

# Operational Data Assimilation at ECMWF and medium term plans

Lars Isaksen

Head of Data Assimilation Section

ECMWF

Reading, UK

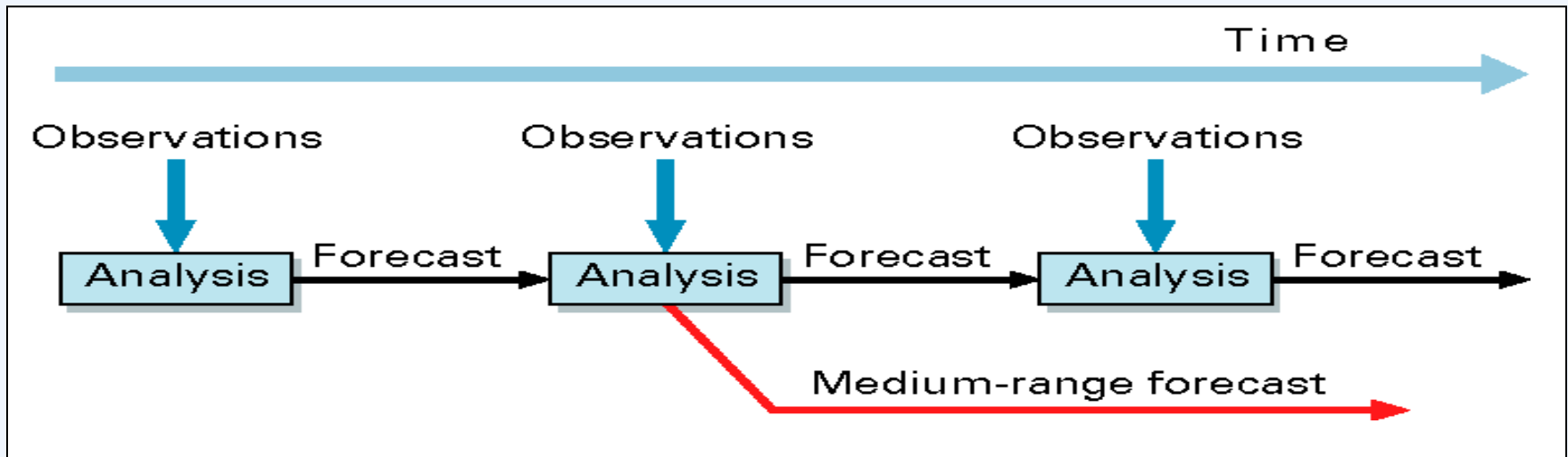
SANGOMA Kick-off Meeting 24-25 November 2011

Acknowledgements for ECMWF staff contributions to my presentation from  
Massimo Bonavita, Deborah Salmond, Yannick Trémolet,  
Kristian Mogensen, Dick Dee and Anne Fouilloux

# Outline

- Operational Data Assimilation at ECMWF
- Resolution and performance
- Hybrids methods: the best of both worlds?
- Ensemble of Data Assimilations (EDA)
- Future DA software: modularity and flexibility
- Scalability issues

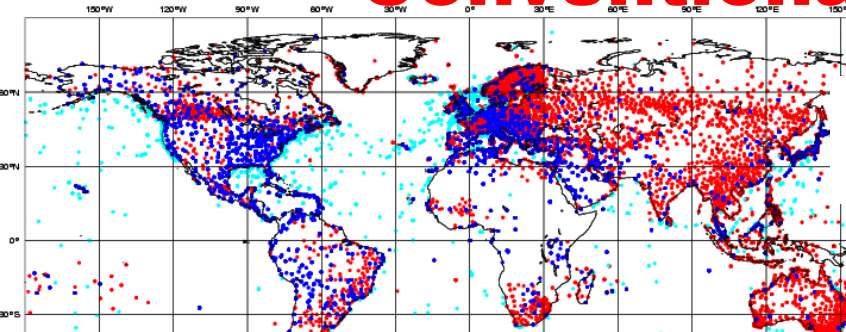
# Data assimilation system (4D-Var)



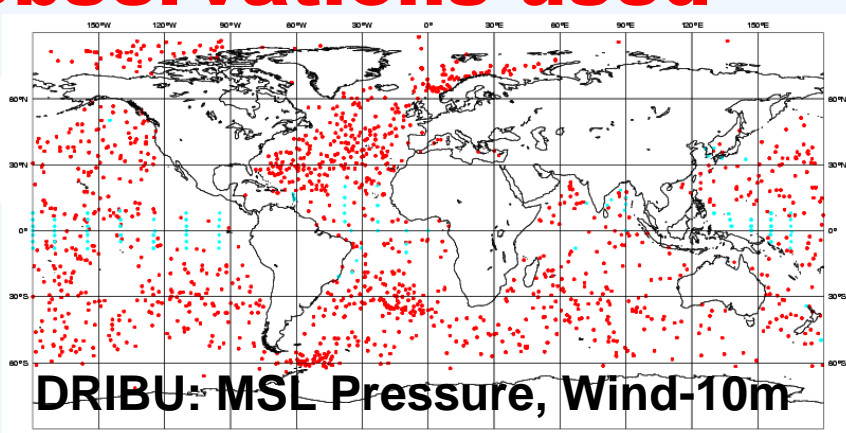
- The observations are used to correct errors in the short forecast from the previous analysis time.
- Every 12 hours we assimilate 7 – 9,000,000 observations to correct the 80,000,000 variables that define the model's virtual atmosphere.
- This is done by a careful 4-dimensional interpolation in space and time of the available observations; this operation takes as much computer power as the 10-day forecast.



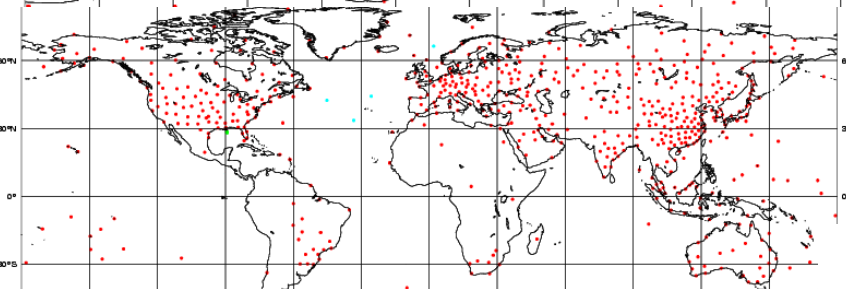
# Conventional observations used



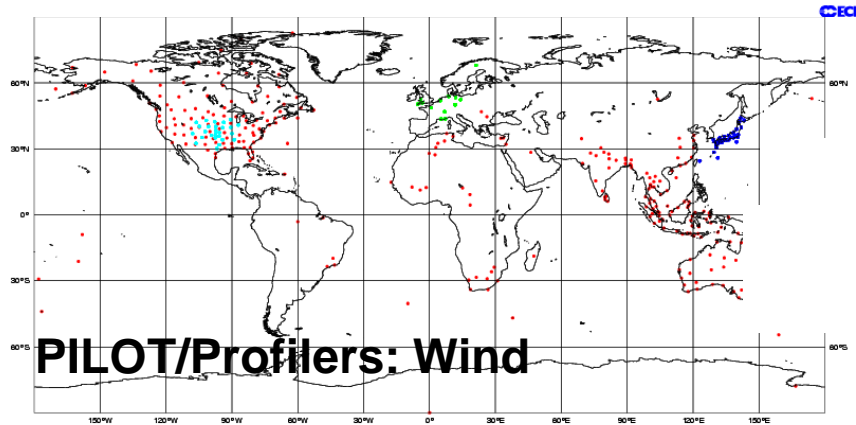
**SYNOP/METAR/SHIP:**  
MSL Pressure, 10m-wind, 2m-Rel.Hum.



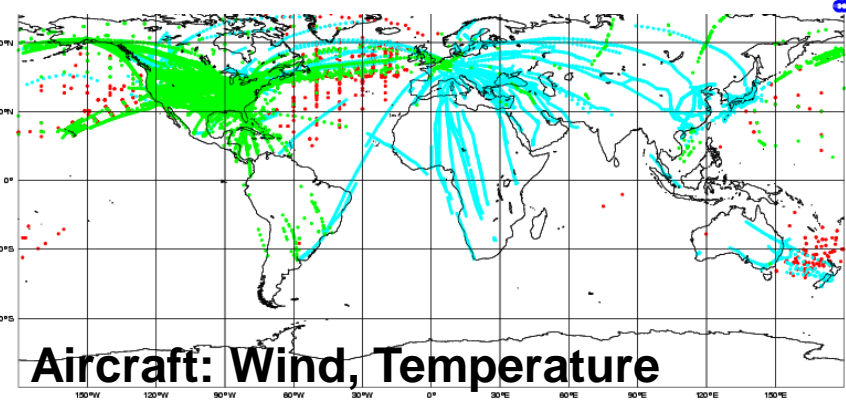
**DRIBU:** MSL Pressure, Wind-10m



**Radiosonde balloons (TEMP):**  
Wind, Temperature, Spec. Humidity



**PILOT/Profilers:** Wind



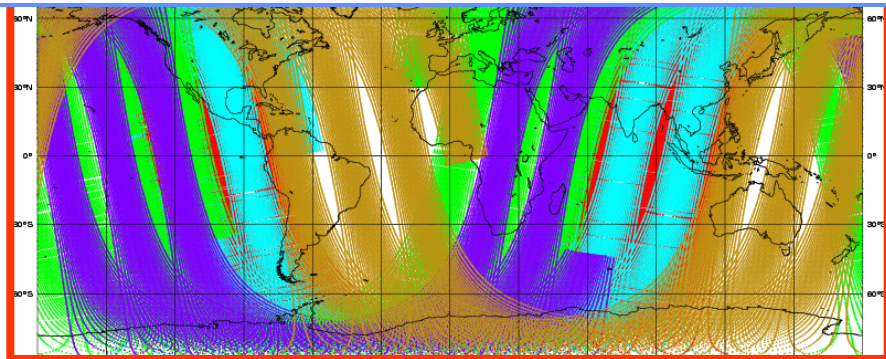
**Aircraft:** Wind, Temperature

**Note: We only use a limited number of the observed variables; especially over land.**

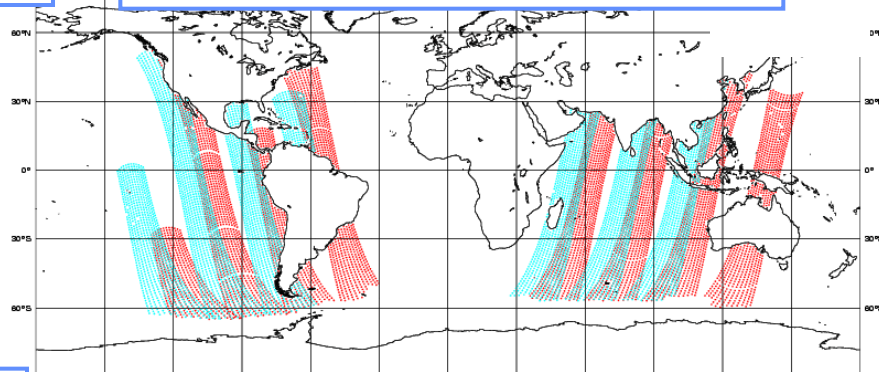


# Satellite data sources used in the operational ECMWF analysis

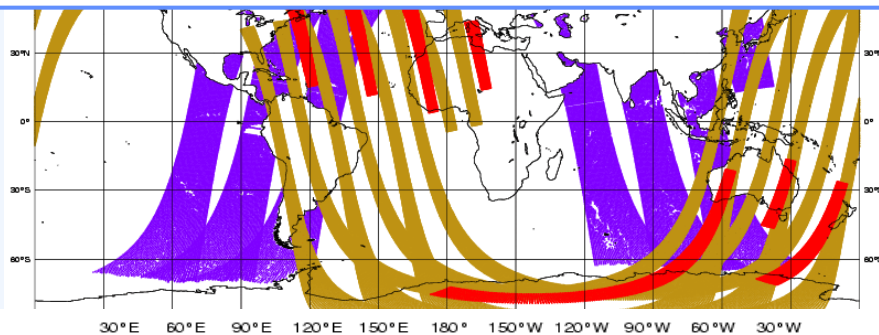
13 Sounders: NOAA AMSU-A/B, HIRS, AIRS, IASI, MHS



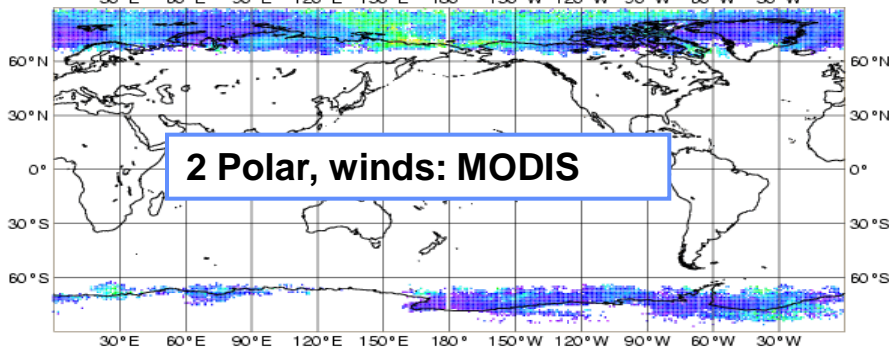
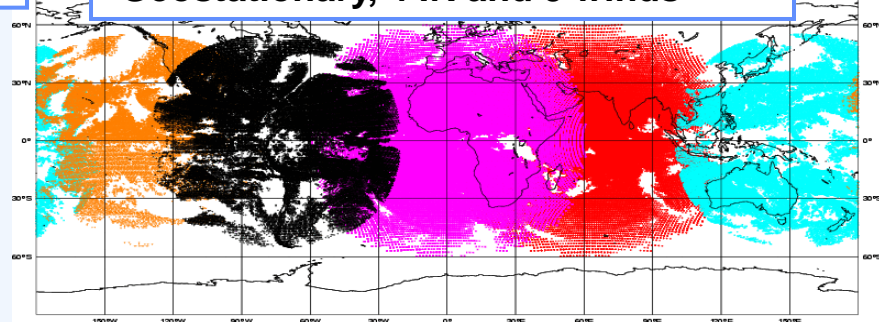
5 imagers: 3xSSM/I, AMSR-E, TMI



3 Scatterometer sea winds: ERS, ASCAT, QuikSCAT

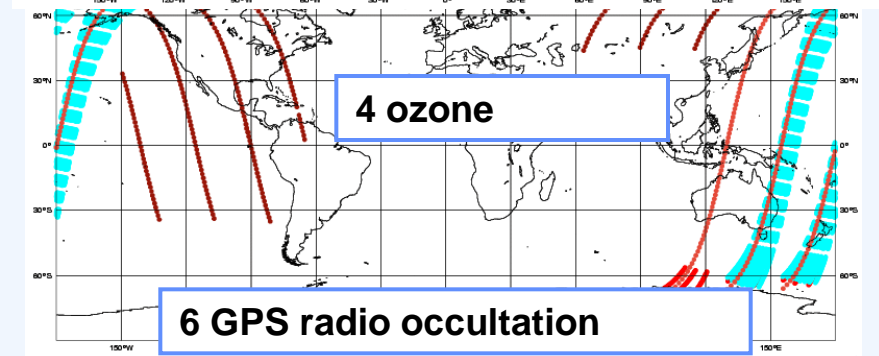


Geostationary, 4 IR and 5 winds



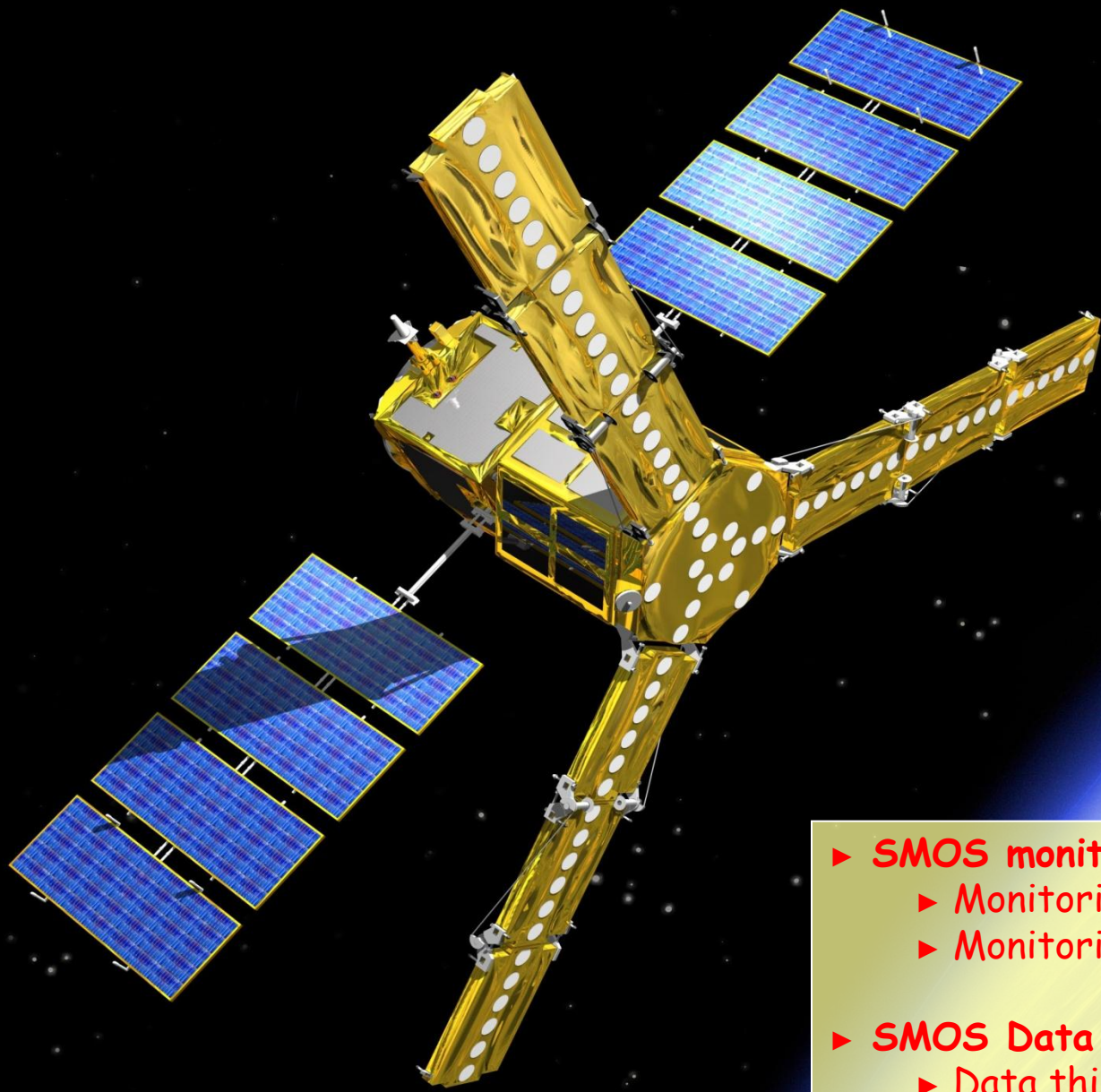
2 Polar, winds: MODIS

4 ozone



6 GPS radio occultation



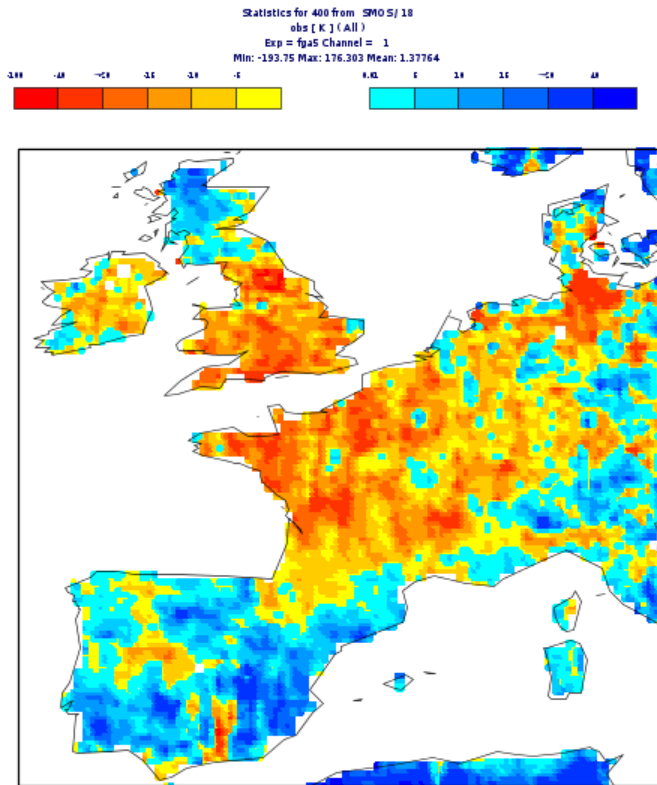


- ▶ **SMOS monitoring (Phase-I)**
  - ▶ Monitoring developments
  - ▶ Monitoring webpage
- ▶ **SMOS Data Assimilation (Phase-II)**
  - ▶ Data thinning
  - ▶ Noise Filtering
  - ▶ Bias correction

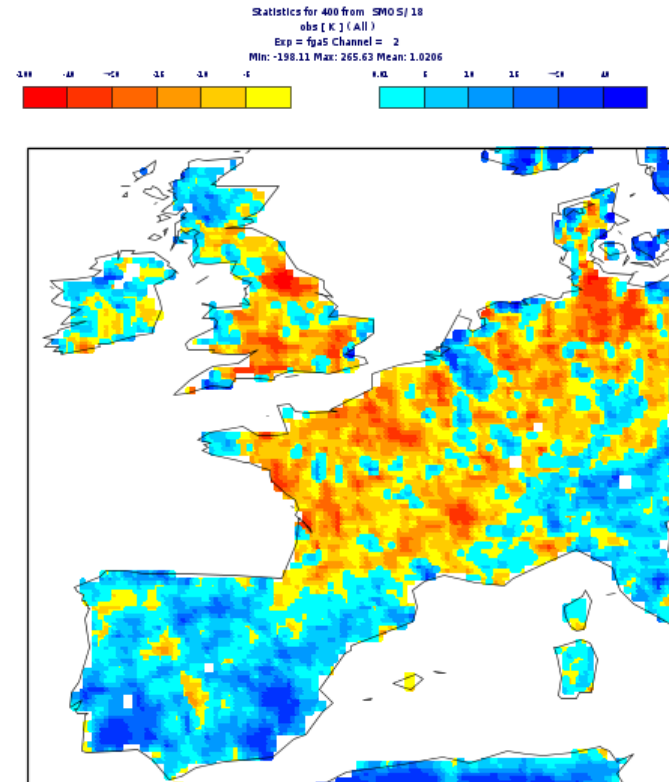
# UK–France droughts vs. Iberian floods

$T_B$  [first – last] week of April 2011

H-pol

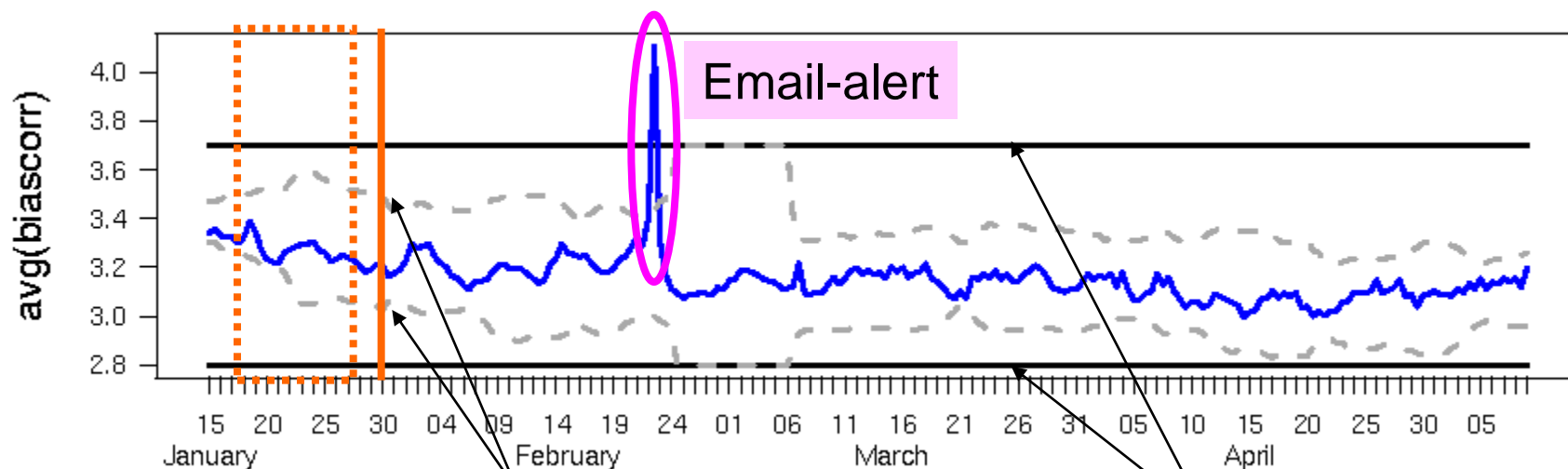


V-pol



Selected statistics are checked against an expected range.

E.g., global mean bias correction for GOES-12 (in blue):



Soft limits (mean  $\pm$  5 stdev of statistic to be checked, calculated from past statistics over a period of 20 days, ending 2 days earlier)

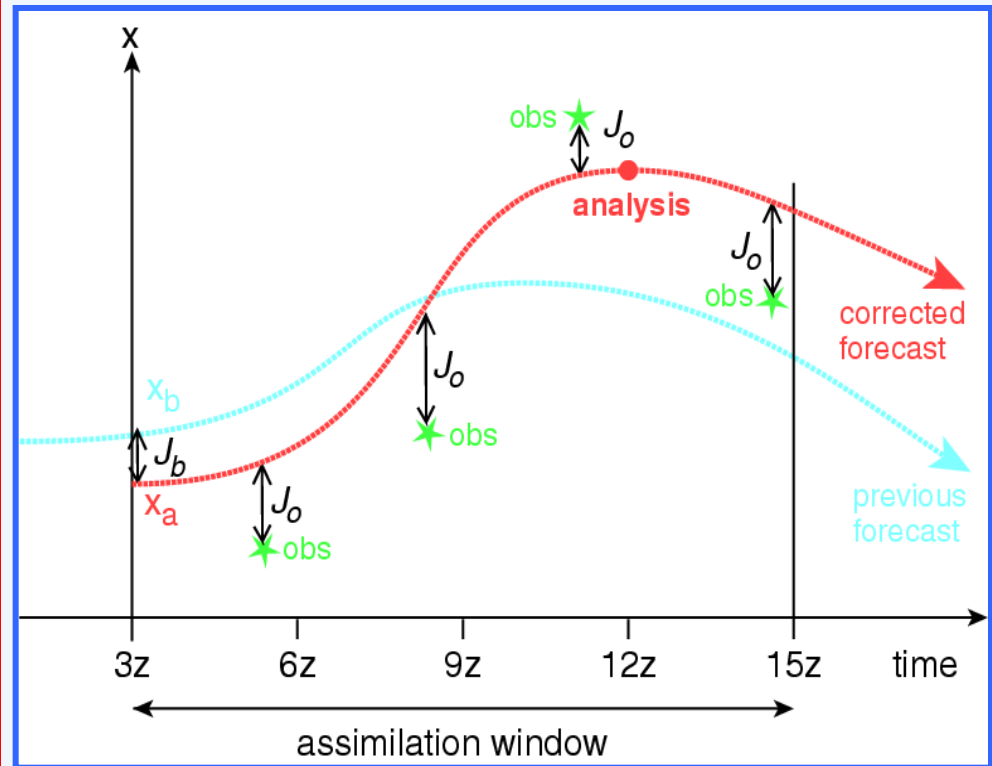
Hard limits (fixed)



# ECMWF uses a 4D-Var assimilation system

All observations within a 12-hour period (~8,000,000) are used simultaneously in one global (iterative) estimation problem

- Observation minus model differences are computed at the observation time using the full forecast model at T1279 (16 km) resolution
- 4D-Var finds the 12-hour forecast evolution that optimally fits the available observations. A linearized forecast model is used in the minimization process based on the adjoint method
- It does so by adjusting **surface pressure**, the upper-air fields of **temperature**, **wind**, **specific humidity** and **ozone**
- The analysis vector consists of **80,000,000** elements at T255 resolution (80 km)

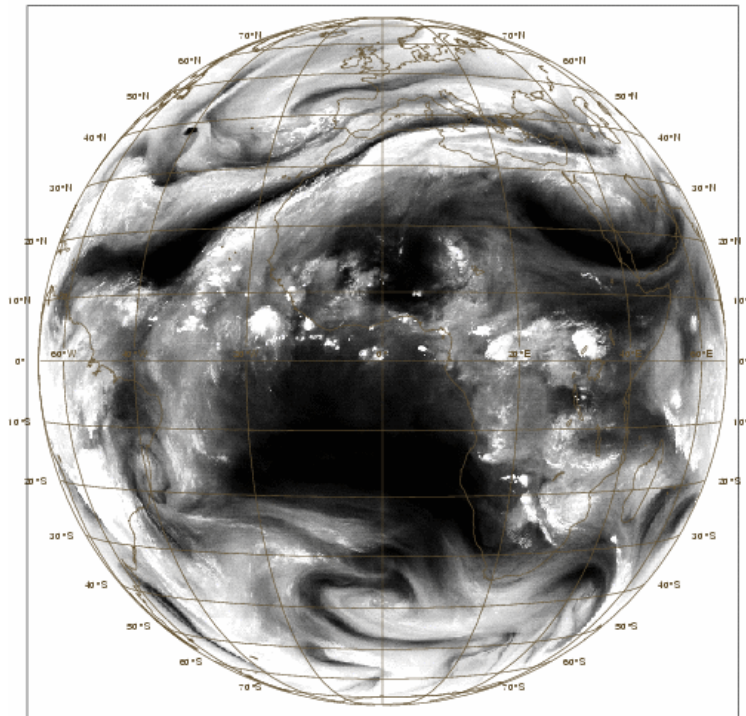


# Variational assimilation of satellite radiances

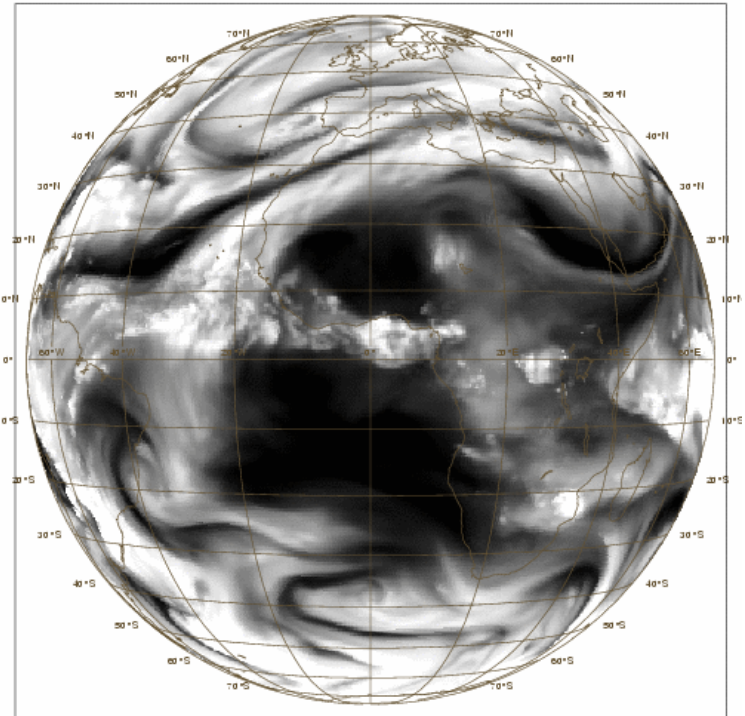
Adjust the model state to improve the match with observations:

$$\mathbf{J}(\mathbf{x}) = (\mathbf{x}_b - \mathbf{x})^T \mathbf{B}^{-1} (\mathbf{x}_b - \mathbf{x}) + [\mathbf{y} - \mathbf{h}(\mathbf{x})]^T \mathbf{R}^{-1} [\mathbf{y} - \mathbf{h}(\mathbf{x})]$$

**y** : observed radiances



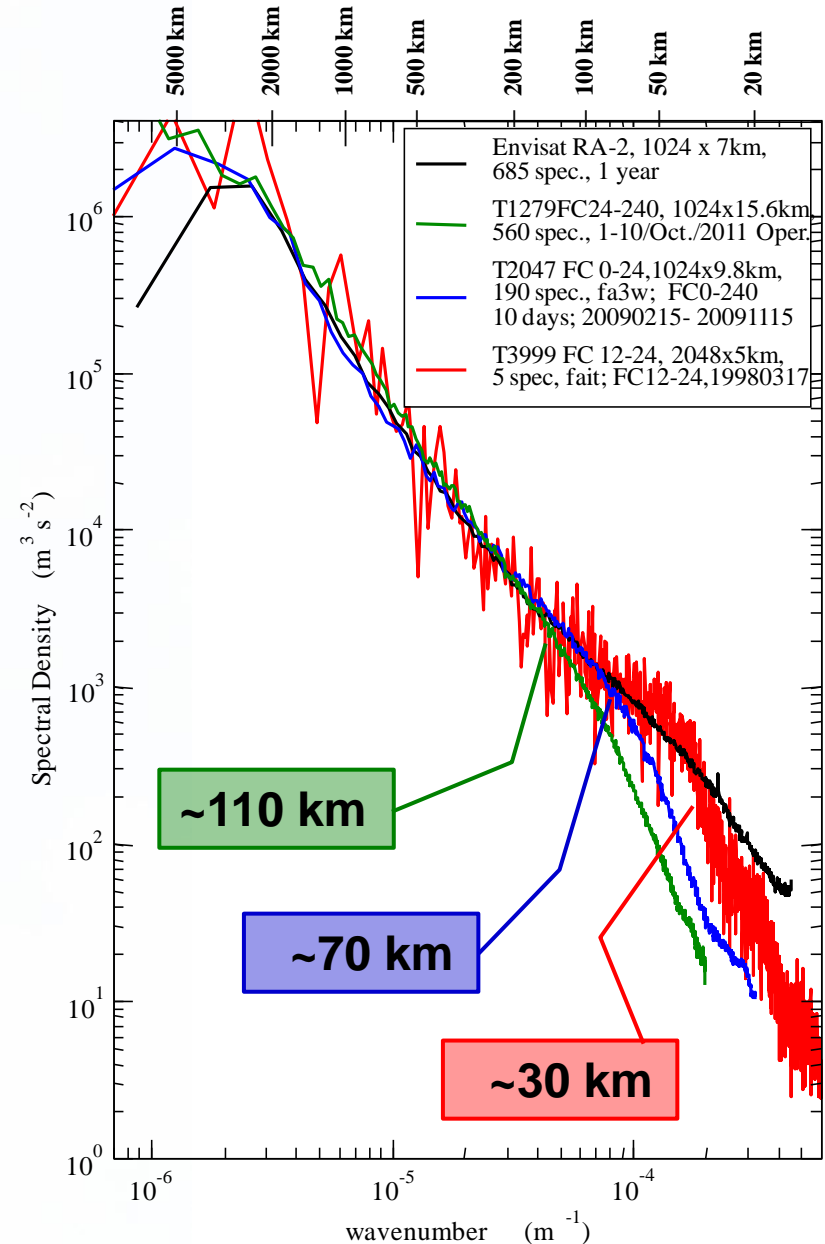
**h(x)** : model-simulated radiances



- ECMWF HPC systems
  - At the moment IBM Power6 (2x9200 cores)
  - Is now being upgraded to IBM Power7 (2x24500 cores)
- Operational Forecast and 4D-Var assimilation configuration
  - We are using the IFS - Integrated Forecast System
  - 10-day T1279L91 Forecast (16 km horizontal grid)
  - 12 hour 4D-Var T1279 outer loop T255/T159 inner loop
  - Operational Ensemble of Data Assimilations (EDA)
    - 10 member 4D-Var T399 outer and T95/T159 inner loop

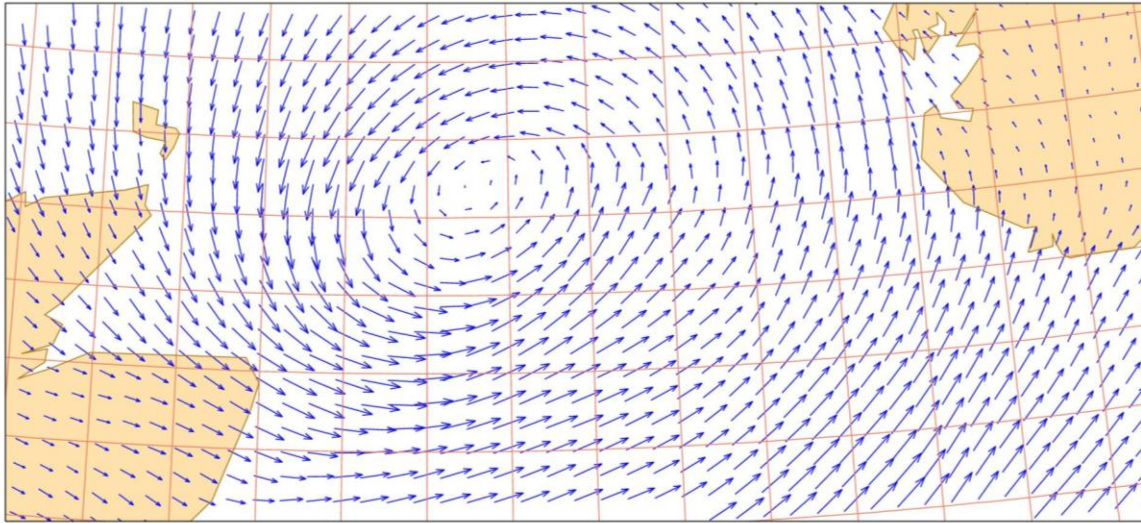
# Impact of Model Resolution

- **Black line: Observations:**  
Envisat Altimeter RA-2 data.
- **Green line: Resolution ~16 km**  
T1279, Current ECMWF operational model resolution.
- **Blue line: Resolution ~10 km**  
T2047, Next ECMWF model resolution ~2014.
- **Red line: Resolution ~5 km**  
T3999, ECMWF model resolution ~2020.

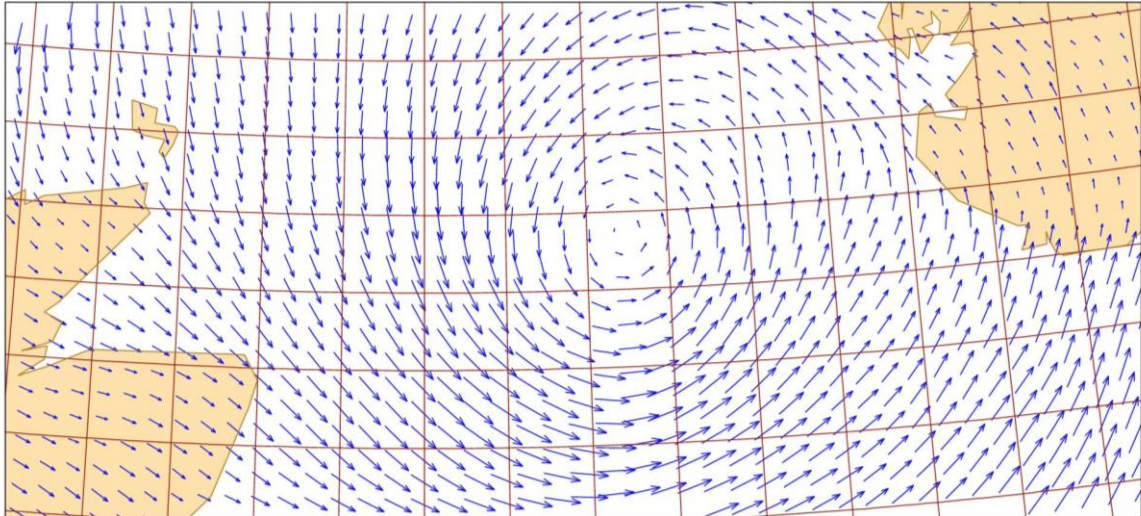


# 10m wind field from ECMWF analysis and 3h forecast T1279 resolution (16km grid) to be used from 2009

ECMWF Analysis VT:Friday 1 February 2008 00UTC Model Level 91 U velocity/V velocity 25.0m/s



Friday 1 February 2008 00UTC ECMWF Forecast t+3 VT: Friday 1 February 2008 03UTC Model Level 91 U velocity/V velocity 25.0m/s



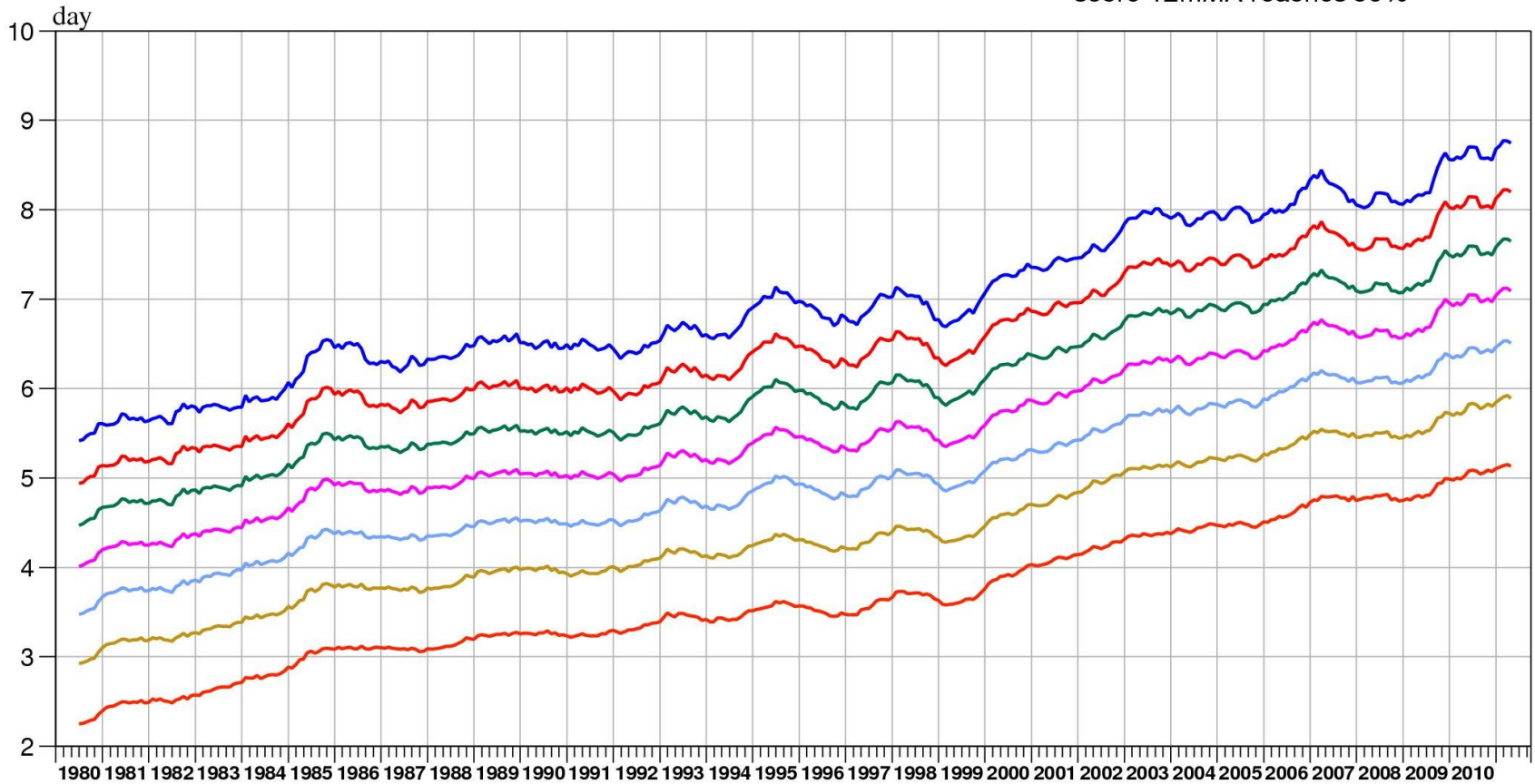
**A low in the North Sea shows the T1279 forecast model's ability to represent scales below 50km resolution**

**Each box on the plot is approx. 50kmx50km**

**Note the distinct land/sea contrast in wind speed that is well presented by the model at sub 50km resolution**

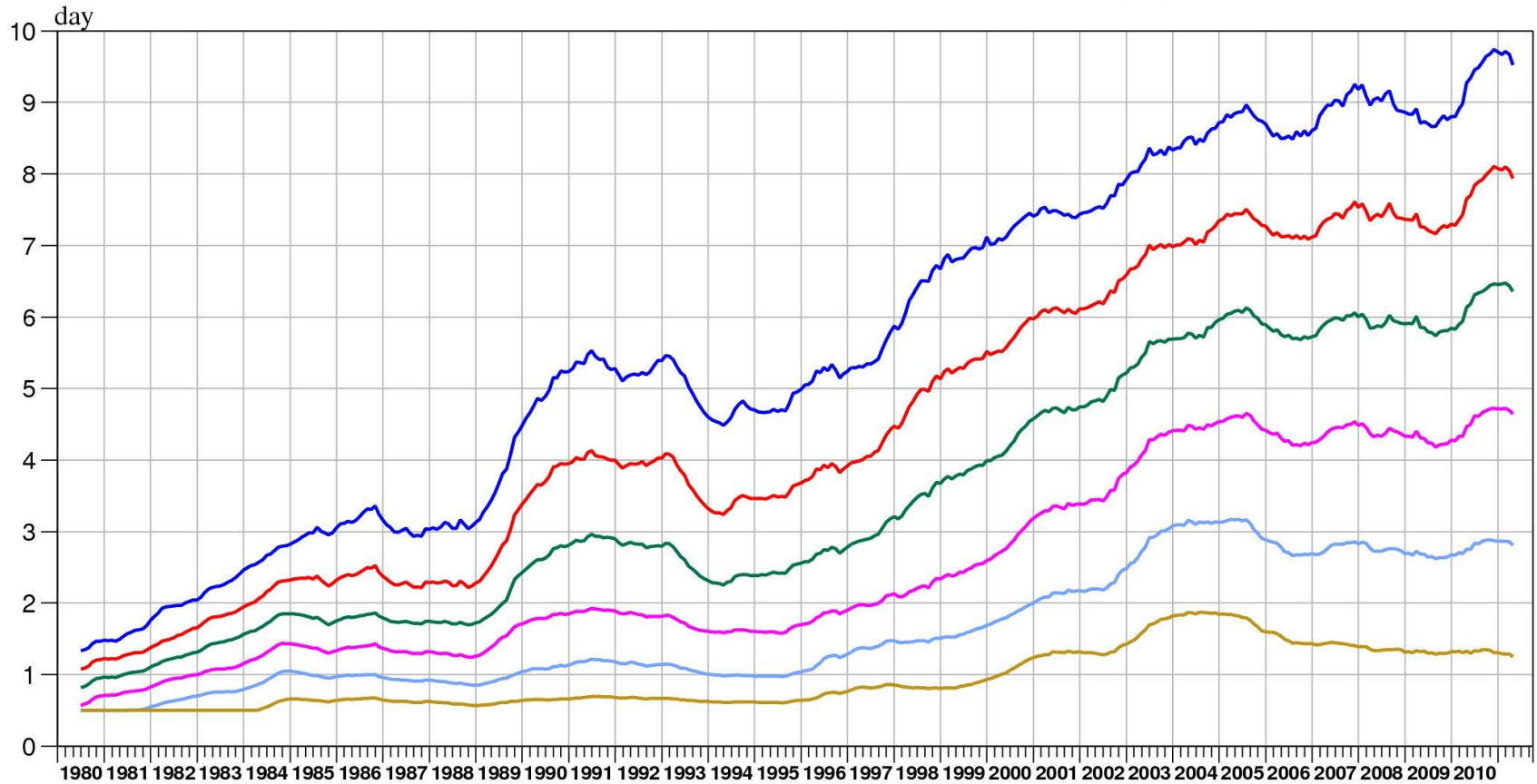
ECMWF forecast verification 12UTC  
 geopotential 500hPa  
 Correlation coefficient of forecast anomaly  
 NH Extratropics Lat 20.0 to 90.0 Lon -180.0 to 180.0  
 (12mMA = 12 months moving average)

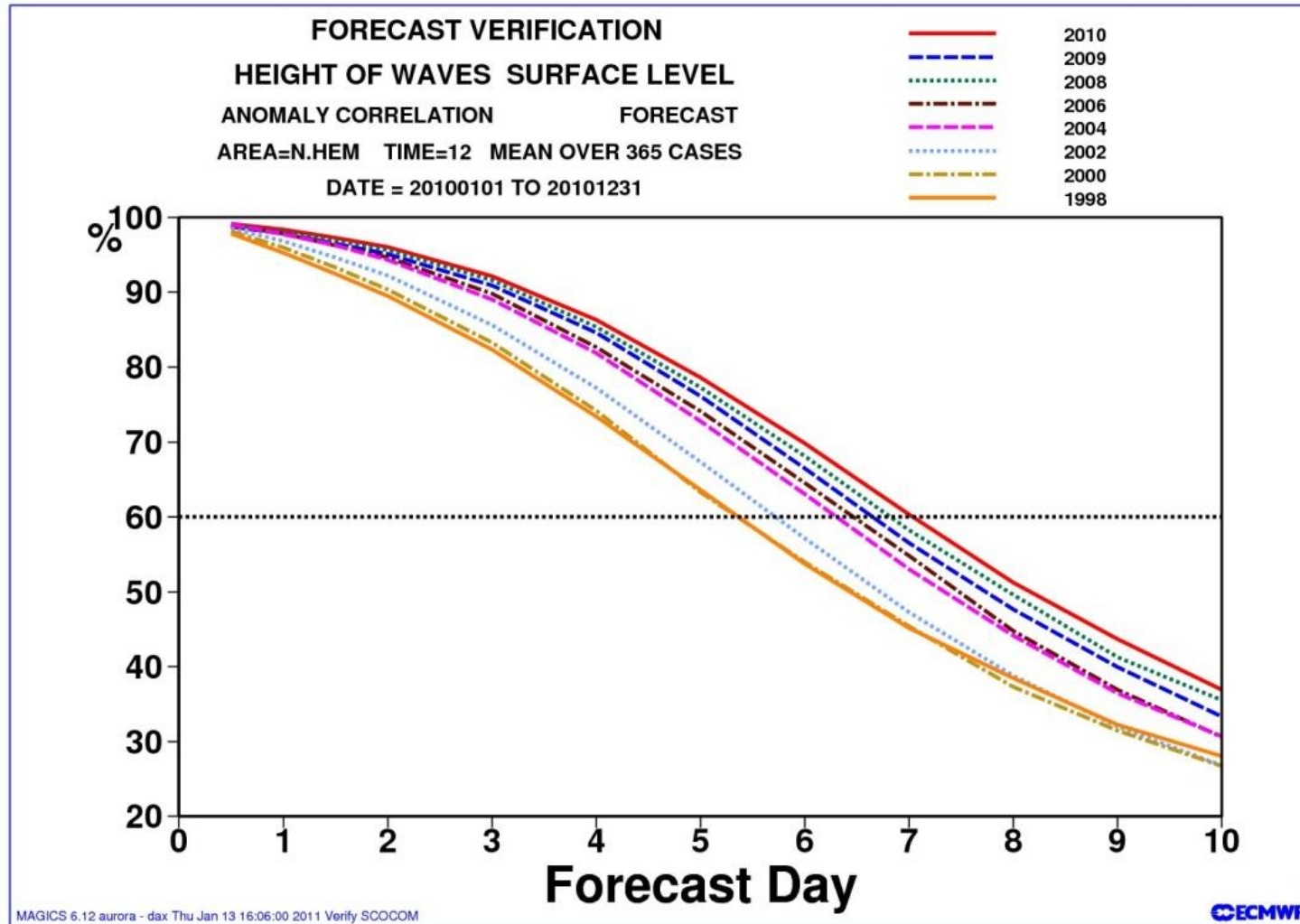
- score 12mMA reaches **60%**
- score 12mMA reaches **65%**
- score 12mMA reaches **70%**
- score 12mMA reaches **75%**
- score 12mMA reaches **80%**
- score 12mMA reaches **85%**
- score 12mMA reaches **90%**



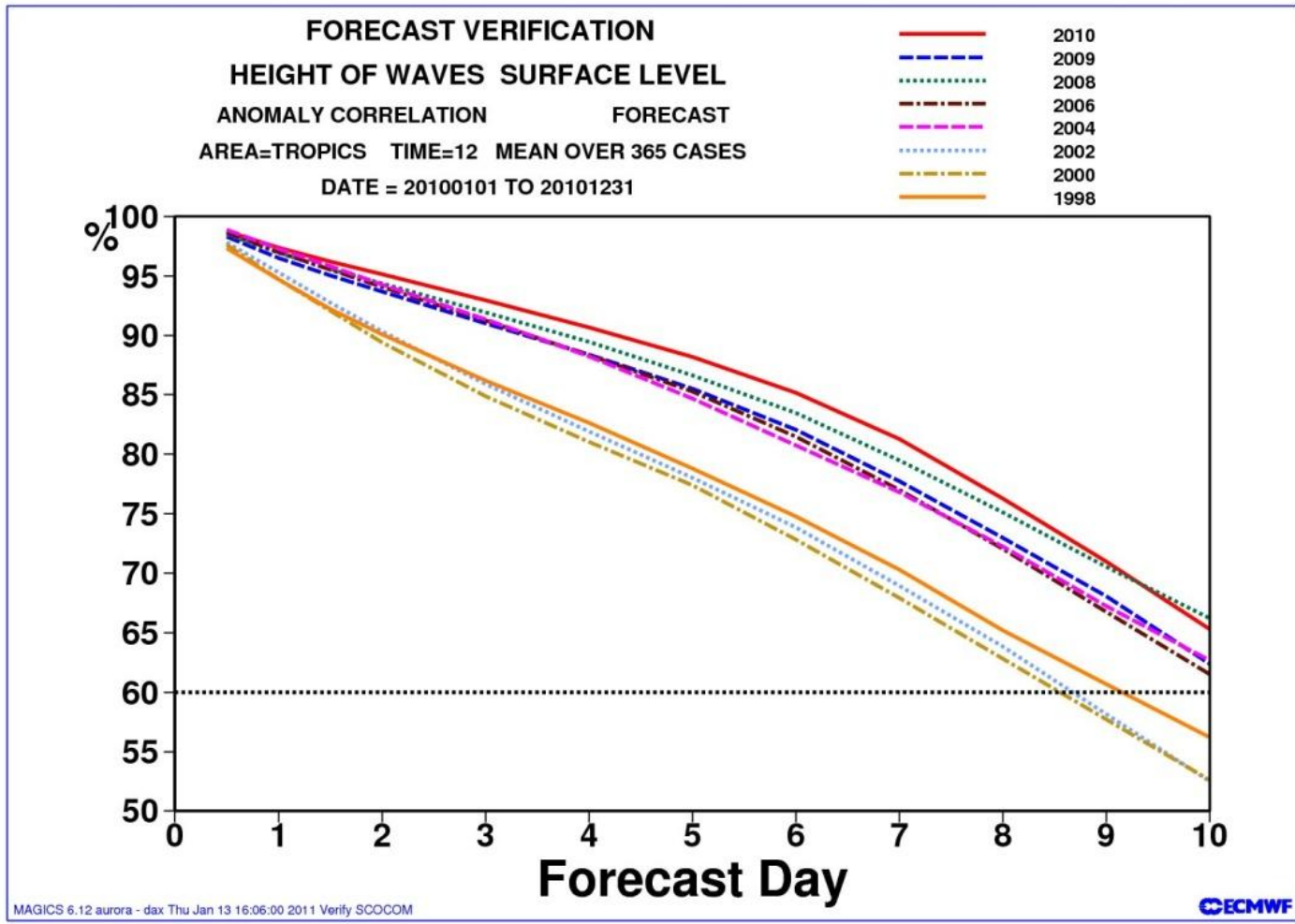
ECMWF forecast verification 12UTC  
vector wind 850hPa  
Correlation coefficient of forecast  
Tropics Lat -20.0 to 20.0 Lon -180.0 to 180.0  
(12mMA = 12 months moving average)

- score 12mMA reaches **65%**
- score 12mMA reaches **70%**
- score 12mMA reaches **75%**
- score 12mMA reaches **80%**
- score 12mMA reaches **85%**
- score 12mMA reaches **90%**

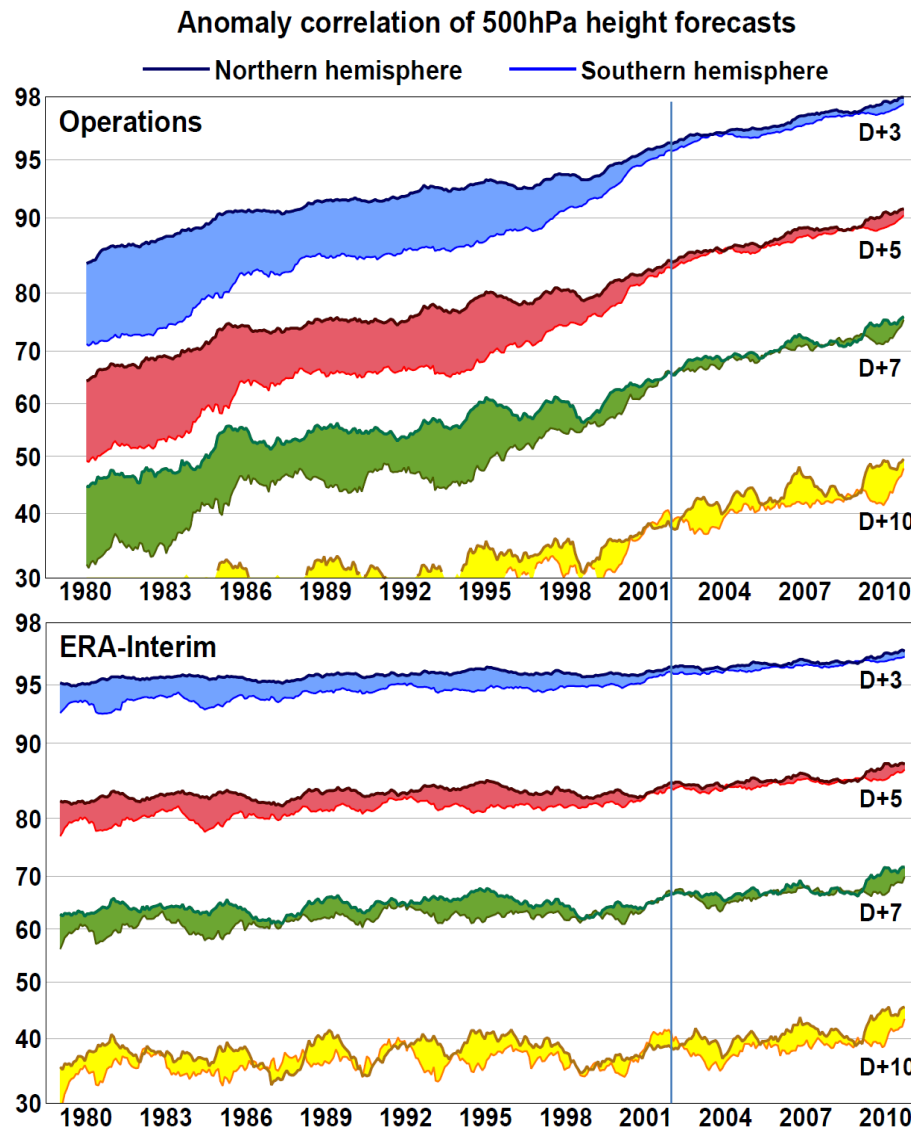








# ERA-interim: Performance compared to operations



# Why use an Ensemble of Data Assimilations (EDA)?

- a) **Kalman Filter** is computationally unfeasible for realistic NWP;
- b) **Non-sequential approx. (4D-Var)** do not cycle state error estimates: work well for short assimilation windows (6-12h), but longer windows have proved more difficult;
- c) **Sequential approx. (EnKF)** cycle low-rank estimates of state error covariances, but analysis increments are confined to perturbations subspace;

....

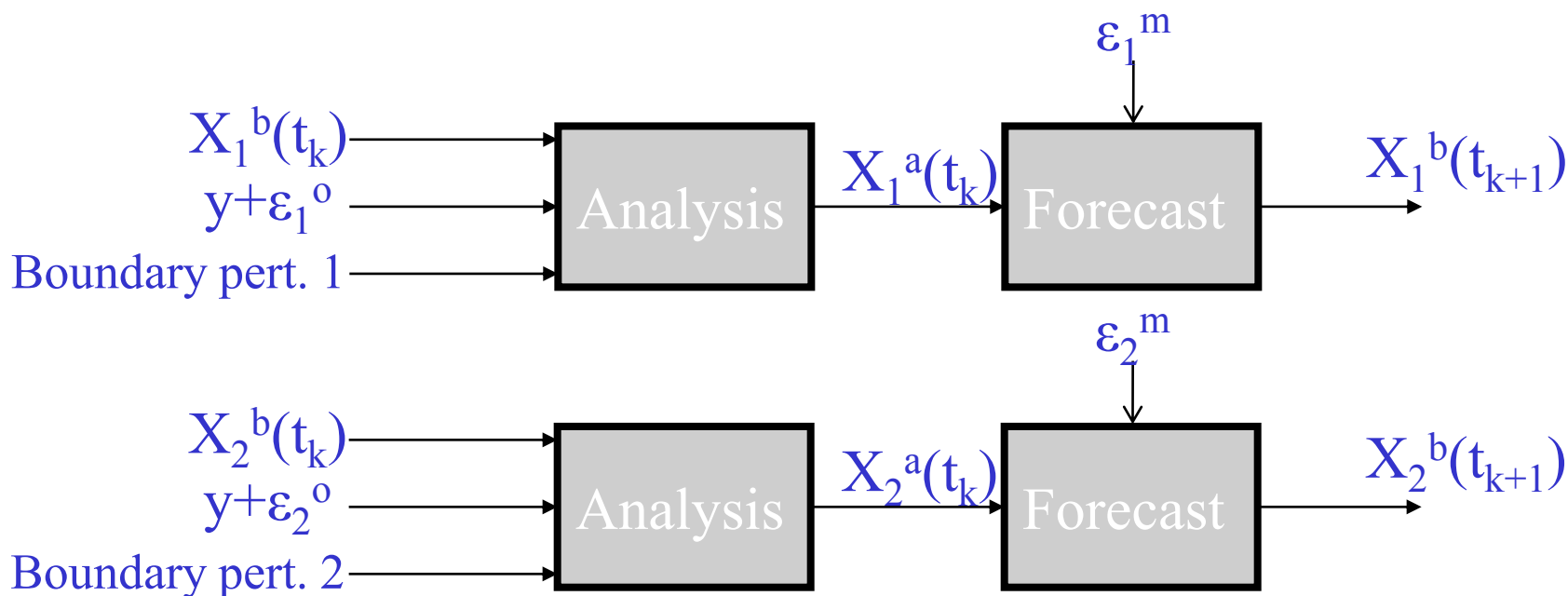
**Hybrid approach:** Use flow-dependent state error estimates (from an EnKF/EDA system) in a 3/4D-Var analysis algorithm

# Hybrid methods: How EDA works

1. We can use **an ensemble of perturbed 4D-Var** to simulate the **errors** of our reference high resolution 4D-Var
2. The ensemble of perturbed DAs should be as similar as possible to the reference DA (i.e., **same or similar  $\mathbf{K}$  matrix**)
3. The applied observation error and model error perturbations must have represent the error covariances ( $\mathbf{R}$ ,  $\mathbf{Q}$ ); however we do not need an explicit covariance model of  $\mathbf{Q}$
4. Important applications: To provide a **flow-dependent sample of background errors** at the initial time of the 4D-Var assimilation window. But also to provide analysis error estimates.

# Hybrid methods: How EDA works

The **Ensemble of Data Assimilations** (EDA, Isaksen *et al.* 2010) can be considered a **flow-dependent extension** of the way the *climatological background error matrix* is estimated (Fisher, 2003).



# Using EDA in 4D-Var to provide uncertainty estimation

**Hybrid approx.:** Use flow-dependent state error estimates (from an EnKF/EDA system) in a 3/4D-Var analysis algorithm

This solution would:

- 1) Integrate flow-dependent state error covariance information into the variational analysis
- 2) Keep the full rank representation of  $\mathbf{B}$  and its implicit evolution inside the assimilation window
- 3) More robust than pure EnKF for limited ensemble sizes and large model errors
- 4) Allow (eventual) localization of ensemble perturbations to be performed in state space;
- 5) Allow for flow-dependent QC of observations

# The operational EDA at ECMWF

- **10** ensemble members using 4D-Var assimilations
- **T399** outer loop, **T95/T159** inner loops. (Reference DA: **T1279** outer loop, **T159/T255/T255** inner loops)
- **Observations** randomly perturbed according to their specified **R**
- **SST perturbed** with realistically scaled structures
- **Model error** represented by stochastic methods (**SPPT**, Leutbecher, 2009)

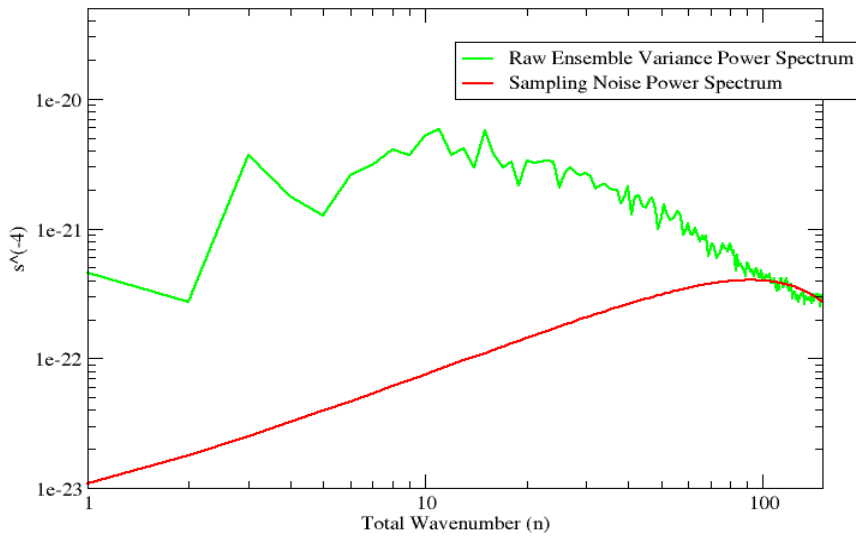
# EDA variances – sampling noise issues

a) **Sampling Noise** due to the small EDA dimensionality ( $N_{eda}=10$ )

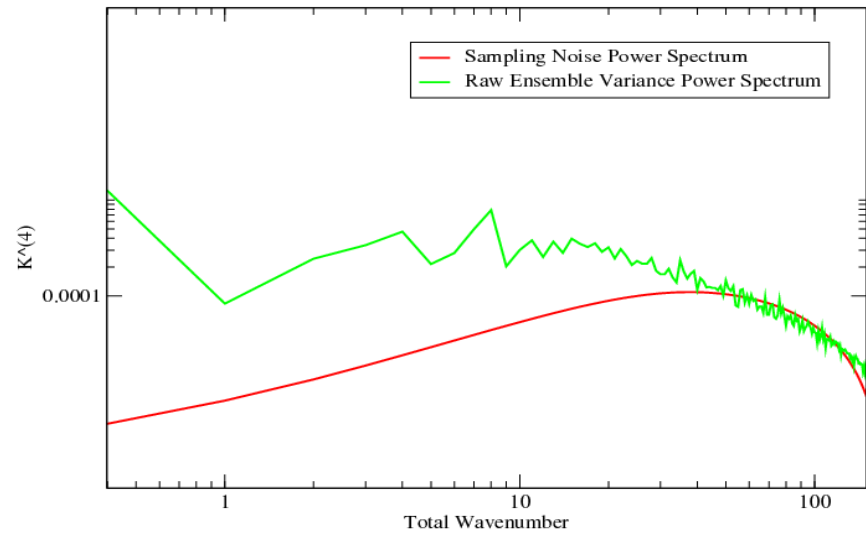
The key insight is to recognise that *sampling noise is small scale with respect to the error variance field* (Raynaud *et al.*, 2008)

We may use a **spectral filter** to disentangle noise error from the signal

Vorticity  
ml 64 (~500hPa)

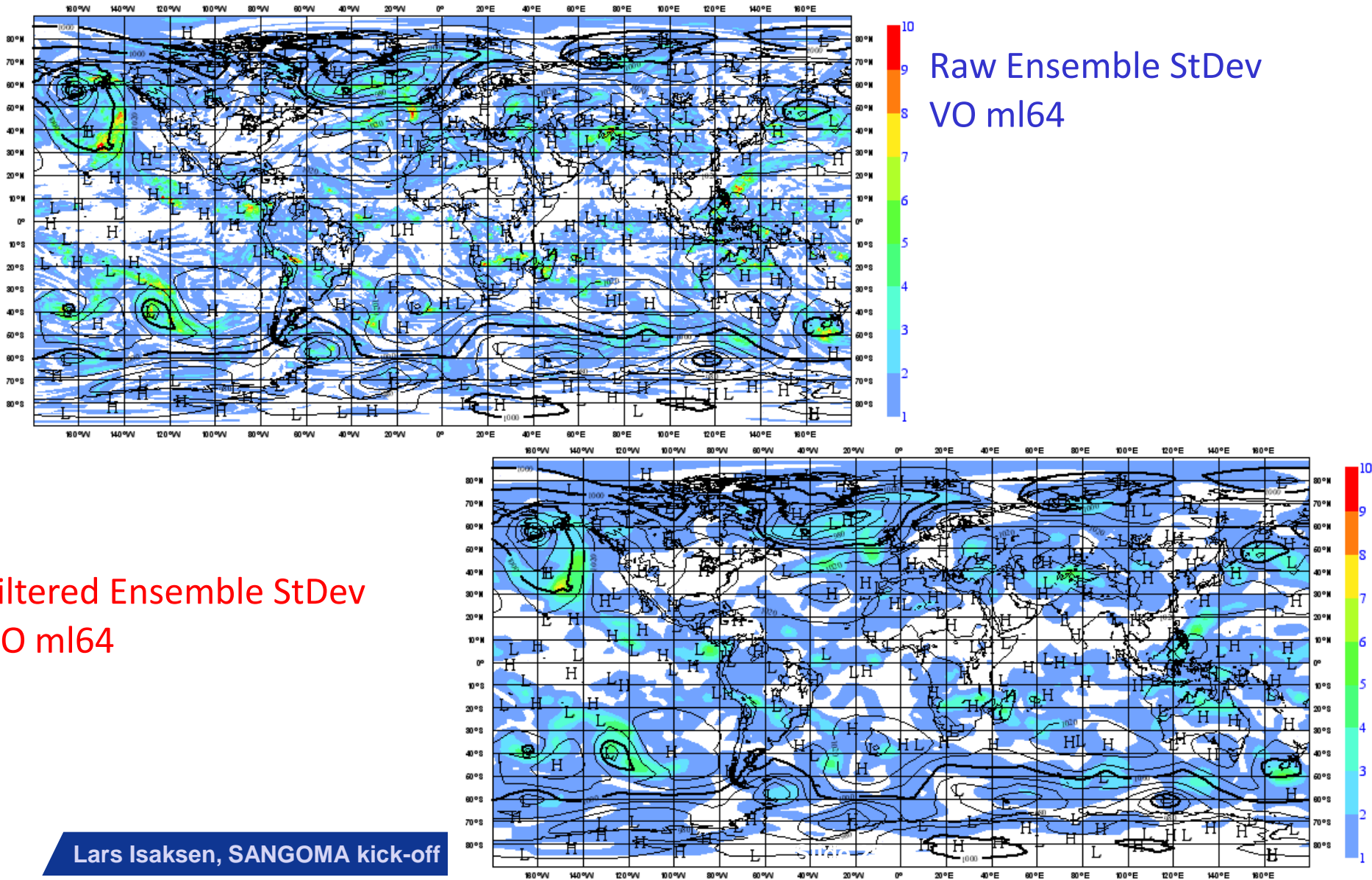


Temperature  
ml49 (~200hPa)





# EDA variances – sampling noise issues

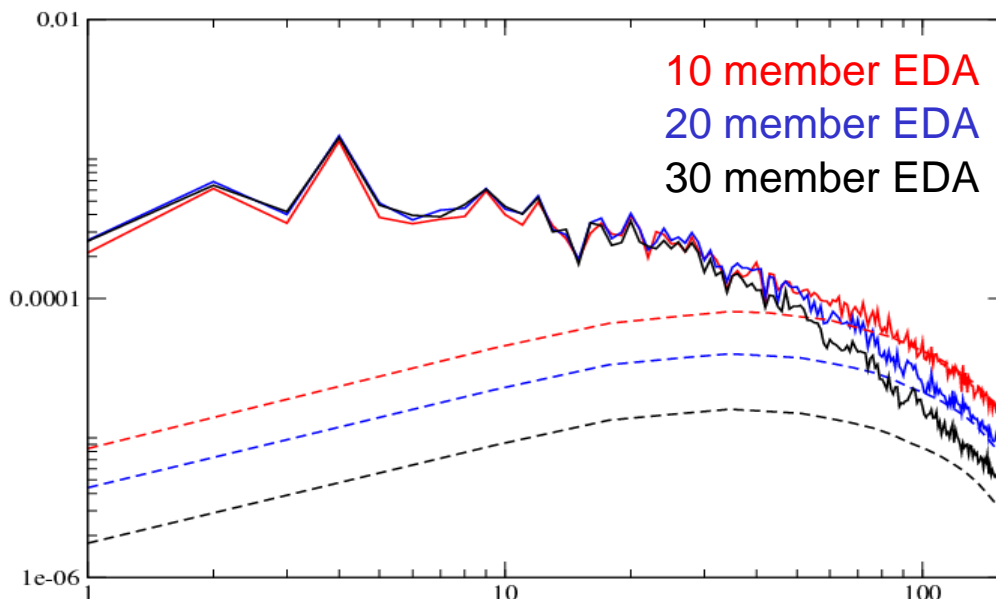


# EDA variances: Ensemble Size

The sampling noise effectively **limits the scales** that we can robustly estimate from the EDA.

The **effective spatial resolution** of the diagnosed errors is much coarser than the nominal EDA resolution (T399) and is primarily determined by the **ensemble size** (Bonavita et al., 2010)

## Temperature ml 49 (~200 hPa)



# EDA variances – systematic errors

- b) **Systematic errors** due to incorrect specification of error sources in the EDA (i.e., mis-specification of  $\mathbf{R}$ ,  $\mathbf{Q}$ , uncertainties in the boundary conditions)

A statistically consistent ensemble should satisfy:

$$\left(1 - \frac{1}{N_{ens}}\right)^{-1} \left\langle \frac{1}{N_{ens}} \sum_{i=1}^{N_{ens}} (x_i - \bar{x})^2 \right\rangle = \left(1 + \frac{1}{N_{ens}}\right)^{-1} \left\langle (\bar{x} - x^*)^2 \right\rangle$$

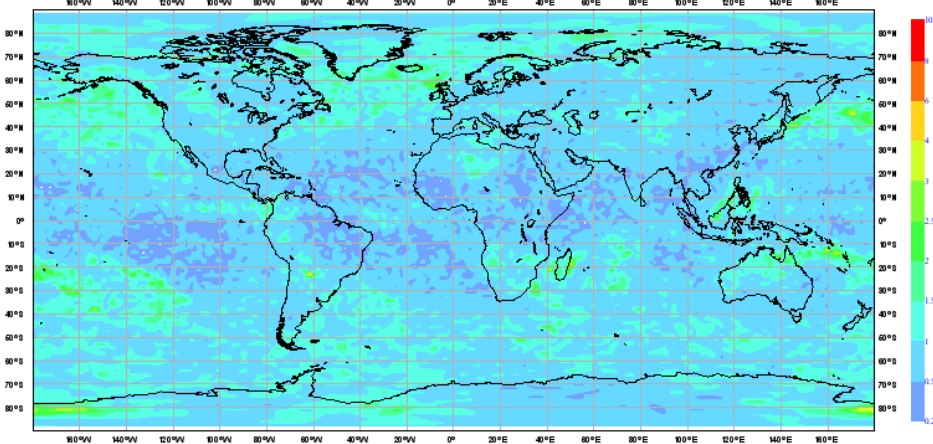
$\langle \text{ensemble variance} \rangle \approx \langle \text{squared ensemble mean error} \rangle$

# What type of errors affect EDA sample stats.?

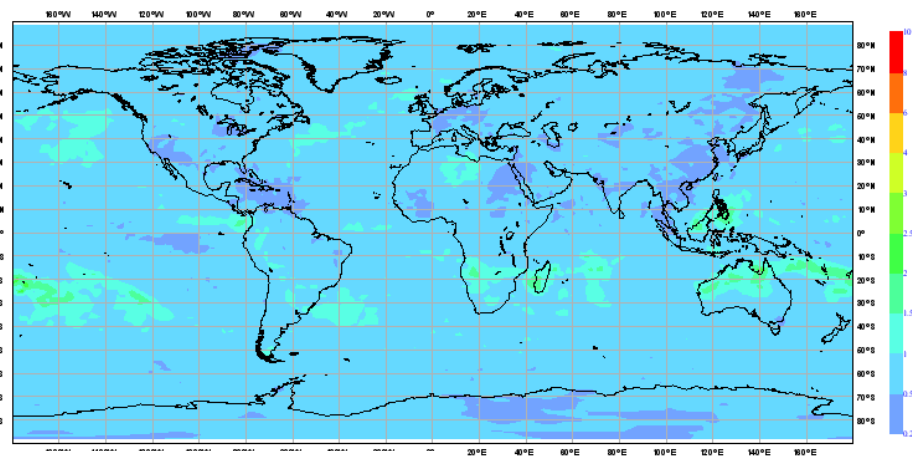
## Vorticity model level 78 (~850hPa)

### Ensemble mean error

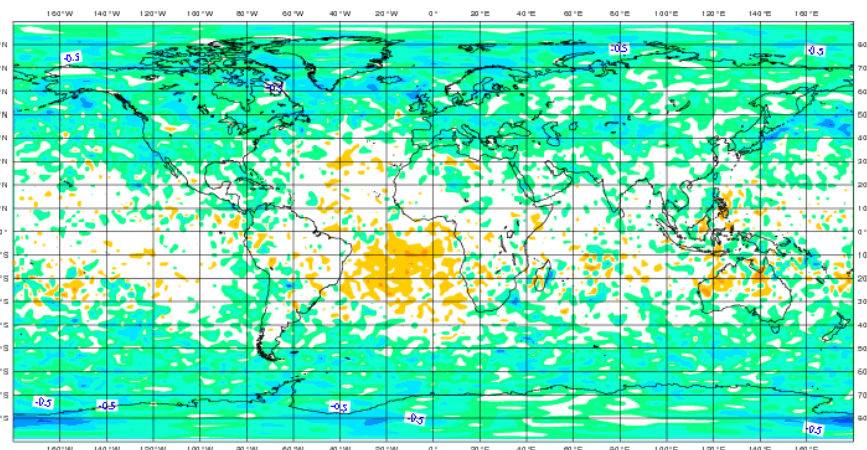
Tuesday 6 January 2009 12UTC ECMWF Forecast t-9 VT: Tuesday 6 January 2009 21UTC Model Level 78 "Vorticity (relative)



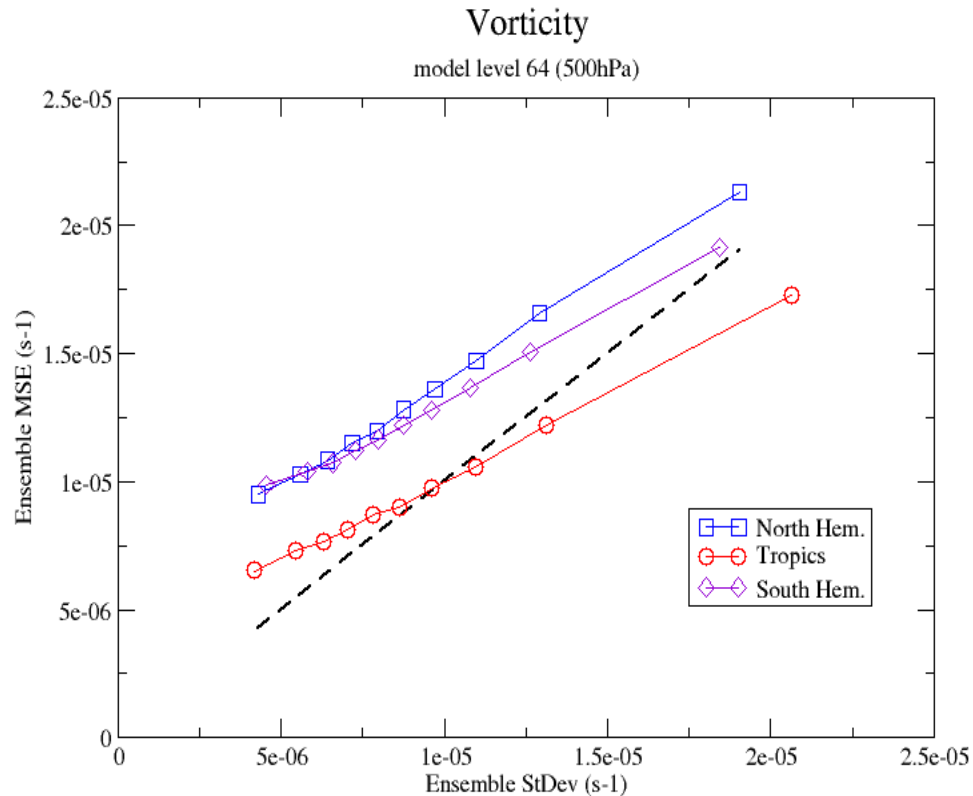
### Ensemble Spread



### Spread - Error



# EDA variances

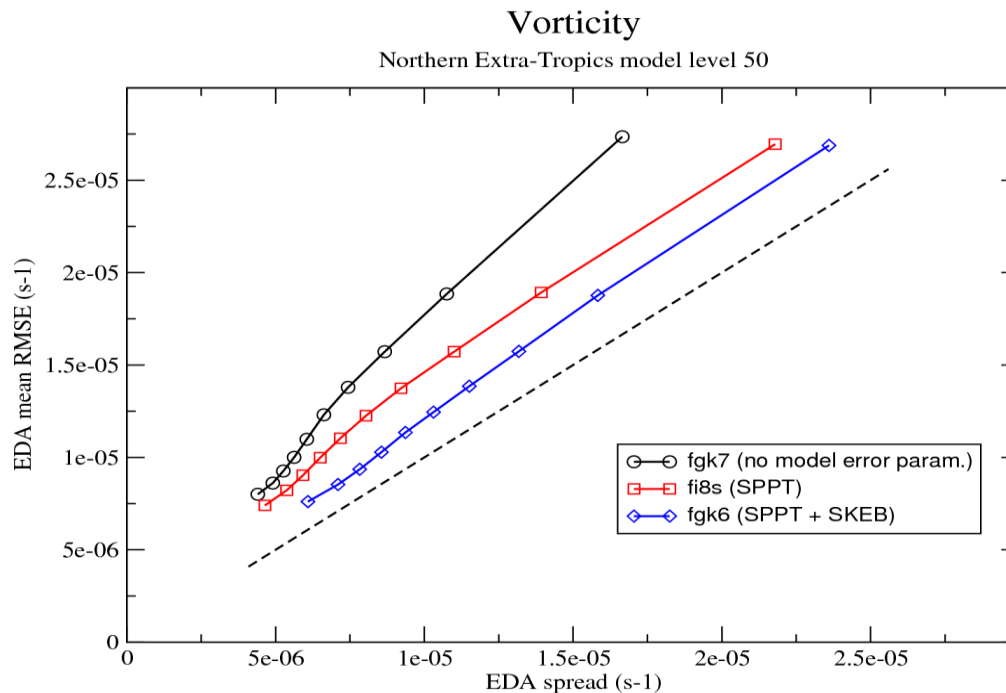


Conditional distribution of the EDA mean background RMS error for given EDA background standard deviation

# EDA variances

“Spread-Skill” regressions of the type shown serve two purposes:

1. **Diagnose** the progress (or lack thereof!) in the modelling of system uncertainties in the EDA
2. **Calibrate on-line** the EDA sample variances to obtain realistic estimates of background errors ([Ensemble Variance Calibration](#), *Kolczynsky et al., 2009, 2011*; *Bonavita et al., 2011*)



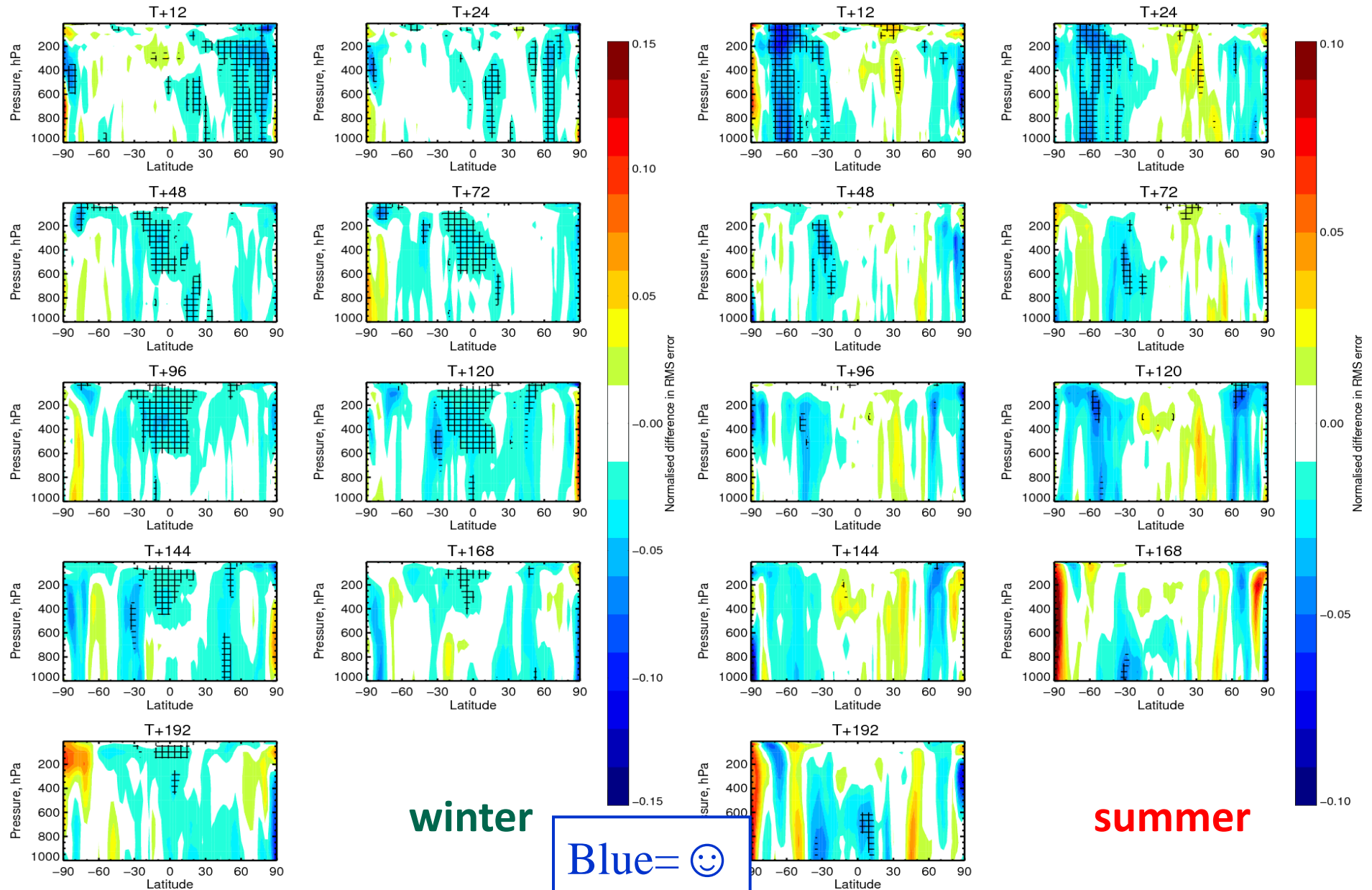
# It works well: Geopotential RMSE reduction

RMS forecast errors in Z(ffg8-fezj), 11-Jan-2010 to 30-Mar-2010, from 72 to 79 samples.

Point confidence 99.5% to give multiple-comparison adjusted confidence 90%. Verified against own-analysis.

RMS forecast errors in Z(ffge-0051), 2-Aug-2010 to 30-Oct-2010, from 83 to 90 samples.

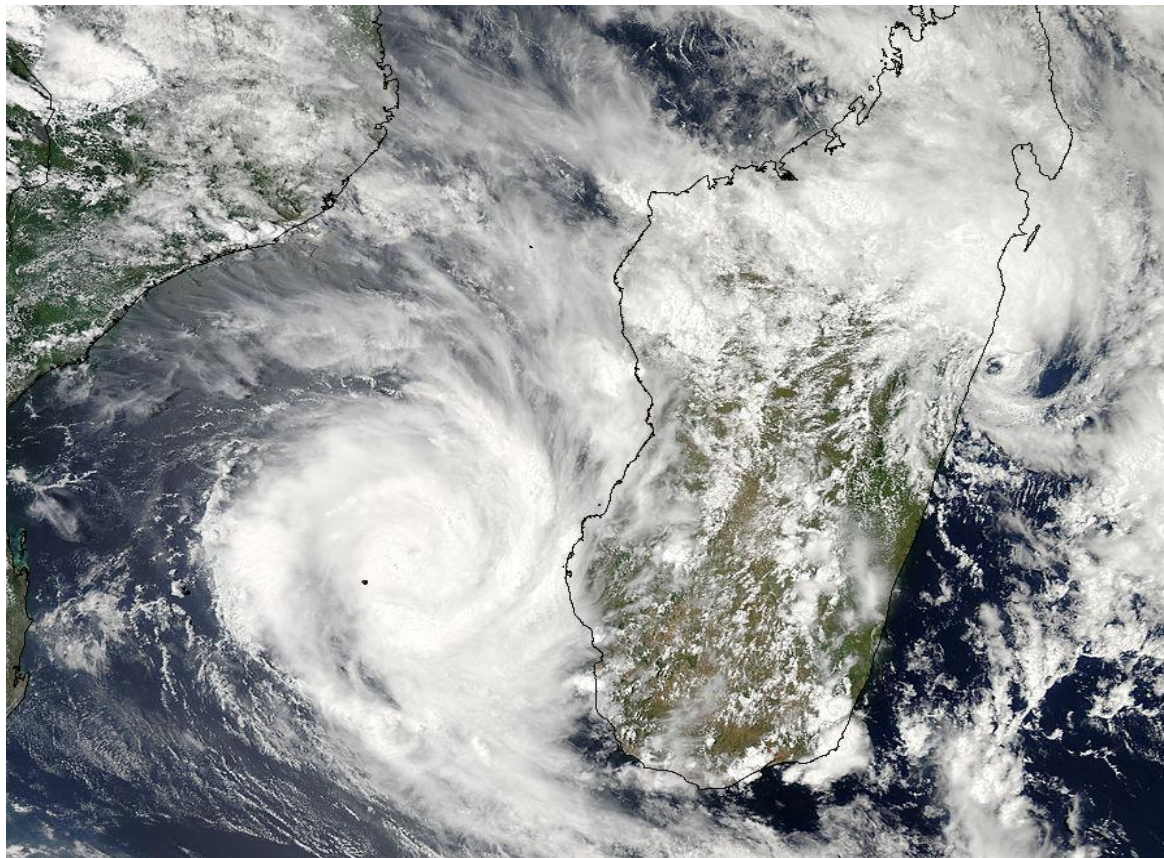
Point confidence 99.5% to give multiple-comparison adjusted confidence 90%. Verified against own-analysis.



# Hybrids next step: EDA Covariance estimation

## Diagnosing the Background Error Correlation Length-Scales

Hurricane Fanele, 20 January 2009



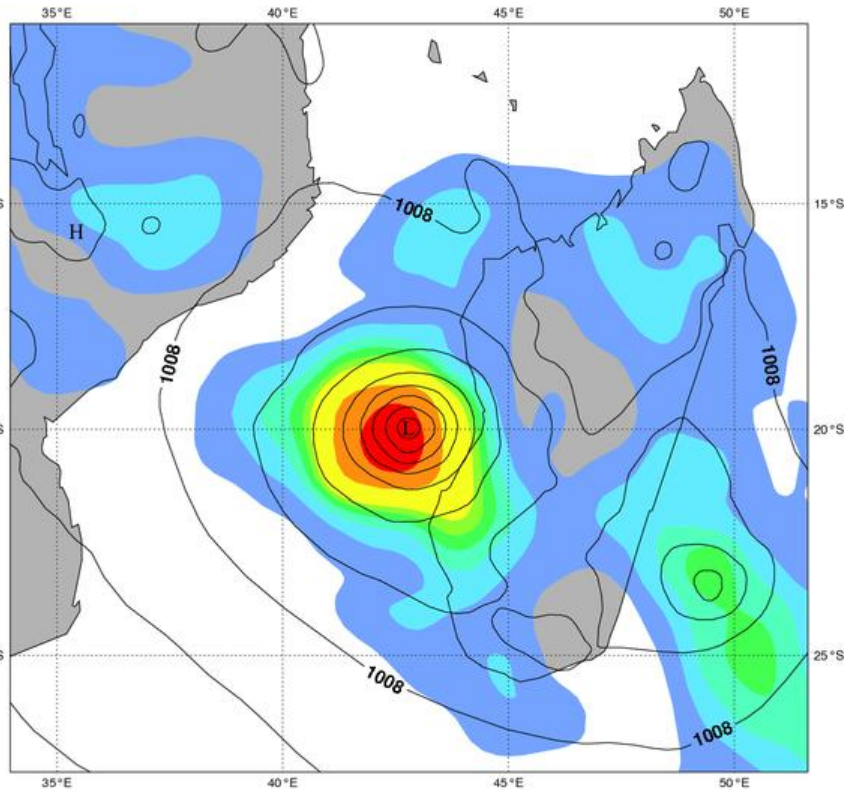


# Hybrids next step: EDA Covariance estimation

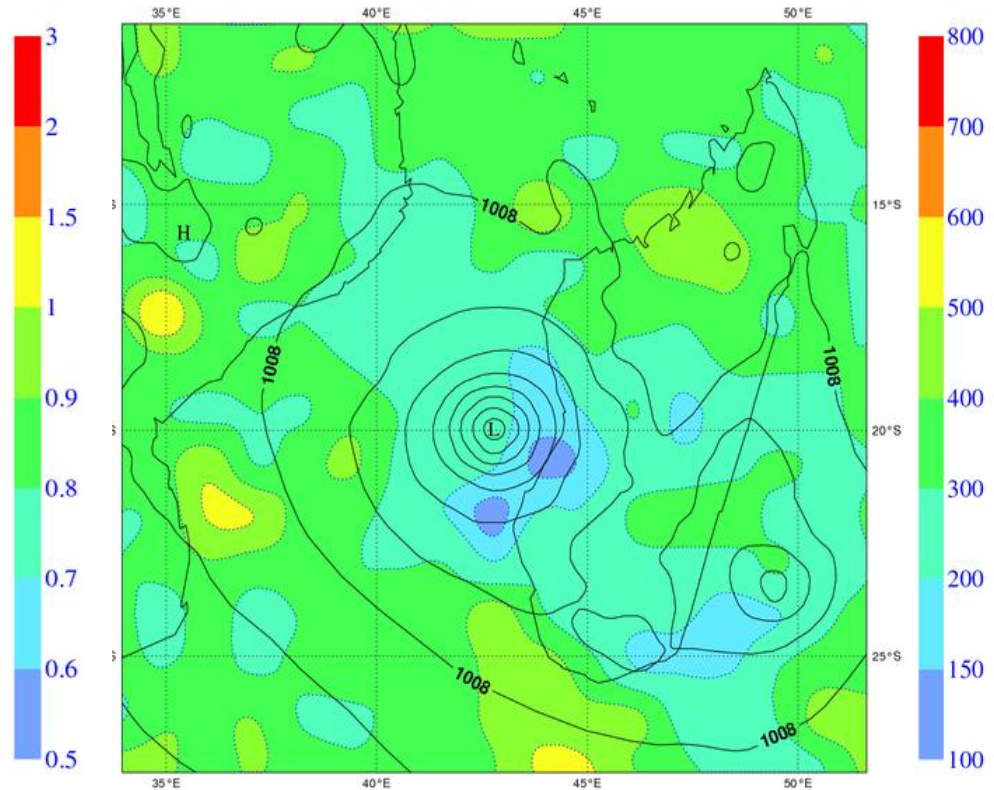
## 20 member EDA

Surf. Press. Background Err. St.Dev.      Surf. Press. BG Err. Correlation L. Scale

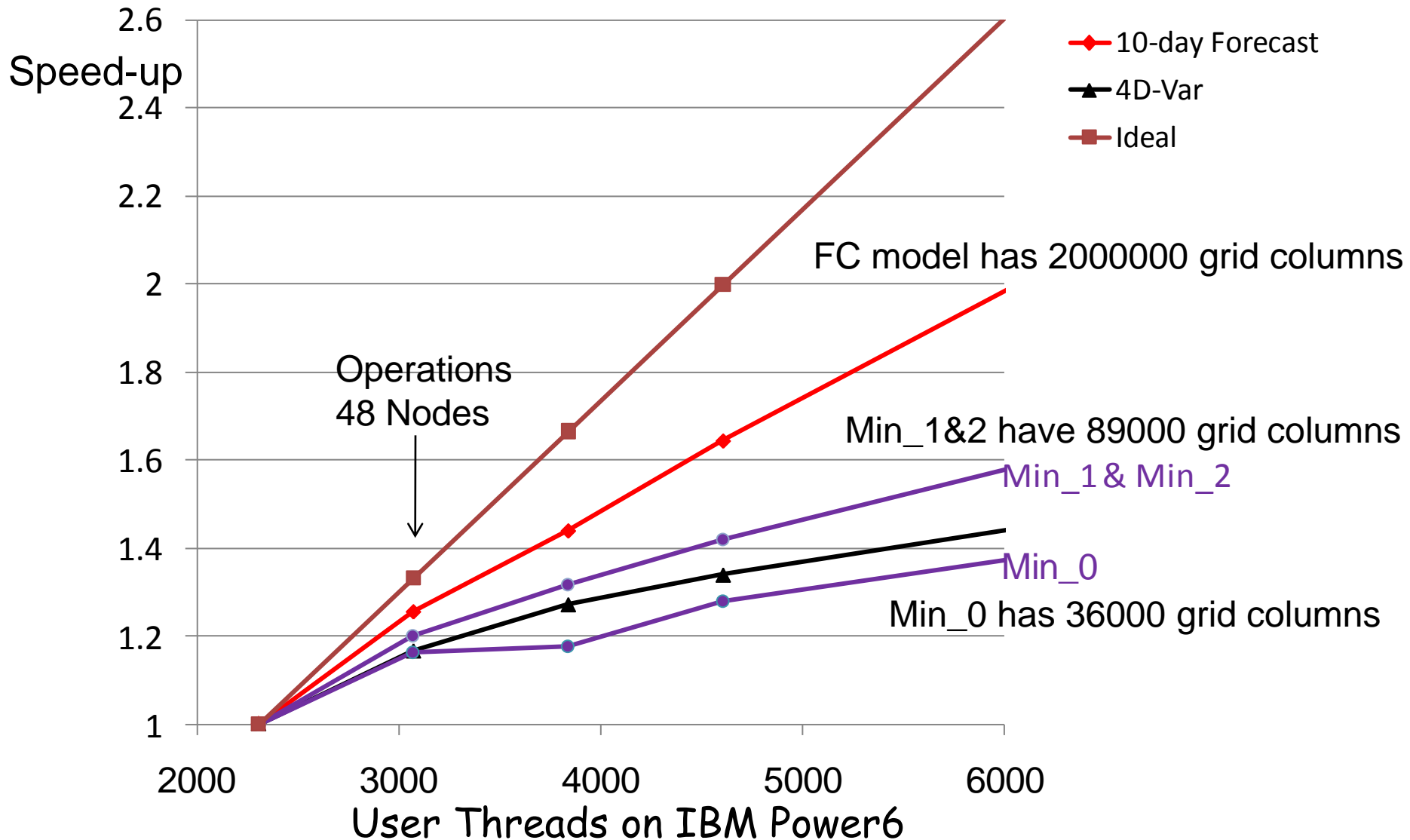
Tuesday 20 January 2009 00UTC ECMWF Forecast t+9 VT: Tuesday 20 January 2009 09UTC Surface: Mean sea level pressure



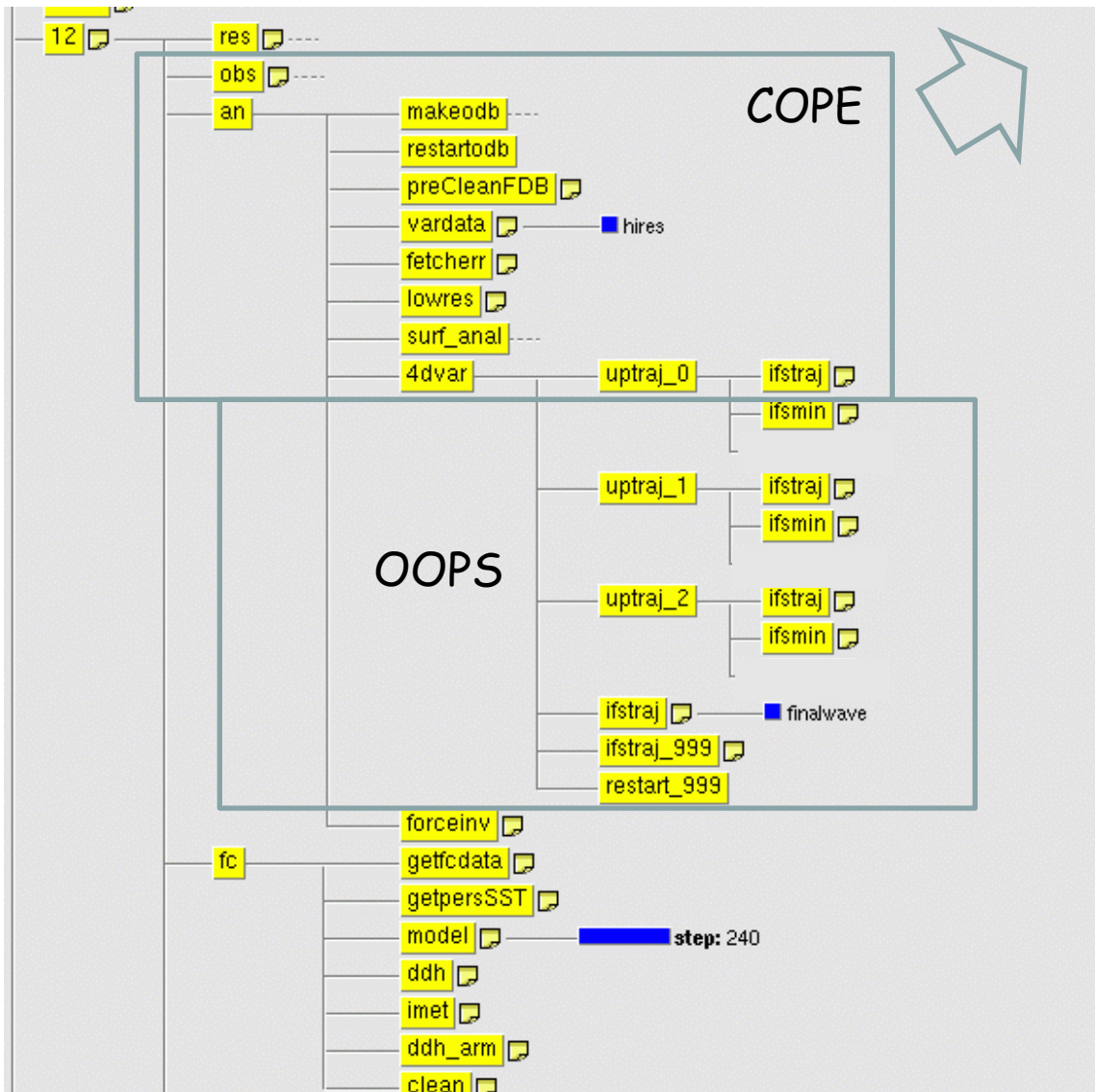
Tuesday 20 January 2009 00UTC ECMWF Forecast t+9 VT: Tuesday 20 January 2009 09UTC Surface: Mean sea level pressure



# Scalability of T1279 Forecast and 4D-Var



# Improving scalability of time critical DA suite



4D-var time window is 12 hours

**Forecast & Outer loop trajectories:**

(Traj\_0,1,2) are using T1279 resolution  
Grid columns =  $2 \times 10^6$

**Three minimizations:**

Min\_0 : T159  
Grid columns = 36000

Min\_1 & 2 : T255  
Grid columns = 89000

Vertical = 91 levels

# The 5 Dimensions of 4D-Var

- The bulk of the 4D-Var algorithm comprises 5 nested loop directions:
  - ① Minimisation algorithm iterations (inner and outer),
  - ② Time stepping of the model (and TL/AD),
  - ③ Latitude, **NPROMA**
  - ④ Longitude, **NPROMA**
  - ⑤ Vertical.
- Only **two** are parallel!
- We need to look at the **other directions** for more parallelism, for example:
  - ▶ Minimisation algorithm:
    - ★ Parallel search directions,
    - ★ Parallel preconditioner and less iterations,
    - ★ Observation space algorithms, saddle point algorithms.
  - ▶ Time stepping:
    - ★ Weak constraint 4D-Var.
- Scalability **cannot** be improved solely by technical or local optimizations!

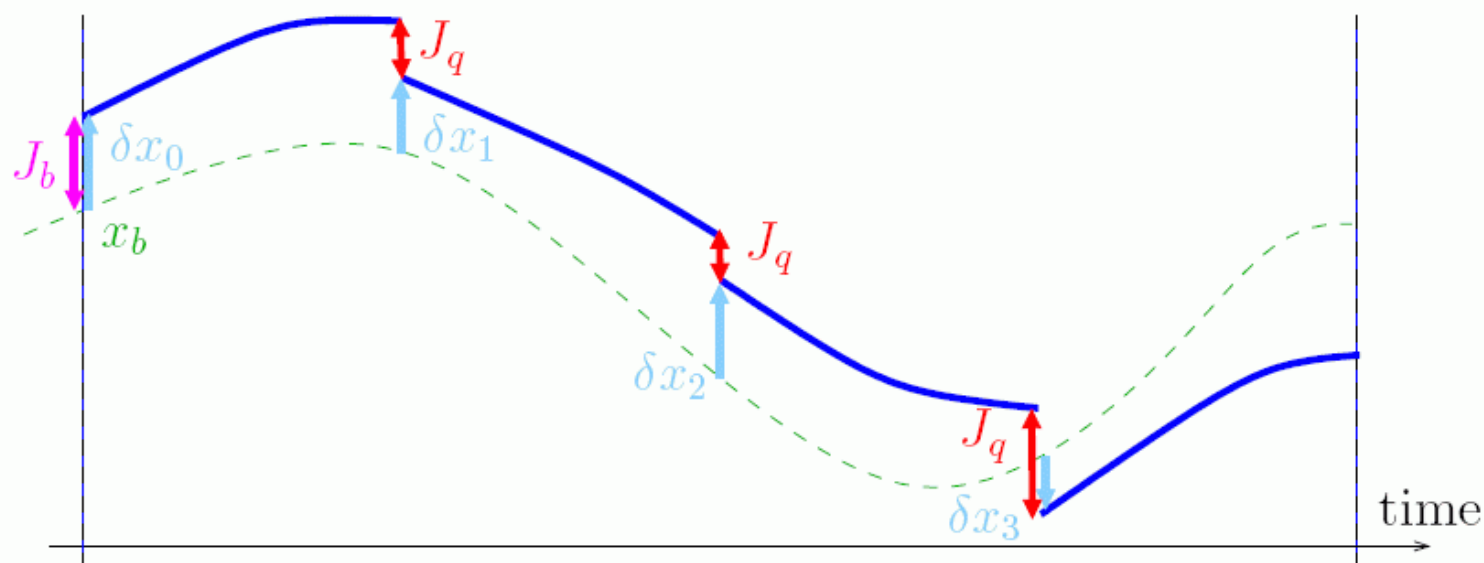
# Object-Oriented Prediction System – The OOPS project

- Data Assimilation algorithms manipulate a limited number of entities (objects):
  - $x$  (State),  $y$  (Observation),
  - $H$  (Observation operator),  $M$  (Model),  $H^*$  &  $M^*$  (Adjoint),
  - $B$  &  $R$  (Covariance matrices), etc.
  - To enable development of new data assimilation algorithms in IFS, these objects should be easily available & re-usable
- More Scalable Data Assimilation
- Cleaner, more Modular IFS

## OOPS → More Scalable Data Assimilation

- One execution instead of many will reduce start-up - also I/O between steps will not be necessary
- New more parallel minimisation schemes
  - Saddle-point formulation  
(Only OOPS has made it possible for Mike Fisher to implement the saddle-point formulation so quickly!!)
- For long-window, weak-constraint 4D-Var: Minimization steps for different sub-windows can run in parallel as part of same analysis

## Weak Constraint 4D-Var



- Model integrations within each time-step (or sub-window) are independent:
  - ▶ Information is not propagated across sub-windows by TL/AD models,
  - ▶  $\mathcal{M}$  and  $\mathcal{H}$  can be run in parallel over the sub-windows.
- Several shorter 4D-Var cycles are coupled and optimised together.
- 4D-Var becomes an elliptic problem and preconditioning becomes more complex.

# Conclusions and Perspectives

- ❑ Use of hybrids consistently **improves deterministic analysis and forecast skill** w.r.to pure sequential (EnKF) and non-sequential (4D-Var) solutions;
- ❑ EDA/EnKF, possibly re-centred around deterministic analysis, provide improved sampling of initial errors for **Ensemble Prediction**
- ❑ We can expect **growing ensemble use in 4D-Var**:
  1. A larger ensemble (both in the EDA and EnKF) improves error characterization and ultimately skill scores;
  2. 4D background error covariances sampled from an EDA/EnKF could be used over the all 4D-Var assimilation window (not only at the start!): **En-4D-Var** (Liu et al., 2008; Buehner et al., 2010). This would remove the need of developing and maintaining a TL and Adjoint version of the forecast model



# Conclusions and Perspectives

- We can expect **growing ensemble use in 4D-Var**:
  3. Weak-constraint Long-window 4D-Var revolves around the estimation of  $Q$ : It is conceivable that an EDA will provide a way of effectively sampling  $Q$
  4. The EnKF is more computationally efficient than an ensemble of 4D-Var analysis (EDA): **if** it can be shown to be as accurate as standard 4D-Var with the full observing system, then it will provide a relatively cheap and efficient way of cycling error estimates in a hybrid system

# Summary of ECMWF's Data Assimilation strategy

- Hybrid DA system: Use EDA information in 4D-Var

Flow dependent background error variances and covariances, and model error in 4D-Var

Provides improved uncertainty estimation

- Long-window weak-constraint 4D-Var

- Unified Ensemble of Data Assimilations (EDA) and Ensemble Prediction System

For estimation of analysis and short range forecast uncertainty that will benefit the deterministic 4D-Var

For estimation of long range forecast uncertainty (the present role of the EPS)

Note: The EDA is a 'stochastic EnKF' with an expensive 4D-Var component. It may be replaced or supplemented by an LETKF system, if beneficial.