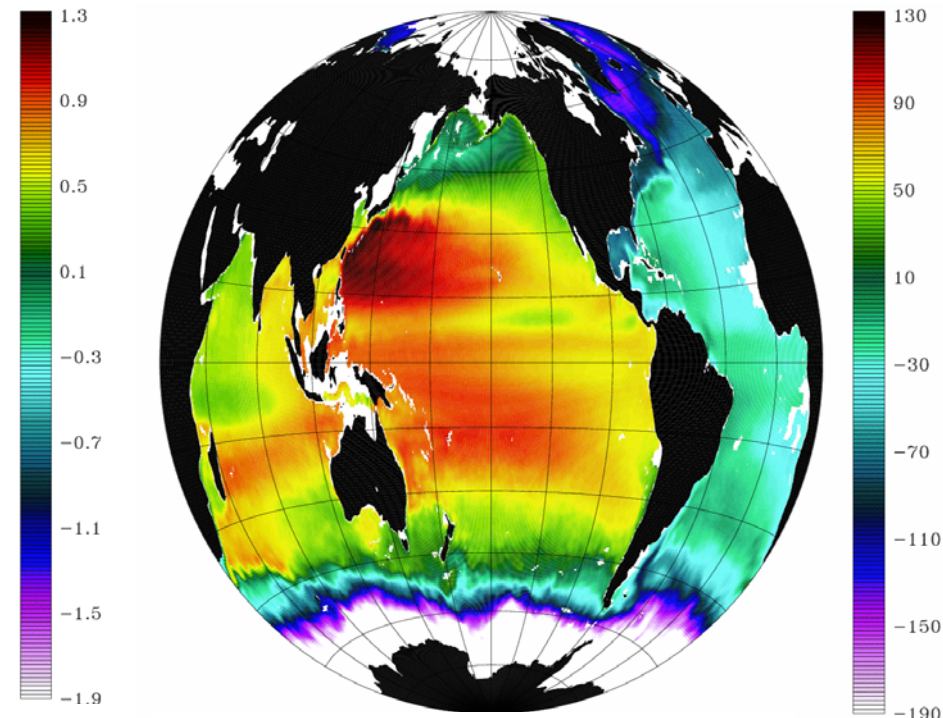
**Model mean SSH differs from observed mean SSH**

Mean sea-level difference between Pacific and Atlantic systematically too small in the model



**OPA model** (Madec et al. 2003)

**Data** (Niiler *et al.*, 2003)

Optimality properties of the KF only in the absence of biases !

❑ **« Optimal analysis is obtained when …the statistics …of error fields associated with forecasts and observations are known and accurately specified. Since these statistics are not generally available, actual implementations of assimilation algorithms are always sub-optimal»**

*Dee and Da Silva*

*Q. J. R. Meteorol. Soc., 1998*

**Error dynamics :**

❑ The forecast error requires ~ $n$ model integrations !!!

$$\mathbf{P}_{i+1}^f = \mathbf{M}\mathbf{P}_i^a\mathbf{M}^T + \mathbf{Q} = \mathbf{M}\left(\mathbf{M}\mathbf{P}_i^a\right)^T + \mathbf{Q}$$

❑ The ~ $n$ model integrations are useless if $\mathbf{Q}$ is poorly known

**Simplification of the Kalman filter**: « Optimal Interpolation »

To save cost and memory requirements, the KF can be simplified drastically by using time-independent « background » covariance matrix

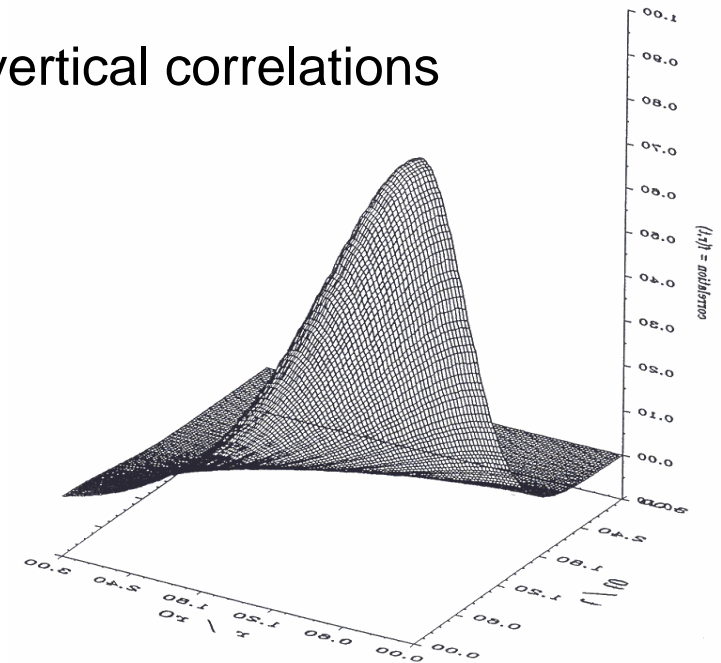$$\mathbf{P}_{i+1}^{f} = \mathbf{B} = \mathbf{D}^{1/2}\,\mathbf{C}\,\mathbf{D}^{1/2}$$

*correlations*

*variances*

❑ $\mathbf{C}$ : expressed as a product of horizontal and vertical correlations

$$\{\mathbf{C}\}_{i,j} = c^{h}(l)\,.\,c^{v}(d)$$

❑ Example: $c^{h}(l) = \left(1 + al + \dfrac{1}{3}a^{2}l^{2}\right)e^{-al}$

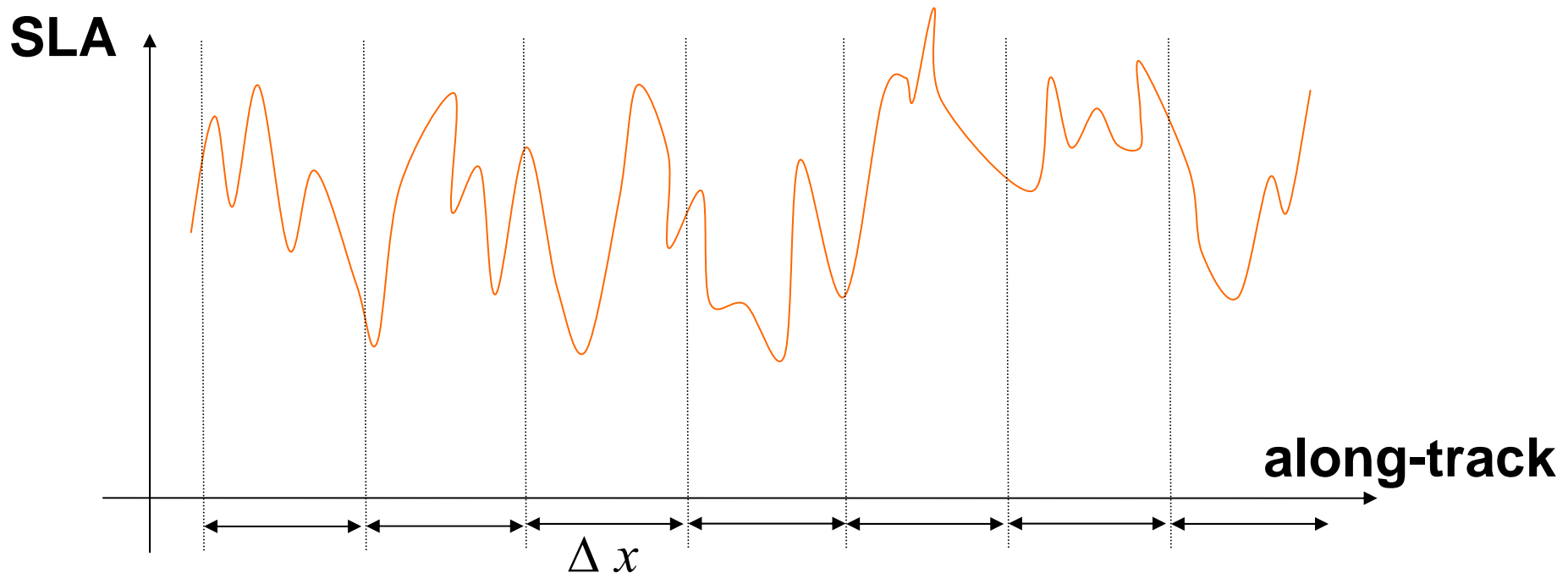❑ **Meridional and zonal correlation scales prescribed in SAM1**

**<u>Steady-state filters</u>** *(Fukumori et al., 1993)*

The KF can be simplified using a time-independent covariance matrix computed as the asymptotic limit of the **Riccati equation**:

$$\mathbf{P}_{i+1}^{f} = \mathbf{M} \underbrace{\left\{ \mathbf{P}_{i}^{f} - \mathbf{P}_{i}^{f} \mathbf{H}^{T} \left[ \mathbf{H}\mathbf{P}_{i}^{f}\mathbf{H}^{T} + \mathbf{R} \right]^{-1} \mathbf{H}\mathbf{P}_{i}^{f} \right\}}_{\mathbf{P}_{i}^{a}} \mathbf{M}^{T} + \mathbf{Q}$$

**The same model with the same set of observations can provide different sequences of « optimal estimates », depending on the « target » field !**
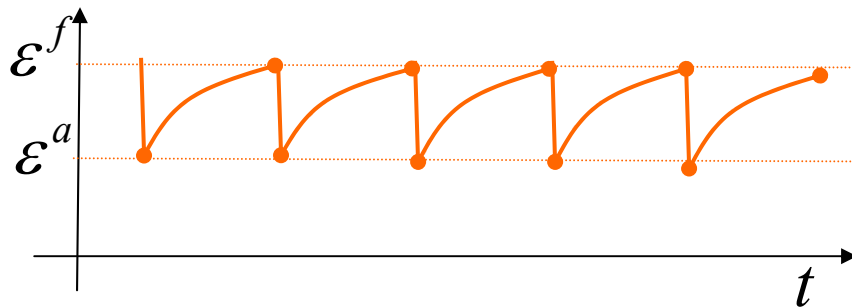
- ❏ Let $M$ : « perfect eddy-resolving » ocean model

    $y$ : « perfect filament-resolving » ocean observations



SLA

along-track

$\Delta x$

**Target = eddies** (~ 50 km)

$$\mathbf{P}_{i+1}^{f} = \mathbf{MP}_{i}^{a}\mathbf{M}^{T} + 0$$

$$\mathbf{K}_{i+1} = \mathbf{P}_{i+1}^{f}\,\mathbf{H}^{T}(\mathbf{HP}_{i+1}^{f}\mathbf{H}^{T} + \mathbf{R}^{\varphi})^{-1}$$



$\mathbf{R}^{\varphi}$ : representativeness error due to the observation of filaments

- Target = model resolution

- KF equations

$$p_{i+1}^f = mp_i^a m + 0 = m^2 p_i^a$$

$$k_{i+1} = p_{i+1}^f (p_{i+1}^f + r^\varphi)^{-1} = m^2 p_i^a (m^2 p_i^a + r^\varphi)^{-1} \qquad (0 < k_{i+1} < 1)$$

$$p_{i+1}^a = (1 - k_{i+1}) p_{i+1}^f = r^\varphi m^2 p_i^a (m^2 p_i^a + r^\varphi)^{-1} \qquad (0 < p_{i+1}^a < r^\varphi)$$

❑ Target = observation resolution

❑ KF equations

$$p_{i+1}^{f} = mp_{i}^{a}m + q^{\varphi} = m^{2}p_{i}^{a} + q^{\varphi}$$

$$k_{i+1} = p_{i+1}^{f}(p_{i+1}^{f} + 0)^{-1} = 1 \quad \Rightarrow \quad x_{i+1}^{a} = x_{i+1}^{f} + k_{i+1}(y_{i+1} - x_{i+1}^{f}) = y_{i+1} \;!$$

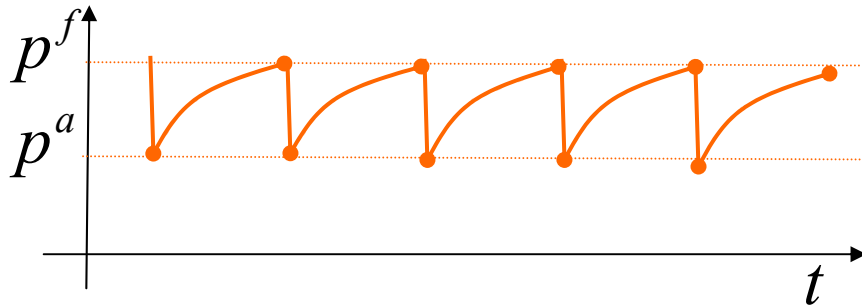$$p_{i+1}^{a} = (1 - k_{i+1})p_{i+1}^{f} = 0$$

**Example : scalar, linear model** $x_{i+1}^f = mx_i^a$

❑ Numerical example: $m = \sqrt{2}$ $q^\varphi = r^\varphi = \varepsilon$

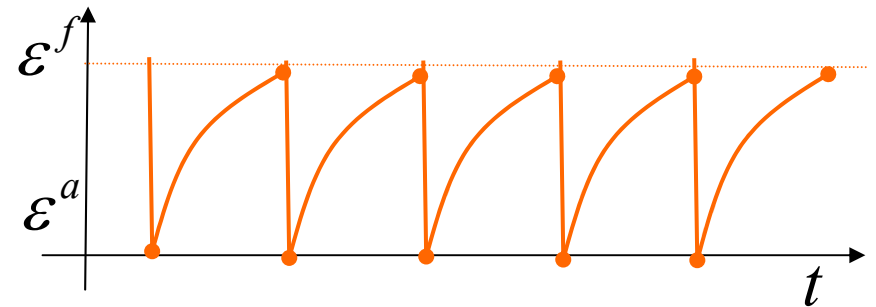**Target = model resolution**

$$p^a = 0.5\,\varepsilon$$

$$p^f = \varepsilon$$

**Target = data resolution**

$$p^a = 0$$

$$p^f = \varepsilon$$

❑ Misleading implementation : $\quad q^{\varphi} = r^{\varphi} = 0$

❑ Initialization : $\quad x_0, p_0$

- A full Kalman filter cannot be implemented into realistic ocean models (error forecast and analysis equations too expensive in CPU and memory requirements)

- « Optimal Interpolation » over-simplifies the propagation of errors by neglecting dynamical principles and statistical information

**Idealized double-gyre model** (Ballabrera *et al., 2001*)