# Data Assimilation: Methods, Concepts and Challenges

Alberto Carrassi

acarrassi@ic3.cat

Climate Forecasting Unit - **CFU**

Catalan Institute for Climate Science - **IC3**

Spain

Exploratory Workshop DADA - 15 October 2012

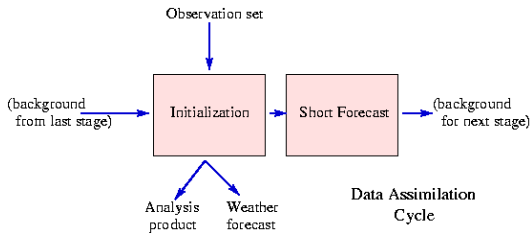Exploring the Use of Data Assimilation for the Detection and Attribution of Climate Change

Data Assimilation is the entire sequence of operations that, starting from the observations and possibly from a statistical/dynamical knowledge about a system, provides an estimate of its state

The main fields of applications in geophysics are:

- initialize weather prediction
- produce reanalysis
- parameter estimation

# Basic Definitions and Problem Statement

OBJECTIVE:

estimate the state of an unknown system based on an imperfect model and a limited set of noisy observations:

$$\mathbf{x}_k = \mathcal{M}_k(\mathbf{x}_{k-1}) + \mu_k \quad k = 1, 2, ...,$$

$$\mathbf{y}_k^o = \mathcal{H}(\mathbf{x}_k) + \varepsilon_k^o \quad k = 1, 2, ...,$$

- $\mathbf{y}^o \in \mathcal{R}^p$ and $\mathbf{x} \in \mathcal{R}^n$ - $p \ll n$ in realistic geophysical applications
- $\{\mu_k\}_{k=1,2...}$ and $\{\varepsilon_k^o\}_{k=1,2...}$ *assumed* to be random error sequences, white in time, and uncorrelated between them
- Collect state estimates and observations as: $\mathbf{X}_k = \{\mathbf{x}_0, \mathbf{x}_1, ..., \mathbf{x}_k\}$ and $\mathbf{Y}_k^0 = \{\mathbf{y}_0^0, \mathbf{y}_1^0, ..., \mathbf{y}_k^0\}$

Smoothing, Filtering or Prediction ?

1. Smoothing $\rightarrow$ Estimate the state at all times $\equiv \mathbf{X}_k$ based on $\mathbf{Y}_k^0$

2. Filtering $\rightarrow$ Estimate the state at the present time $\equiv \mathbf{x}_k$ based on $\mathbf{Y}_{k-1}^0$

3. Prediction $\rightarrow$ Estimate the state at future times $\equiv \mathbf{x}_{k>l}$ based on $\mathbf{Y}_l^0$

# Probabilistic Approach

In the probabilistic framework, problems (1)-(2)-(3) are expressed as the estimation of the corresponding conditional probability density functions:

1. Smoothing → Estimate $\mathcal{P}(\mathbf{X}_k|\mathbf{Y}_k^0)$

2. Filtering → Estimate $\mathcal{P}(\mathbf{x}_k|\mathbf{Y}_{k-1}^0)$

3. Prediction → Estimate $\mathcal{P}(\mathbf{x}_{k>l}|\mathbf{Y}_l^0)$

The PDFs $\mathcal{P}$ fully characterise the estimation problem!

The error PDFs associated to all the information sources read:

- $\mathcal{P}(\mathbf{x}_0)$ PDF of the initial conditions - Prior

- $\mathcal{P}(\mu_k) = \mathcal{P}(\mathbf{x}_k - \mathcal{M}_k(\mathbf{x}_{k-1})) = \mathcal{P}(\mathbf{x}_k|\mathbf{x}_{k-1})$ - Model Error PDF

- $\mathcal{P}(\varepsilon_k^o) = \mathcal{P}(\mathbf{y}_k^0 - \mathcal{H}(\mathbf{x}_k)) = \mathcal{P}(\mathbf{y}_k|\mathbf{x}_k)$ - Observational Error PDF
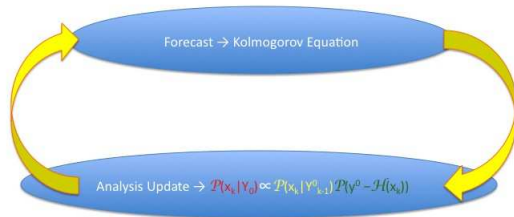
# Probabilistic Approach

With Bayes's rules....
SMOOTHING

$$\mathcal{P}(\mathbf{X}_k|\mathbf{Y}_k^0) \propto \mathcal{P}(\mathbf{x}_0)\Pi_{i=1}^{k}\mathcal{P}(\mathbf{x}_i - \mathcal{M}_i(\mathbf{x}_{i-1}))\mathcal{P}(\mathbf{y}_i^0 - \mathcal{H}(\mathbf{x}_i))$$
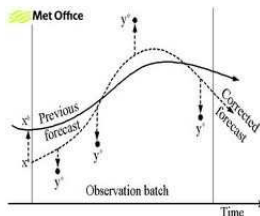
FILTERING



In high-dimensional nonlinear systems the full Bayesian formulation is not affordable

Note: The Particle Filters attempt to solve this problem and their potential application in geoscience has received much attention in recent years. See van Leuween, 2009 (MWR) for a review and the Philippe Naveau's talk in this workshop
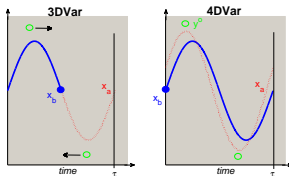
# 4D-Variational Assimilation



Initial condition, observational and model errors are all Gaussian and mutually uncorrelated $\Longrightarrow$ solving the SMOOTHING problem leads to the 4DVar formulation, *i.e.* minimise a penalty function as:

$$2J = \sum_{i=1}^{k} \mu_i^T \mathbf{Q}_i^{-1} \mu_i + \sum_{i=1}^{k} [\mathbf{y}_i^0 - \mathcal{H}(\mathbf{x}_i)]^T \mathbf{R}_i^{-1} [\mathbf{y}_i^0 - \mathcal{H}(\mathbf{x}_i)] + (\mathbf{x_0} - \mathbf{x_b})^T \mathbf{B}^{-1} (\mathbf{x_0} - \mathbf{x_b})$$

- **B** - Background error covariance matrix
- **R** - Observational error covariance matrix
- **Q** - Model error covariance matrix

# 4D-Variational Assimilation

- The sequence (trajectory) $\mathbf{X}_k$ which minimizes $J$ is the maximum likelihood estimator of the PDF $\mathcal{P}(\mathbf{X}_k|\mathbf{Y}_k^0)$

- It provides the *"best"* possible fit to the observations, given the initial guess and the *imperfect* model

- The strong-constraint 4DVar makes the assumption of perfect model and the latter is appended as a strong-constraint when doing the minimization

- The minimization of $J$ can be done in principle by solving the associated Euler-Lagrange (EL) equations (Le Dimet and Talagrand, 1986 Tellus)

- The *Method of Representer* is an efficient way to solve the EL eqs for linear dynamics (Bennett, 1982, chapter 5)

- *Descent Methods* are used in the case of large nonlinear systems (Talagrand and Courtier, 1987 QJRMS)

- The choice of the *Control Variable* defines the size of the problem to be solved and characterises different formulations of the 4DVar (see *e.g.* Tremolet, 2006 QJRMS; Bocquet, 2009 MWR)

- **B is implicitly evolved** within the assimilation window but it is not available for the next analysis cycle

- When observations are assimilated (as they were) at the same time the 3DVar is recovered

- 4DVar (under "strong" simplified assumptions) is operational in several weather services, among them MetOffice and ECMWF.

# Sequential Assimilation

Under the same hypotheses of Gaussianity and mutual uncorrelation of errors the filtering problem reduces to the estimation of the mean and covariance.

## ANALYSIS UPDATE EQUATIONS

$$\mathbf{x}_k^a = \mathbf{x}_k^f + \mathbf{K}_k \left[ \mathbf{y}_k^o - \mathcal{H}_k(\mathbf{x}_k^f) \right]$$

$$\mathbf{P}_k^a = [\mathbf{I} - \mathbf{K}_k \mathbf{H}_k] \mathbf{P}_k^f$$

$$\mathbf{K}_k = \mathbf{P}_k^f \mathbf{H}_k^T (\mathbf{H}_k \mathbf{P}_k^f \mathbf{H}_k^T + \mathbf{R}_k)^{-1}$$

- $\mathbf{x}_k^a$ - Analysis state at time $t_k$
- $\mathbf{x}_k^f = \mathcal{M}(\mathbf{x}_{k-1}^f)$ - Forecast state at time $t_k$
- $\mathbf{P}^f$ - Forecast error covariance matrix
- $\mathbf{R}$ - Observational error covariance matrix
- $\mathbf{K}$ - Kalman gain matrix
- The analysis $\mathbf{x}^a$ is optimal in the sense that it minimizes the analysis error variance
- When all errors are Gaussian the minimum variance estimate is also the maximum likelihood estimate (out of unimodality maximum likelihood estimators are of questionable relevance)
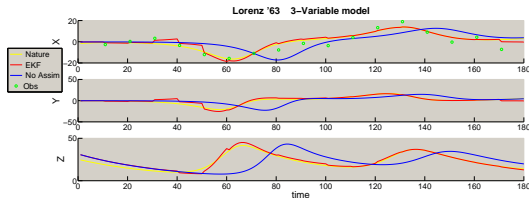
# Kalman Filter (KF) and Extended KF

For linear dynamics and observational operator the KF provides a closed set of estimation equations (Kalman, 1960). The forecast step equations read:

$$\mathbf{x}_k^f = \mathbf{M}\mathbf{x}_{k-1}^f + \mu_k$$

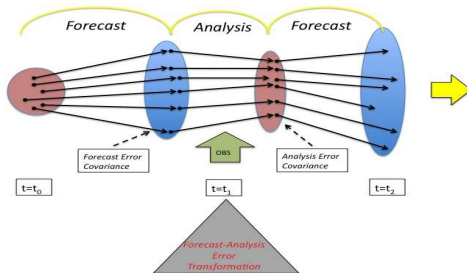$$\mathbf{P}_k^f = \mathbf{M}_k \mathbf{P}_{k-1}^a \mathbf{M}_k^T + \mathbf{Q}_k$$

Extension to nonlinear dynamics - Extended Kalman Filter

- The extended Kalman Filter (EKF) is a first order approximation of the KF
- The tangent linear model is used to forward propagate the forecast uncertainty (*i.e.* the error covariance)
- The full nonlinear model is used to evolve the state estimate
- The analysis update is the same as in the standard KF
- The introduction of the EKF in geoscience is due to Ghil and Malanotte-Rizzoli (1991) *AdvGeohys*
- The EKF response to different degree of nonlinearity has been studied in Miller, Ghil & Gauthiez (1994) *JAS*
- The EKF is almost-operational for ECMWF soil analysis (de Rosnay *et al.*, 2012 QJRMS)



Lorenz '63   3–Variable model

# Ensemble Based Data Assimilation Algorithms

In the ensemble-based DA the forecast/analysis error covariances are approximated using an ensemble of $M$ model trajectories



- Ensemble based covariances $\mathbf{P}^{f,a} = \frac{1}{M-1} \sum_{i=1}^{M} (\mathbf{x}_i^{f,a} - \bar{\mathbf{x}}^{f,a})(\mathbf{x}_i^{f,a} - \bar{\mathbf{x}}^{f,a})^T$
- For the approach to be suitable in geoscience $M \ll n$
- Flow dependent description of the forecast error
- The Kalman gain is computed in the observation space reducing the computational cost at a rate given by the ratio between the number of observations and the system size, $p/n$, $\mathbf{P}^f \mathbf{H}^T = \frac{1}{M-1} \sum_{i=1}^{M} (\delta \mathbf{x}_i^f \mathbf{H} \delta \mathbf{x}_i^f) \ldots$
- Provide automatically a set of initial conditions for ensemble prediction schemes.
- The choice of the forecast-analysis transformation characterises the ensemble-based algorithms.

# Stochastic or Deterministic ?

Ensemble data assimilation algorithms can be divided into Stochastic and Deterministic
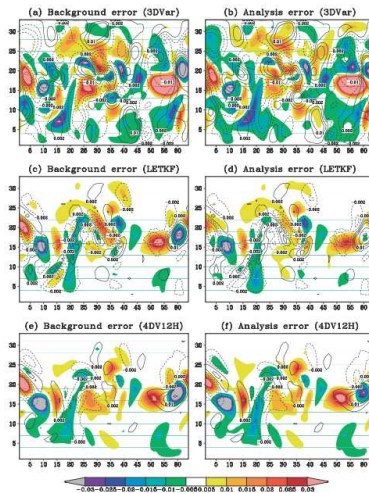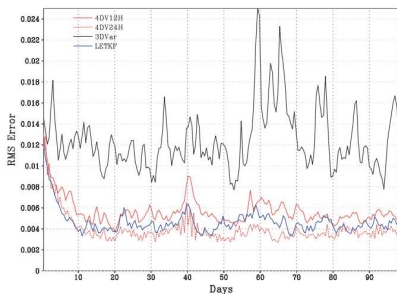
## Stochastic (Monte-Carlo approach)

- In this class of algorithms the observations are treated as a random ensemble by adding noise at each analysis update
- Each ensemble trajectory assimilates a different realization of the observation vector and undergoes an independent analysis update
- The standard Ensemble Kalman Filter (EnKF) belongs to this family (see *e.g.* Houtekamer and Mitchell, 1998 MWR)
- The EnKF has proved efficiency in a number of geophysical applications (see Evensen, 2003 Ocean Dyn for a review)

## Deterministic (Square-Root approach)

- In this class of algorithms the step $P^f \rightarrow P^a$ is made through a linear transformation $T$
- It avoids the introduction of extra noise at the analysis update
- $T$ is usually defined under the constraint that $P^a$ matches some desired value (*i.e.* the EKF one, the Hessian of a penalty function)
- The solution (a square-root matrix) is not unique and the particular choice characterises the algorithm (see Tippet *et al.*, 2003 MWR).
- Algorithms belonging to this family: ETKF, LETKF, EnSRF, MLEF (see Whitaker and Hamill, 2002 MWR; Bishop *et al.*, 2001 MWR; Hunt *et al.*, 2007 Physica D)

# Ensemble-based or Variational: the comparison

- Results with a Quasi-Geostrophic model by Rotunno and Bao, 1996
- Ensemble-based scheme $\Longrightarrow$ Local Ensemble Transform Kalman Filter (Hunt *et al.*, 2007 Physica D)



From Yang,Corazza,Carrassi,Kalnay & Miyoshi, 2009 MWR

# Dealing with Geophysical Systems

When dealing with realistic Atmosphere/Ocean dynamics DA faces a number of obstacles....

- Huge dimension $\Rightarrow$ Computationally suitable solutions ...

- The Atmosphere and the Ocean are example of nonlinear chaotic systems

- Chaos implies (among other things !) high flow dependent variability of error dynamics $\Rightarrow$ Flow-dependent description of the error entering the estimation is required

- Sources of nonlinearities: model $\mathcal{M}$, obs operator $\mathcal{H}$, first guess $\mathbf{B}$

- Nonlinearities push out of Gaussianity $\Rightarrow$ Non-Gaussian analysis frameworks (for a complete review see Bocquet, Pires & Wu, 2010 MWR)

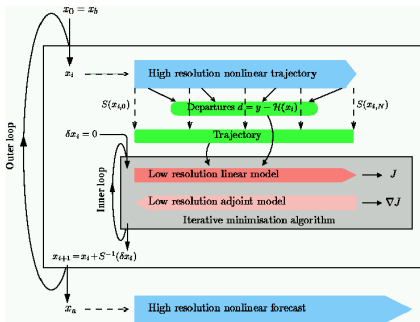Solutions have been proposed in the framework of both Variational and Ensemble schemes ....

# How to deal with geophysical systems: Variational

Main Drawbacks of Variational Approach:

1. Non-Quadratic cost-function in 4DVar
2. with possible Multiple Minima
3. maximum likelihood approach questionable
4. No flow-dependent error description

Proposed Solutions:

- Problem (1) and (2) are alleviated in the Incremental 4DVar (Courtier et al., 1994 QJRMS).
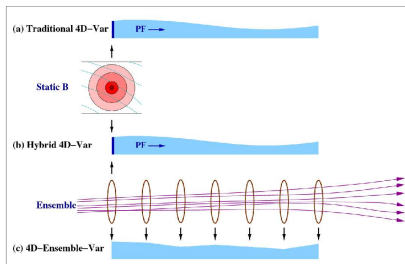


From Andersson et al., 2005 ECMWF-Tech.Rep. 479

# How to deal with geophysical systems: Variational

- Problem (4) is implicitly overcome with the Long Window 4DVar but ... problems (1)-(3) can be made worst
- Problems (1)-(3) are partly solved using the Weak-Constraint 4DVar but ... appropriate model error covariances need to be prescribed and the size of the control variable too big (see *e.g.* Trémolet, 2006 QJRMS)
- Hybrid 3/4DVar-Ensemble algorithms attempt to tackle all problems at the same time (see Barker and Clayton, 2011 ECMWF Ann. Seminar for a review and for details on the operational implementation at MetOffice).

Example: ETKF ↔ 4DVar at MetOffice (from Barker and Clayton, 2011 ECMWF Ann. Seminar)

Two hybrid strategies:

- Hybrid 4DVar operational at MetOffice (Use a combination of static and ensemble cov at the initial time)
- 4D-Ensemble-Var mid-long term development (Use ensemble cov within the entire assimilation window ⇒ No need for Tangent/Adjoint model) See Buehner *et al.*, 2010 MWR



From Barker and Clayton, 2011 ECMWF Ann. Seminar

# How to deal with geophysical systems: Ensemble Schemes
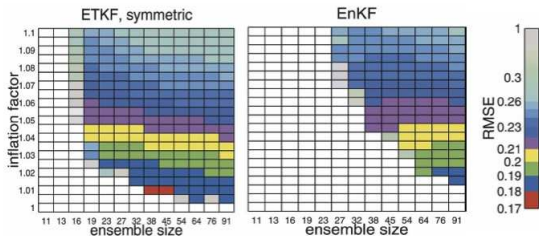
Main Drawbacks of Ensemble Based Approach:

1. Sampling Error ($M\ O(100)$)
2. Use only observations at analysis time
3. Only the Gaussian approximation of the flow-dependent $\mathbf{P}^f$ is accounted for at the analysis update

Proposed Solutions:

- Sampling errors (problem (1)) are mitigated using Covariance Localization $\Rightarrow$ Effective increase the rank of $\mathbf{P}^f$; but:
  - dynamical consistency is broken
  - the actual optimal size for the localization is time-dependent $\Rightarrow$ Flow-Dependent Covariance Localization (Bishop and Hodyss, 2011)

- Variance Underestimation (still problem (1)) $\Rightarrow$ Multiplicative or Additive Inflation
  Multiplicative Inflation (See e.g. Anderson and Anderson, 1999 MWR):
  - $\mathbf{P}^f \rightarrow (1 + \alpha)\mathbf{P}^f$
  - keep the same rank/structure of $\mathbf{P}^f$, only the explained variance is modified
  - the inflation can be made adaptive $\Leftrightarrow$ more inflation where/when required: based on Kalman gain (Sacher and Bartello, 2008 MWR), on analysis error variance (Whitaker and Hamill, 2012 MWR)
  Additive Inflation:
  - add random noise to $\mathbf{P}^f$ or $\mathbf{P}^a$
  - the process introduce new structures in the error space spanned by the ensemble covariances
  - a combined additive/multiplicative scheme has been proposed by Zhang et al., 2004 MWR
  - an ensemble based algorithm without the need of inflation has been proposed recently (Bocquet, 2011 NPG)

- An Hybrid approach is used to deal with problem (2) $\Rightarrow$ Several ensemble schemes introduce the time dimension to assimilate observations simultaneously over a given reference period (see e.g. Hunt et al., 2004 Tellus; Zhang and Zhang, 2012 MWR)

- Solution to problem (3) $\Rightarrow$ Particle Filters but ....

# exploiting chaos... → optimal ensemble size

Can the ensemble size be designed based on the system dynamical properties ?



Adapted from Sakov and Oke, 2008 MWR

- Lorenz 1996 Model
- Main finding: The ETKF converges to low error level when $N_{ens} \geq KY_{dim}$ reaches the model subspace dimension
- This behavior is deeply different from the EnKF whose performance improves indefinitely when $Ens_{size} \to \infty$
- With another deterministic filter (MLEF, Zupanski, 2005 MWR), without inflation, Carrassi et al., 2009 (Tellus) found a similar behavior (error saturation when $N_{ens} \geq KY_{dim}$)
- Bocquet, 2011 (NPG) introduced a new deterministic filter (ETKF-N) that does not need inflation as long as $Ens_{size} > N^+$
- In deterministic filters ensemble perturbations reflect the intrinsic system error dynamics and have to be intended as factorization of the system's error covariance rather than its Monte Carlo approximation as in the EnKF
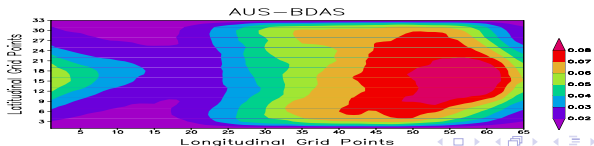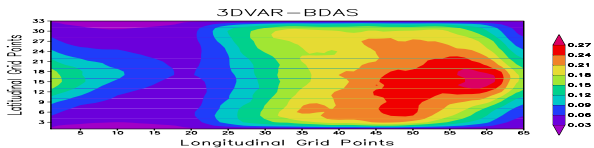
# exploiting chaos...→ Target Observations

## Assimilation in the Unstable Subspace (Trevisan and Uboldi, 2004 JAS)

Application with target observations – Strategy: Breeding on the Data Assimilation System BDAS (Carrassi et al., 2007 Tellus)
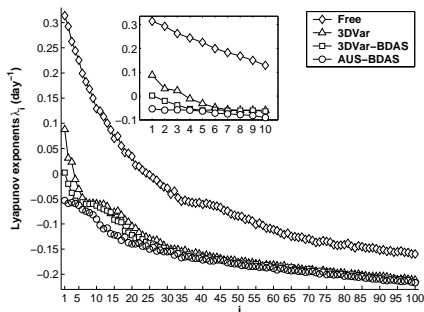
- Quasi-geostrophic atmospheric model (Rotunno and Bao, 1996 MWR)
- Perfect model setup - Observation Dense area (1-20 Longitude) - Target Area, one obs between 21-64 Longitude

| Experiment | Ocean Obs Type | Ocean Obs Location | Ocean Obs Assimil | RMS Error |
|------------|----------------|--------------------|--------------------|-----------|
| LO | - | - | - | 0.462 |
| FO | sounding | fixed (x=42, y=16) | 3DVar | 0.338 |
| RO | sounding | random | 3DVar | 0.311 |
| 3DVar-BDAS | sounding | BDAS | 3DVar | 0.184 |
| AUS-BDAS | temperature | BDAS | AUS | 0.060 |

# exploiting chaos...→ DA as nonlinear stability problem

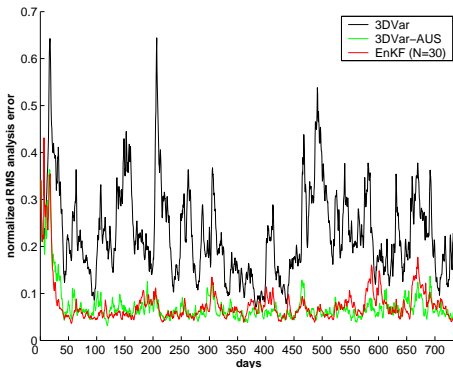**Can efficient DA methods be constructed to achieve the asymptotic stabilization of the system ?**



From Carrassi, Ghil, Trevisan & Uboldi, 2008 CHAOS

- DA provides a stabilizing effect (compare 3DVar with free system Lyapunov spectrum) but ...
- if the DA is designed to kill the instabilities, the estimation error is efficiently reduced

# exploiting chaos...→ Hybrid 3DVar - AUS

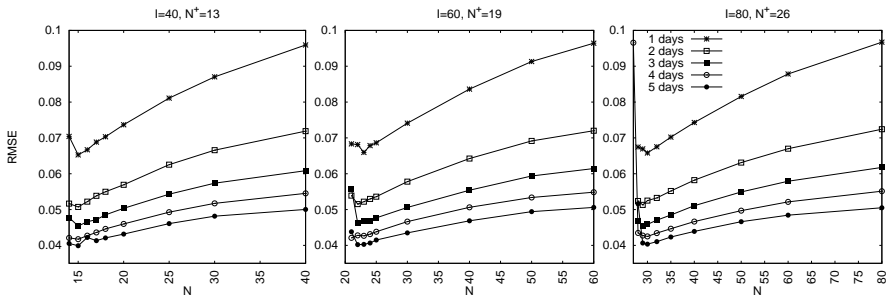## Enhancing the performance of a 3DVar by using AUS
### Comparison with EnKF



Adapted from Carrassi, Trevisan, Descamps, Talagrand & Uboldi, 2008 NPG

- A network of randomly distributed obs (vertical soundings)
- 3DVar-AUS: (1) AUS assimilate the obs able to control an unstable mode; (2) 3DVar process the remaining obs
- 3DVar-AUS comparable to EnKF with only one BDAS mode ⇒ Reduced computational cost and implementation on a pre-existing 3DVar scheme

# exploiting chaos...→ 4DVar-AUS

4DVar-AUS: The analysis increment is confined in the unstable and neutral subspace by applying to 4DVar the AUS constraint
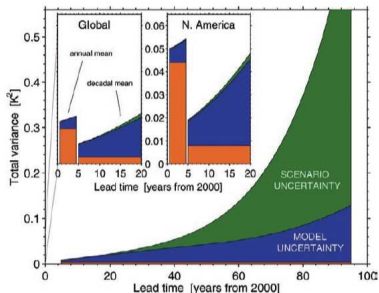


From Trevisan, D'Isidoro & Talagrand, 2010 QJRMS

- Lorenz 40 variables
- The assimilation is performed in a subspace of dimension $N = N^0 + N^+$
- When $N = n$, the standard 4DVar is recovered
- The error of 4DVar-AUS is smaller than the error of 4DVar, particularly for short assimilation windows, when the errors in the stable directions are not yet damped
- It exists an optimal subspace dimension for the assimilation that is approximately equal to $N^+ + N^0$
- 4DVar-AUS does not need tangent/adjoint model
- See Trevisan & Palatella 2011 NPG for the EKF-AUS and Palatella, Carrassi & Trevisan, 2012 JPA for a review of the AUS algorithms
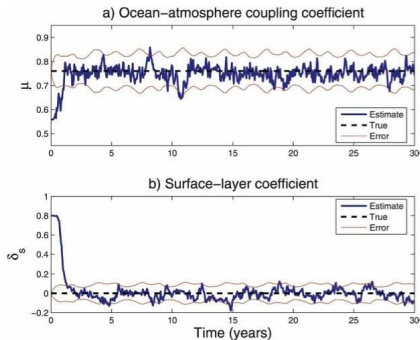
# Climate Prediction Error

Relative importance of error sources in climate prediction



From Hawkins & Sutton, 2009 BAMS

- For time horizons of many decades or longer, the dominant sources of uncertainty at regional or larger spatial scales are model uncertainty and scenario uncertainty.
- For time horizons of a decade or two, the dominant sources of uncertainty on regional scales are model uncertainty and internal variability.
- The importance of internal variability increases at smaller spatial scales and shorter time scales.

# Climate Prediction and parameter estimation



a) Ocean–atmosphere coupling coefficient
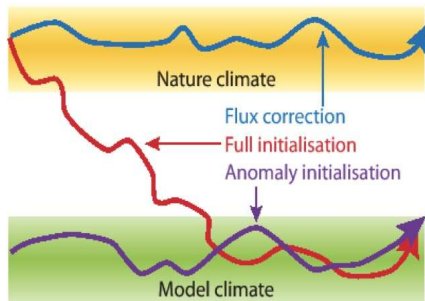
b) Surface–layer coefficient

From Kondrashov, Sun & Ghil, 2008 MWR

- Intermediate Atmosphere-Ocean Coupled Model of the tropical Pacific Ocean
- Prognostic Upper Ocean coupled with a Diagnostic Atmosphere
- Uncertain parameters: relative coupling coefficient $\mu$ and surface layer coefficient $\delta_s$
- State Augmented EKF for state and parameter estimation
- Synthetic Observations of SST

Parameter Estimation is nowadays of central importance in Data Assimilation (see talk by J. Ruiz)

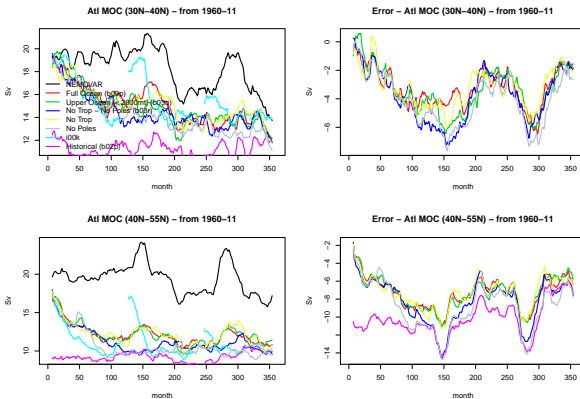# Initialization of long-term predictions



From Magnusson *et al.*, 2012 ECMWF Tech Memo 676

- Full-Field Initialization - FFI
- Anomaly Initialization: Observed anomalies are added to model climatology - AI
- FFI reduces RMSE in the short-term but it requires a bias-reduction approach to reduce the drift
- AI maintains the model trajectory on its own attractor so that drift is reduced but bias is strong
- A surface flux correction is designed a-priori to reduce model biases (Magnusson *et al.*, 2012 ECMWF Tech Memo 676)

# Nudging Experiments with Ec-Earth climate model
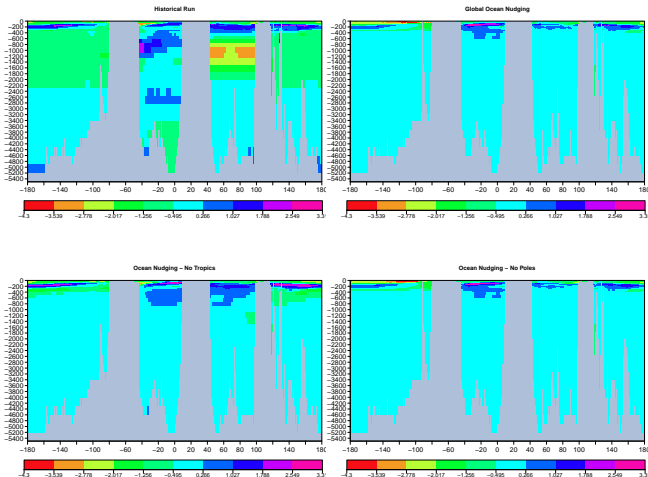
Atlantic Meridional Overturning Circulation



Acknowledgments to V. Guemas, D. Volpi and F. Doblas-Reyes of the CFU team

- Ocean Nudging toward Temperature and Salinity
- Relaxation time: $1/10$ days$^{-1}$ (below mix-layer and 800 mt); $1/2$ years$^{-1}$
- Observation DataSet: NEMOVAR reanalysis
- Classical Nudging Approach is adopted here. See Auroux & Blum, 2008 NPG for a novel advanced nudging technique.

# Nudging Experiments with Ec-Earth climate model
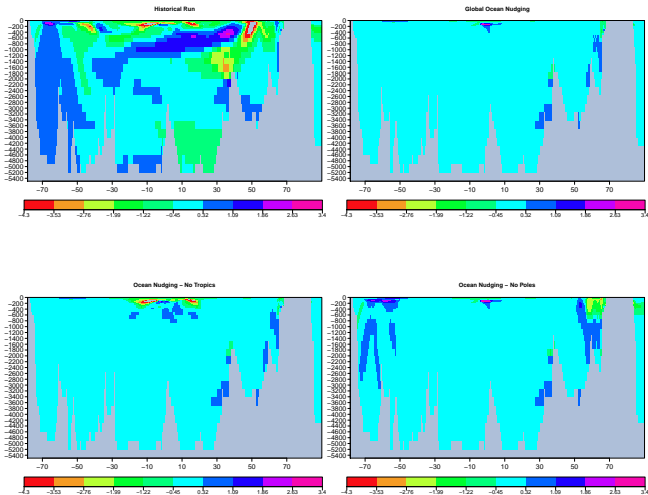
Reconstruction of Ocean Temperature



Equatorial Vertical Section of Temperature Error. Averages are taken over the period 1985-1990.

# Nudging Experiments with Ec-Earth climate model



Longitudinal (30W) Vertical Section of Temperature Error. Averages are taken over the period 1985-1990. Acknowledgments to V. Guemas, D. Volpi and F. Doblas-Reyes of the CFU team
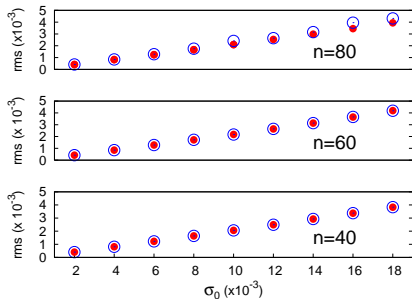
## Prospects

- Data assimilation for system possessing many scale of motions. Lorenc & Payne, 2007 QJRMS for variational scheme. What about ensemble-based algorithms ?
- Treatment of model error ...
- Design of computationally suitable algorithms for the initializations of seasonal-to-decadal prediction
- Parameter estimation approach for the online reduction of bias in long-term forecast
- Advanced nudging techniques are they feasible for initializing climate predictions ?

# exploiting chaos...→ EKF-AUS

EKF-AUS: The analysis is performed in a manifold of dimension $N = N^0 + N^+$



EKF (full circles), EKF-AUS (empty circles). From Trevisan & Palatella, 2011 NPG

- Lorenz 40 variables
- EKF-AUS belongs to the family of square-root implementations of the Extended Kalman Filter
- The assimilation is performed in a manifold of dimension $N = N^0 + N^+$.
- When $N = n$ the standard EKF is recovered
- When $N = N^+ + N^0$ the reduced form, with Assimilation in the Unstable Subspace (EKF-AUS) is obtained.
- See Palatella, Carrassi & Trevisan, 2012 JPA for a review on the AUS algorithms