

Detection and Correction of Model Bias During Data Assimilation

Dick Dee
ECMWF, on leave from
NASA Global Modeling and Assimilation Office (**GMAO**)

ECMWF Seminar
September 8, 2003

Outline

- Formulation of the analysis problem
- Manifestations of model bias during data assimilation
- Sequential schemes for estimating persistent biases
- Mixed results with a global atmospheric assimilation system
- Sequential estimation of non-persistent biases
- Simple experiments with a global ocean assimilation system

Acknowledgements

- Arlindo da Silva, Ricardo Todling, Banglin Zhang (NASA-GMAO)
- Tony McNally, Per Kållberg, Thomas Jung, Antonio Garcia Mendez (ECMWF)
- Steven Pawson (NASA-GMAO)
- Gennady Chepurin, Jim Carton (University of Maryland)

The Analysis Problem

$$\text{Maximize } p(\mathbf{x}|\mathbf{y}, \mathbf{x}^f) \text{ where } \begin{cases} \mathbf{y} = \mathbf{h}(\mathbf{x}) + \mathbf{e}^o & \text{(observations)} \\ \mathbf{x}^f = \mathbf{x} + \mathbf{e}^f & \text{(first guess)} \end{cases}$$

In case of Gaussian errors in the first guess and observations:

$$\text{Minimize } J(\mathbf{x}) = \frac{1}{2}(\mathbf{x}^f - \mathbf{x})^T \mathbf{P}^{-1}(\mathbf{x}^f - \mathbf{x}) + \frac{1}{2}(\mathbf{y} - \mathbf{h}(\mathbf{x}))^T \mathbf{R}^{-1}(\mathbf{y} - \mathbf{h}(\mathbf{x}))$$

$$\text{where } \begin{cases} \mathbf{P} = \langle \mathbf{e}^f (\mathbf{e}^f)^T \rangle & \text{(first-guess error covariances)} \\ \mathbf{R} = \langle \mathbf{e}^o (\mathbf{e}^o)^T \rangle & \text{(observation error covariances)} \end{cases}$$

- assumptions: no bias, no correlations between first-guess and observation errors
- the evaluation of $\mathbf{h}(\mathbf{x})$ may involve time integration
- data assimilation requires cycling (in time), and updating of \mathbf{P}
- the first guess can be thought of as a special set of observations

The Nonlinear Analysis Equation

The gradient is $\nabla_{\mathbf{x}} J = - \left(\frac{\partial \mathbf{h}}{\partial \mathbf{x}} \right)^T \mathbf{R}^{-1} (\mathbf{y} - \mathbf{h}(\mathbf{x})) - \mathbf{P}^{-1} (\mathbf{x}^f - \mathbf{x})$

so the minimizing solution $\mathbf{x} = \mathbf{x}^a$ satisfies

$$\mathbf{x}^a - \mathbf{x}^f = \mathbf{P} \left(\frac{\partial \mathbf{h}}{\partial \mathbf{x}} \bigg|_{\mathbf{x}=\mathbf{x}^a} \right)^T \mathbf{R}^{-1} (\mathbf{y} - \mathbf{h}(\mathbf{x}^a))$$

This equation is nonlinear, but in any case $\mathbf{x}^a - \mathbf{x}^f \in \text{Range}(\mathbf{P})$

Therefore, the background error covariances constrain the solution:

- all information not already in the background is filtered by \mathbf{P}
- must have a 'rich' and meaningful background error covariance model
- for this reason, the choice of control variable is absolutely critical

Solution Algorithms

$$\mathbf{x}^a - \mathbf{x}^f = \mathbf{P} \left(\frac{\partial \mathbf{h}}{\partial \mathbf{x}} \Big|_{\mathbf{x}=\mathbf{x}^a} \right)^T \mathbf{R}^{-1} (\mathbf{y} - \mathbf{h}(\mathbf{x}^a))$$

All practical analysis algorithms are based on repeated linearization of the observation operator (possibly including time integration), as in

$$\mathbf{h}(\mathbf{x}^a) = \mathbf{h}(\mathbf{x}^f) + \mathbf{H}(\mathbf{x}^a - \mathbf{x}^f) \quad \text{where} \quad \mathbf{H} = \frac{\partial \mathbf{h}}{\partial \mathbf{x}} \Big|_{\mathbf{x}=\mathbf{x}^f}$$

Substitution in the nonlinear analysis equation then gives

$$[\mathbf{I} + \mathbf{P}\mathbf{H}^T\mathbf{R}^{-1}\mathbf{H}] (\mathbf{x}^a - \mathbf{x}^f) = \mathbf{P}\mathbf{H}^T\mathbf{R}^{-1}(\mathbf{y} - \mathbf{h}(\mathbf{x}^f))$$

- involves solution of a linear system of equations in state space
- basis for implementations at NCEP (SSI), ECMWF (4DVAR)

Alternatively, by the Sherman-Morrison-Woodbury formula,

$$\mathbf{x}^a - \mathbf{x}^f = \mathbf{P}\mathbf{H}^T [\mathbf{H}\mathbf{P}\mathbf{H}^T + \mathbf{R}]^{-1} (\mathbf{y} - \mathbf{h}(\mathbf{x}^f))$$

- involves solution of a linear system of equations in observation space
- basis for implementations at DAO (PSAS), NRL (NAVDAS)

Systematic Errors and Data Assimilation

The (standard) analysis is of the form

$$\mathbf{x}^a - \mathbf{x}^f = \mathbf{K} [\mathbf{y} - \mathbf{h}(\mathbf{x}^f)]$$

with \mathbf{K} some close-to-linear operator, e.g., $\mathbf{K} = \mathbf{P}\mathbf{H}^T [\mathbf{H}\mathbf{P}\mathbf{H}^T + \mathbf{R}]^{-1}$

Errors are defined as

$$\mathbf{e}^a = \mathbf{x}^a - \mathbf{x}, \quad \mathbf{e}^f = \mathbf{x}^f - \mathbf{x}, \quad \mathbf{e}^o = \mathbf{y} - \mathbf{h}(\mathbf{x})$$

Systematic errors show up in the mean observed-minus-guess residuals:

$$\langle \mathbf{y} - \mathbf{h}(\mathbf{x}^f) \rangle \approx \langle \mathbf{e}^o \rangle - \mathbf{H} \langle \mathbf{e}^f \rangle$$

in the mean analysis increments:

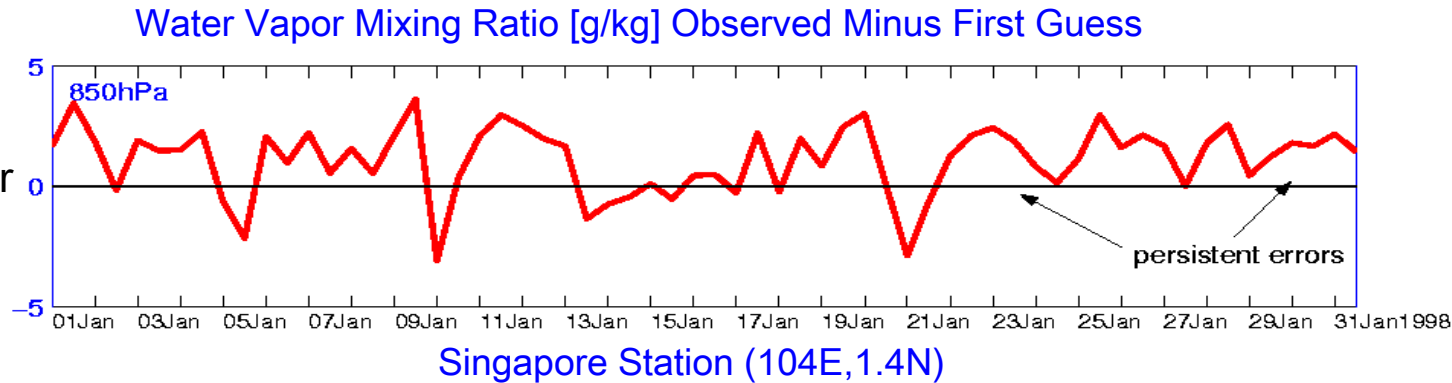
$$\langle \mathbf{x}^a - \mathbf{x}^f \rangle \approx \mathbf{K} \langle \mathbf{e}^o \rangle - \mathbf{K}\mathbf{H} \langle \mathbf{e}^f \rangle$$

and, ultimately, in the analysis itself:

$$\langle \mathbf{e}^a \rangle \approx \mathbf{K} \langle \mathbf{e}^o \rangle + [\mathbf{I} - \mathbf{K}\mathbf{H}] \langle \mathbf{e}^f \rangle$$

Humidity Observed-Minus-First Guess (GEOS 2)

Station data show persistent dry bias in the tropical lower troposphere...

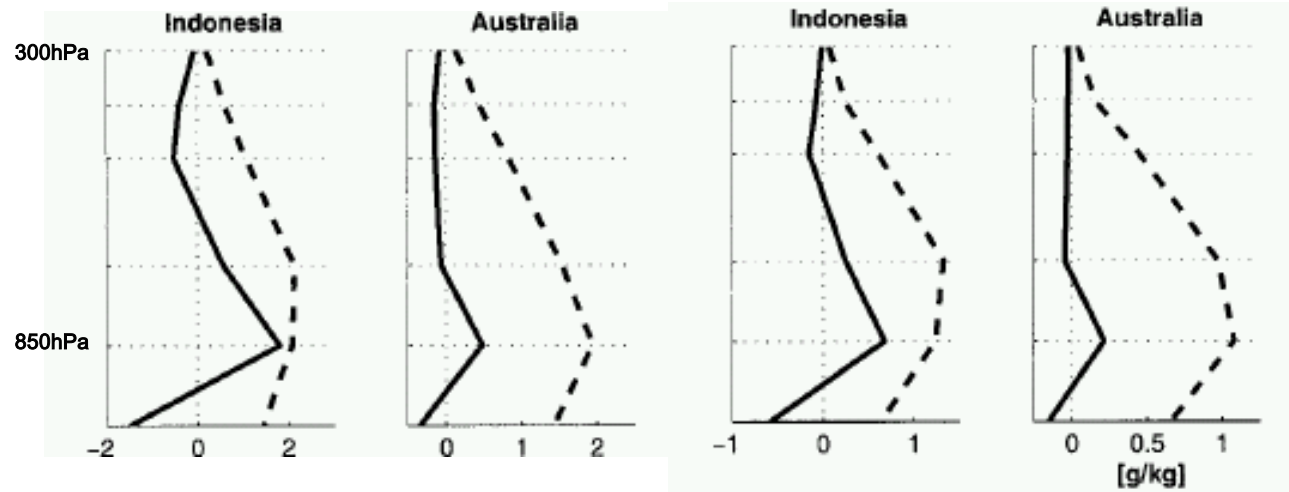


$$\langle y - h(x^f) \rangle \approx \langle e^o \rangle - H \langle e^f \rangle$$

Jan 1998 mean (solid) and std dev (dashed)

Observed-minus-forecast

Observed-minus-analyzed

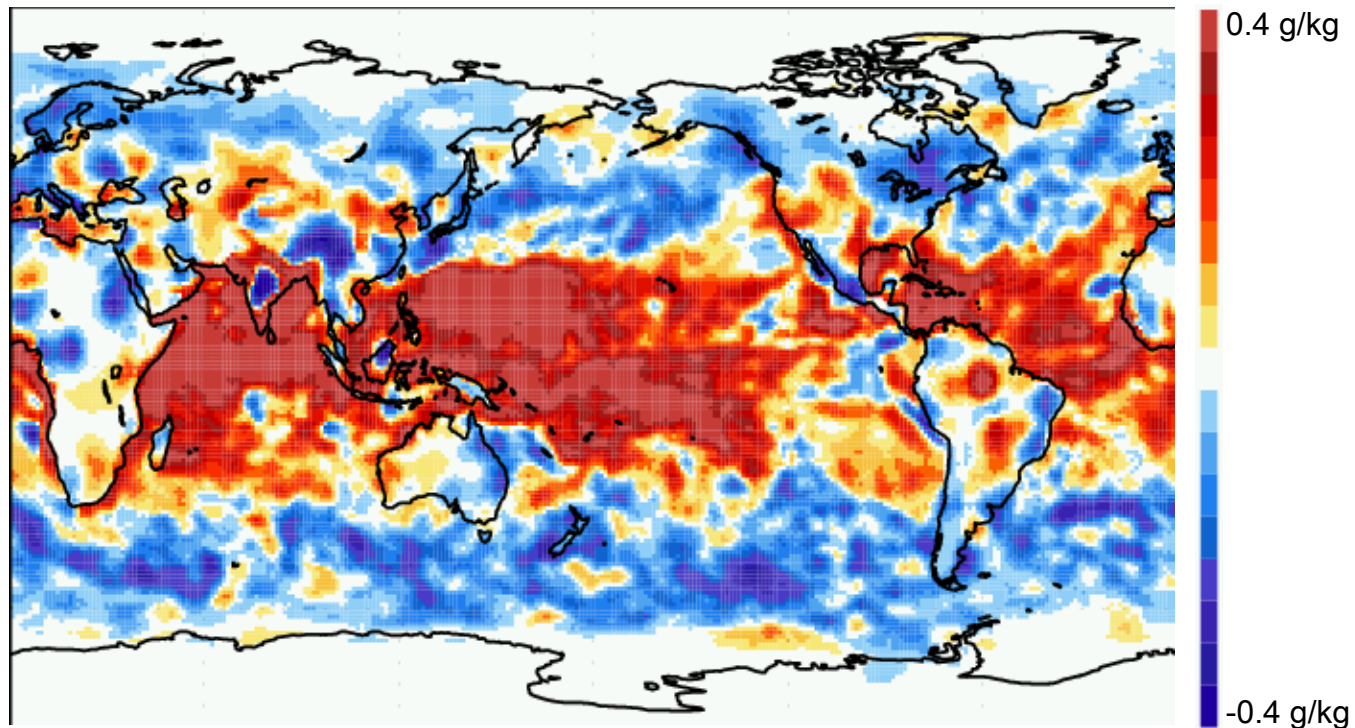


...and this clearly shows up in the monthly statistics of the data residuals for tropical stations

Non-Zero Mean Humidity Increments (fvDAS)

$$\langle x^a - x^f \rangle \approx K \langle e^o \rangle - KH \langle e^f \rangle$$

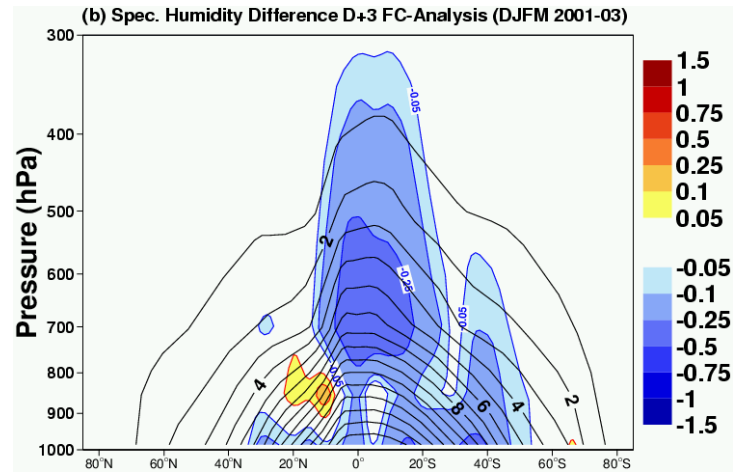
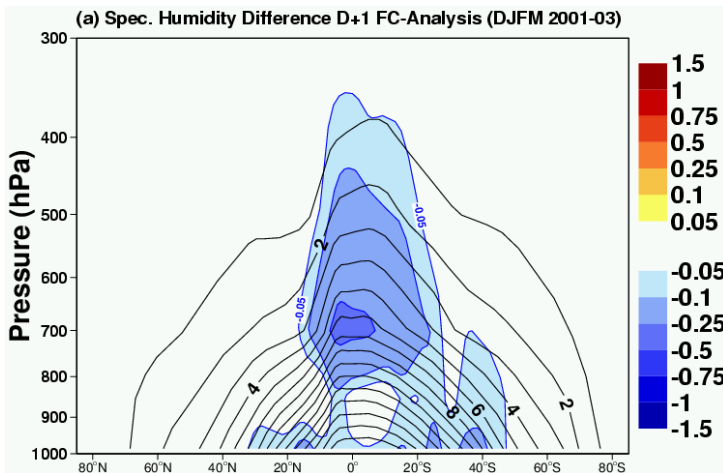
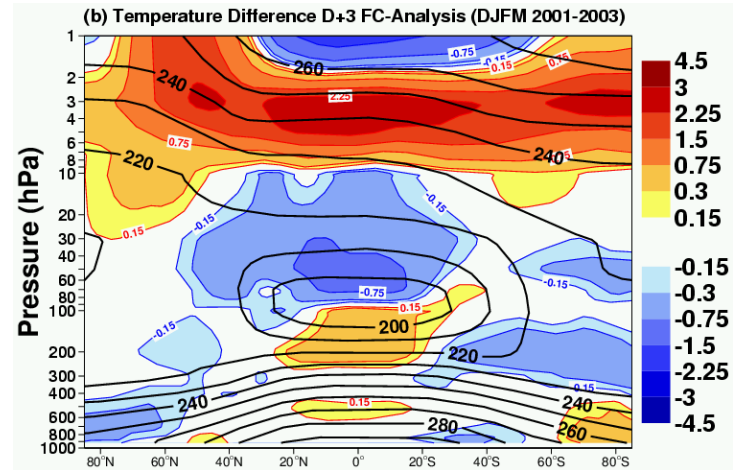
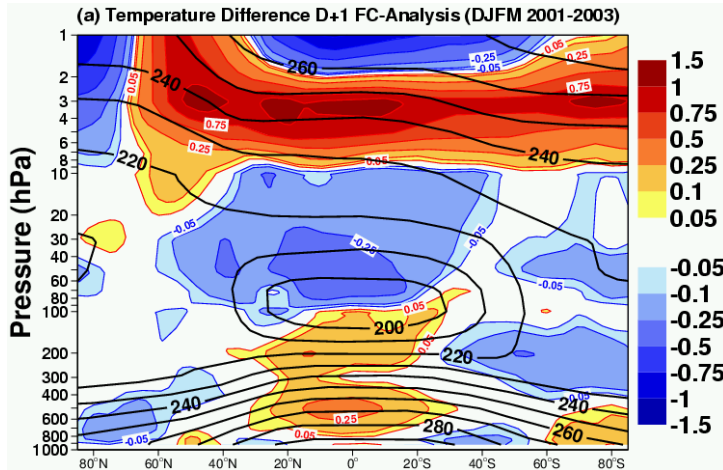
Jan2002 mean analysis increment for specific humidity in fvDAS, layer 4 (~850hPa)



Based on rawinsonde, TOVS, and SSM/I (TPW) observations

Systematic Errors in the ECMWF Operational System

Zonal mean forecast errors (verified against analyses) reflect persistent model bias



(Thanks to Thomas Jung)

Stratospheric Temperature Bias

$$\langle y - h(x^f) \rangle \approx \langle e^o \rangle - H \langle e^f \rangle$$

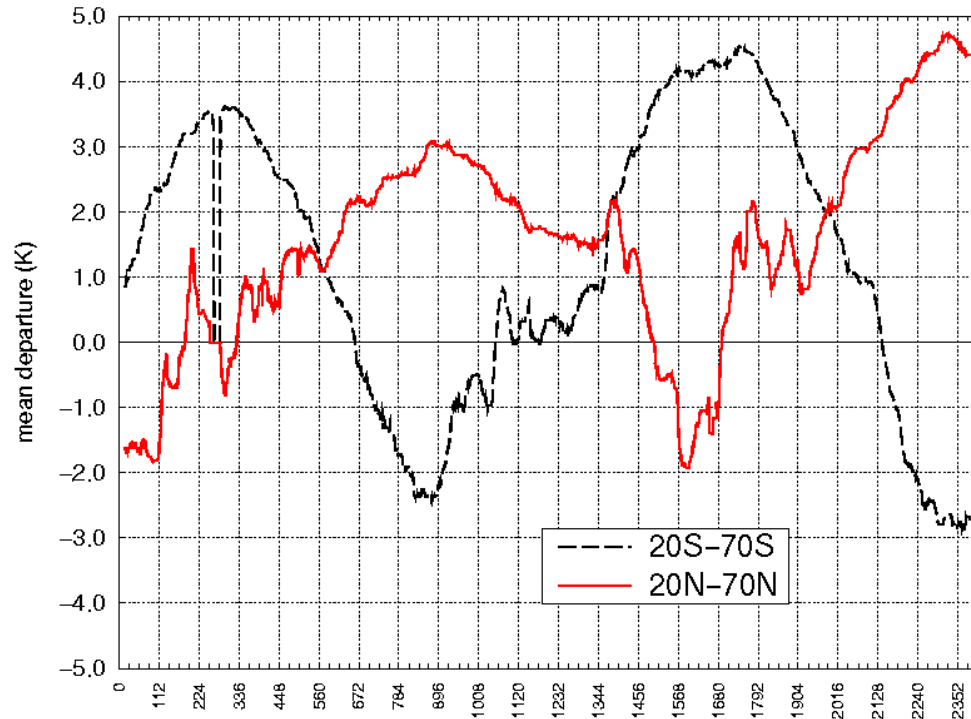
Very large mean temperature residuals in AMSU channel 14

This channel is most sensitive in the upper stratosphere

Comparison with HALOE and LIDAR data shows this is due to model bias

AMSUA-14 mean Obs-Fg departure

1998-11-01 to 2000-06-26

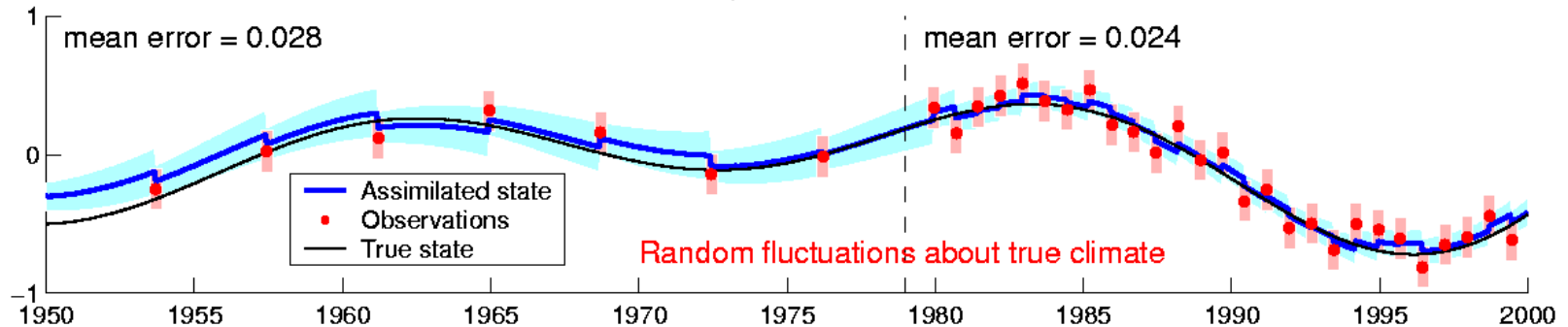


(Thanks to Tony McNally)

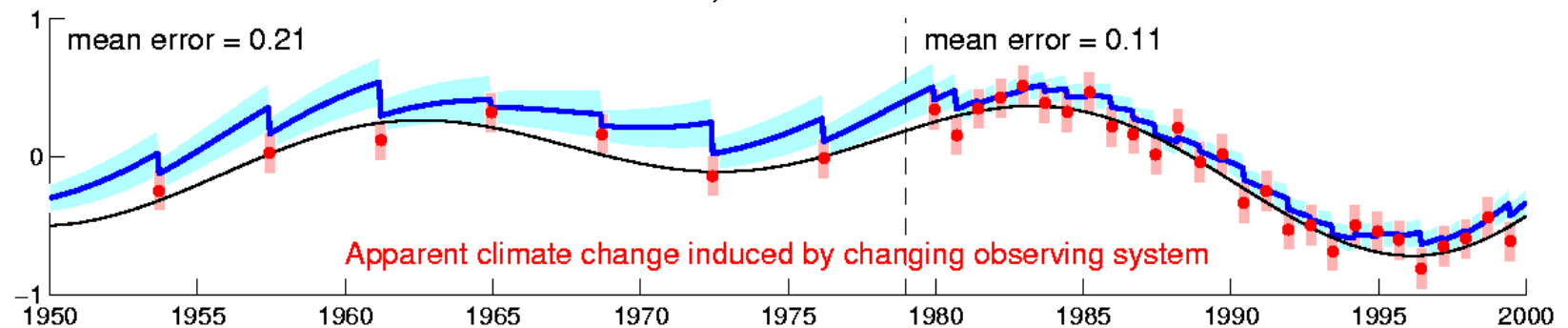
Impact of Model Bias on Climate Parameters

$$\langle e^a \rangle \approx \mathbf{K} \langle e^o \rangle + [\mathbf{I} - \mathbf{K}\mathbf{H}] \langle e^f \rangle$$

Unbiased model, unbiased observations

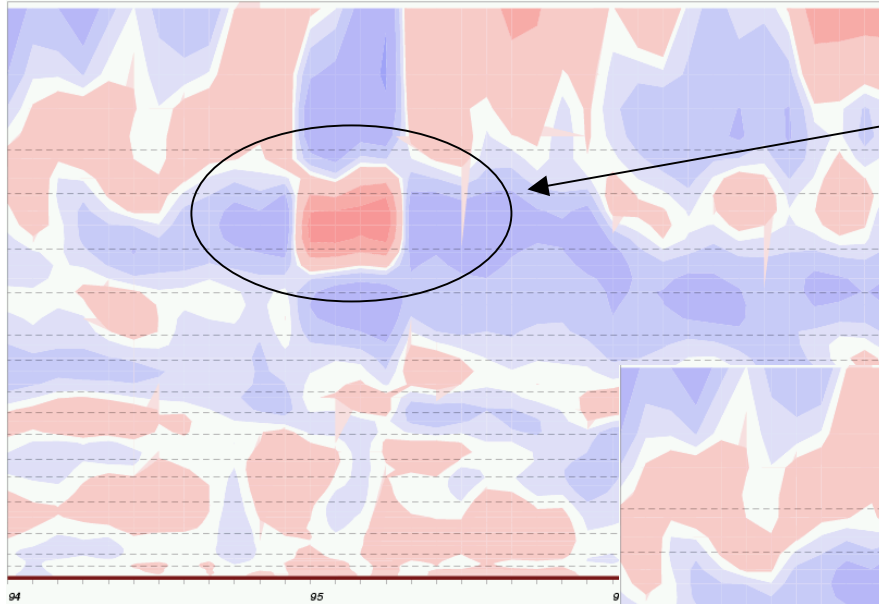


Biased model, unbiased observations



Model Warm Bias in ERA-40 (analyses)

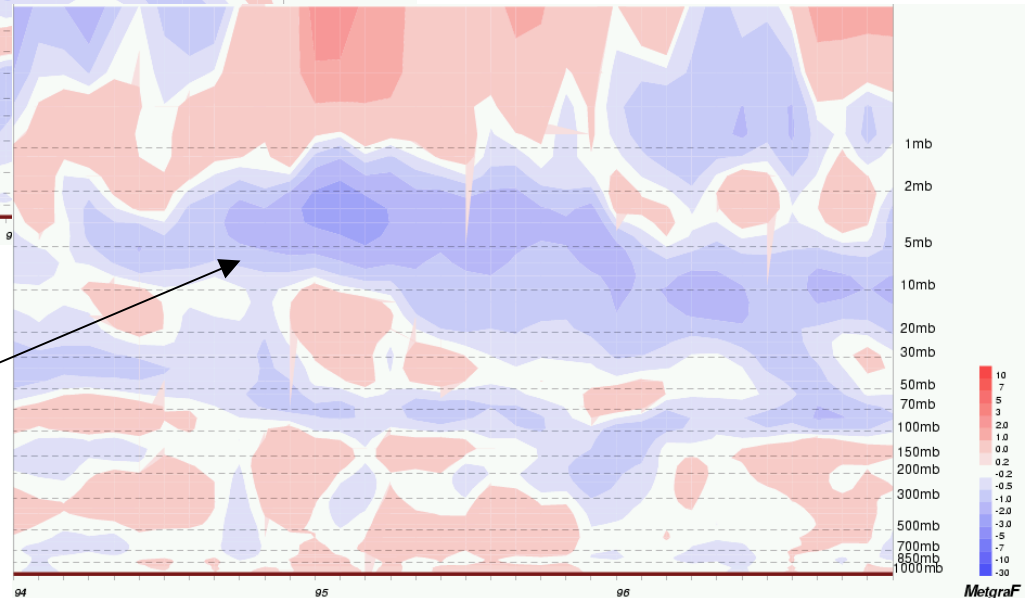
anomaly to ERA-40 10-year monthly mean. region: Global 090S- 090N 000E- 360E
Temperature differences in C.



Preliminary assimilation without SSU between Jan-Apr 1995

$$\langle e^a \rangle \approx \mathbf{K} \langle e^o \rangle + [\mathbf{I} - \mathbf{KH}] \langle e^f \rangle$$

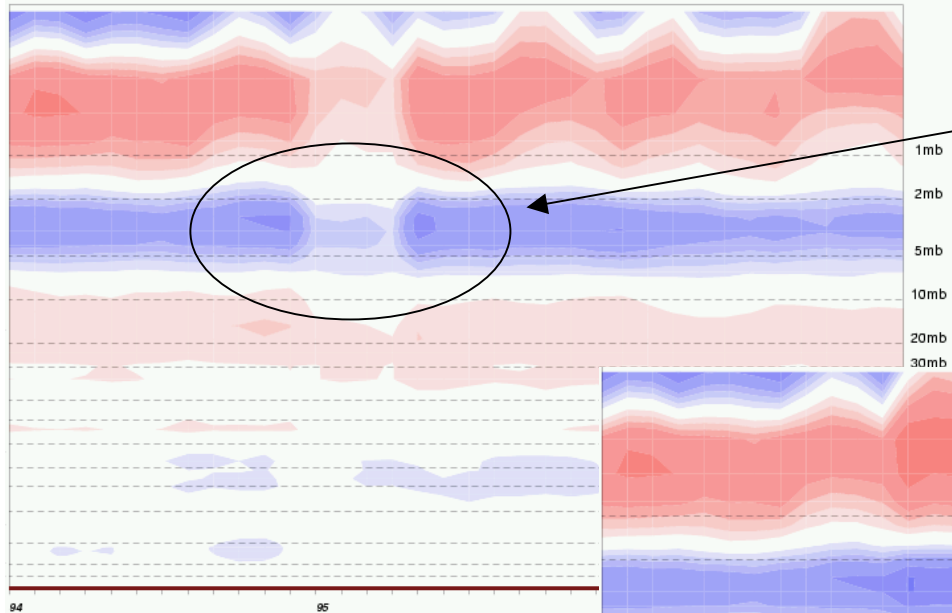
Final ERA-40 production (including SSU)



(Thanks to Per Kållberg)

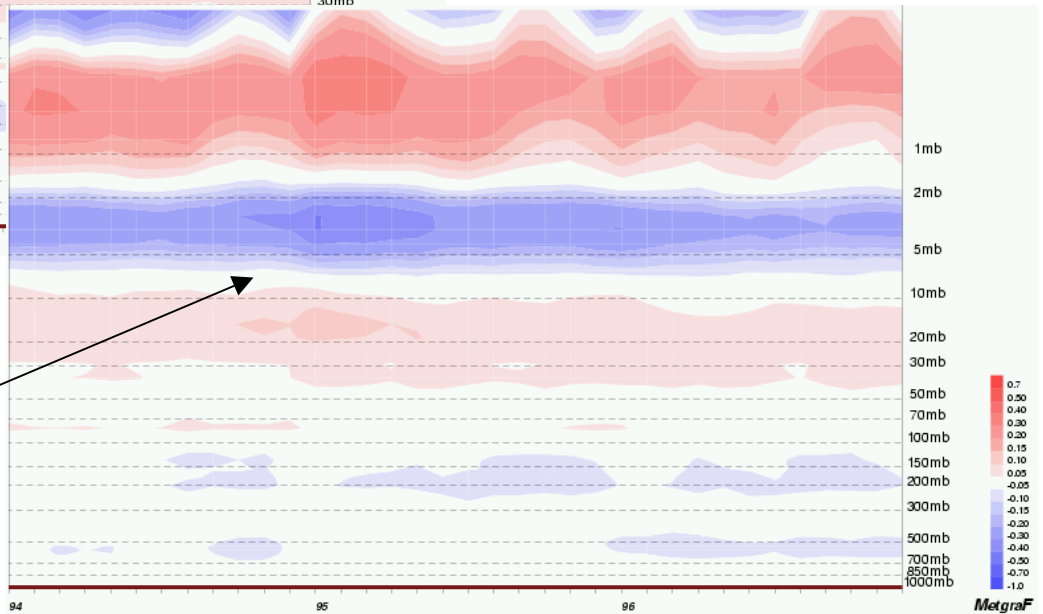
Model Warm Bias in ERA-40 (increments)

analysis increments with 0018 from monthly means region: Global 090S-090N 000E-360E
Temperature increments in C.



Preliminary assimilation with passive SSU from Jan-Apr 1995

$$\langle x^a - x^f \rangle \approx K \langle e^o \rangle - KH \langle e^f \rangle$$

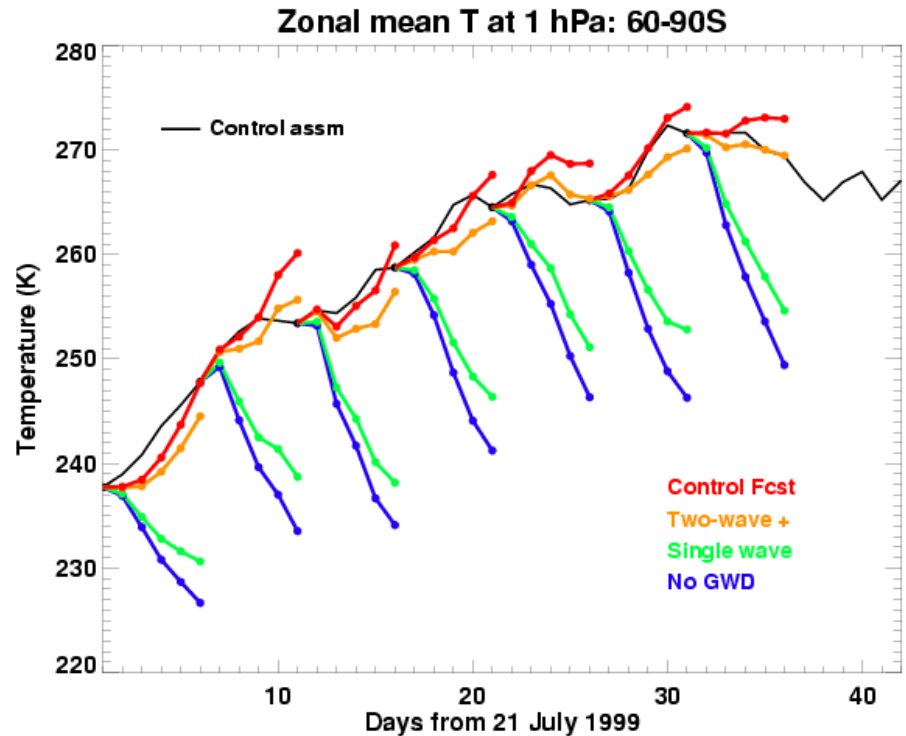


Final ERA-40 production (including SSU)

(Thanks to Per Kållberg)

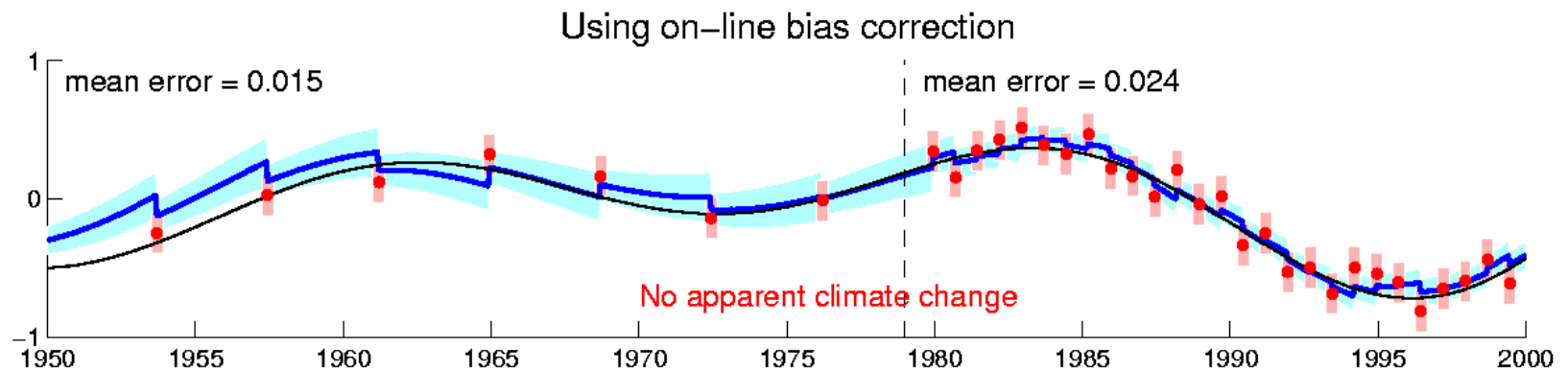
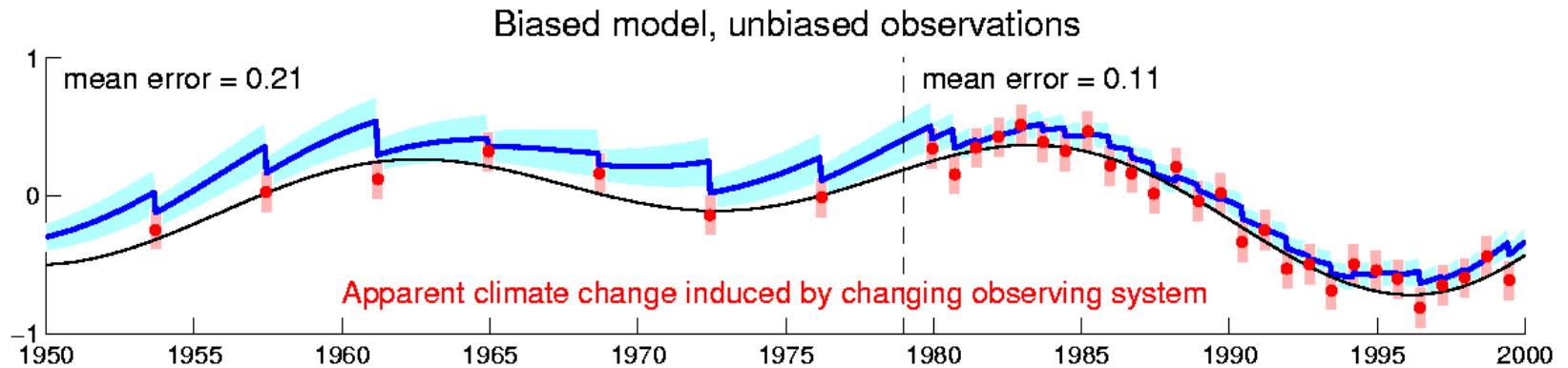
Stratospheric Bias Due To Gravity-Wave Drag

- Model climate is very sensitive to gravity-wave-drag scheme
- Stratospheric bias develops very quickly during the integration
- Clearly we wish to remove this bias by improving the model, but in the mean time...
- ... should the analysis simply assume that there is no bias?



(Thanks to Steven Pawson)

Correction of Persistent Bias in the First Guess (1)



Correction of Persistent Bias in the First Guess (2)

Standard assumptions include $\langle \mathbf{e}^f \rangle = 0$ (*no model bias*)

Instead, assume that $\mathbf{e}^f = \mathbf{b} + \tilde{\mathbf{e}}^f$, $\langle \tilde{\mathbf{e}}^f \rangle = 0$

with \mathbf{b} constant (or slowly varying) in time (*persistent model bias*)

Then the following algorithm produces bias estimates and unbiased analyses:

$$\begin{aligned}\hat{\mathbf{b}}_k &= \hat{\mathbf{b}}_{k-1} - \mathbf{L} \left[\mathbf{y}_k - \mathbf{H}(\mathbf{x}_k^f - \hat{\mathbf{b}}_{k-1}) \right] \\ \mathbf{x}_k^a &= (\mathbf{x}_k^f - \hat{\mathbf{b}}_k) + \mathbf{K} \left[\mathbf{y}_k - \mathbf{H}(\mathbf{x}_k^f - \hat{\mathbf{b}}_k) \right]\end{aligned}$$

- the standard analysis equation corresponds to $\hat{\mathbf{b}} = 0$
- the first equation produces an ‘analysis’ of the bias, given a ‘forecast’ of the bias
- for this algorithm, the bias estimation costs as much as the standard analysis
- one could use only selected observations for bias estimation
- this scheme is related to ‘separate bias estimation’ (Friedland 1969)

Covariance Modeling for the Bias Estimation

It can be shown that the unbiased analysis equations are optimal when

$$\mathbf{K} = \mathbf{P}^f \mathbf{H}^T [\mathbf{H} \mathbf{P}^f \mathbf{H}^T + \mathbf{R}]^{-1}$$
$$\mathbf{L} = \mathbf{P}^b \mathbf{H}^T [\mathbf{H} \mathbf{P}^b \mathbf{H}^T + \mathbf{H} \mathbf{P}^f \mathbf{H}^T + \mathbf{R}]^{-1}$$

where

$$\mathbf{P}^f = \langle \tilde{\mathbf{e}}^f (\tilde{\mathbf{e}}^f)^T \rangle \quad \text{unbiased first-guess error covariances}$$

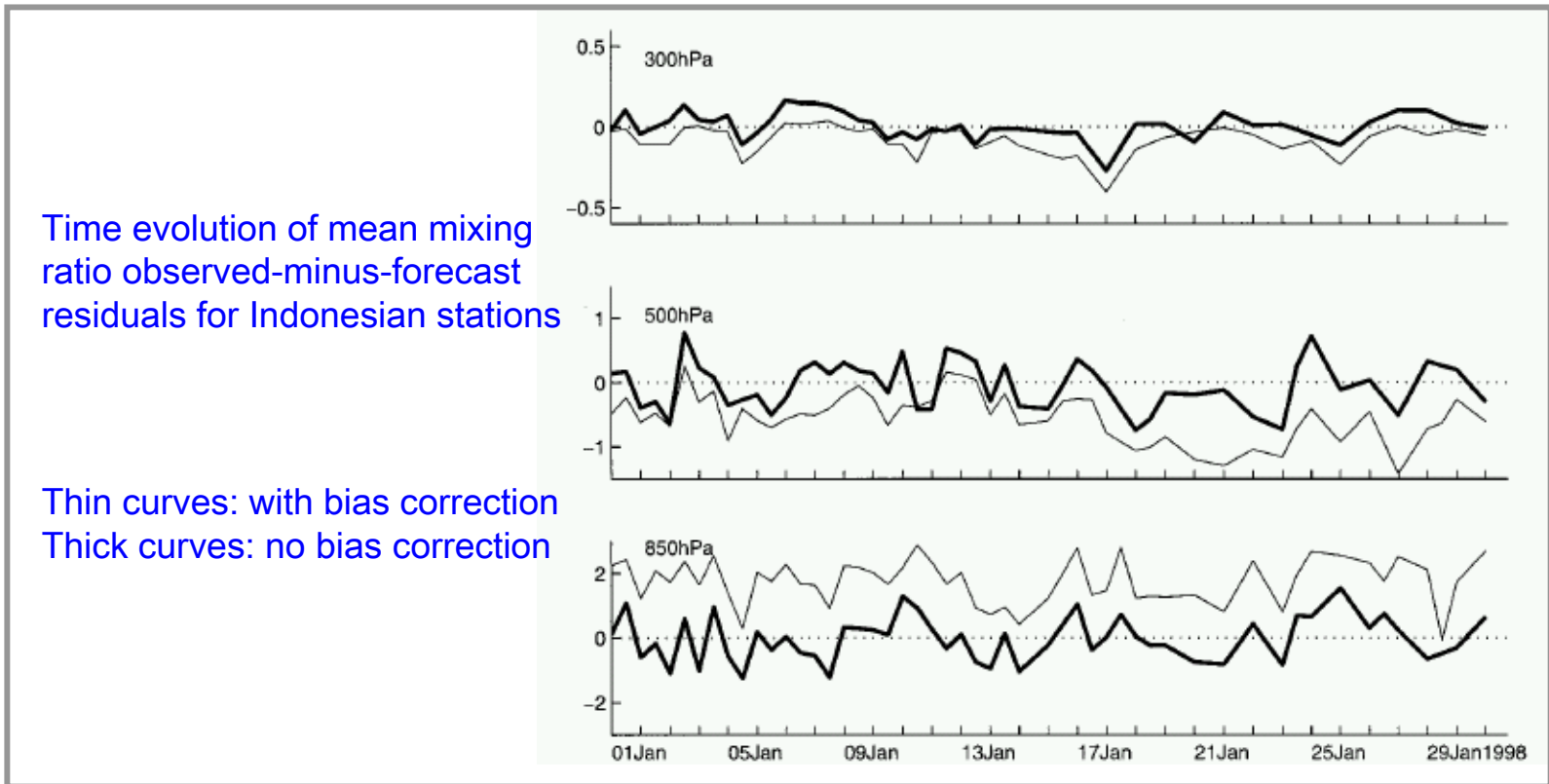
$$\mathbf{P}^b = \langle (\hat{\mathbf{b}} - \mathbf{b})(\hat{\mathbf{b}} - \mathbf{b})^T \rangle \quad \text{bias estimate error covariances}$$

In practice, \mathbf{P}^b must be modeled – e.g., $\mathbf{P}^b = \alpha \mathbf{P}^f$

- α is related to the time scale upon which the ‘persistent’ bias is allowed to vary
- the bias estimation does not have to be optimal
- however, the algorithm is very sensitive to problems with the multivariate balance implied by the covariance models
- some parameters in \mathbf{P}^b can be tuned from data

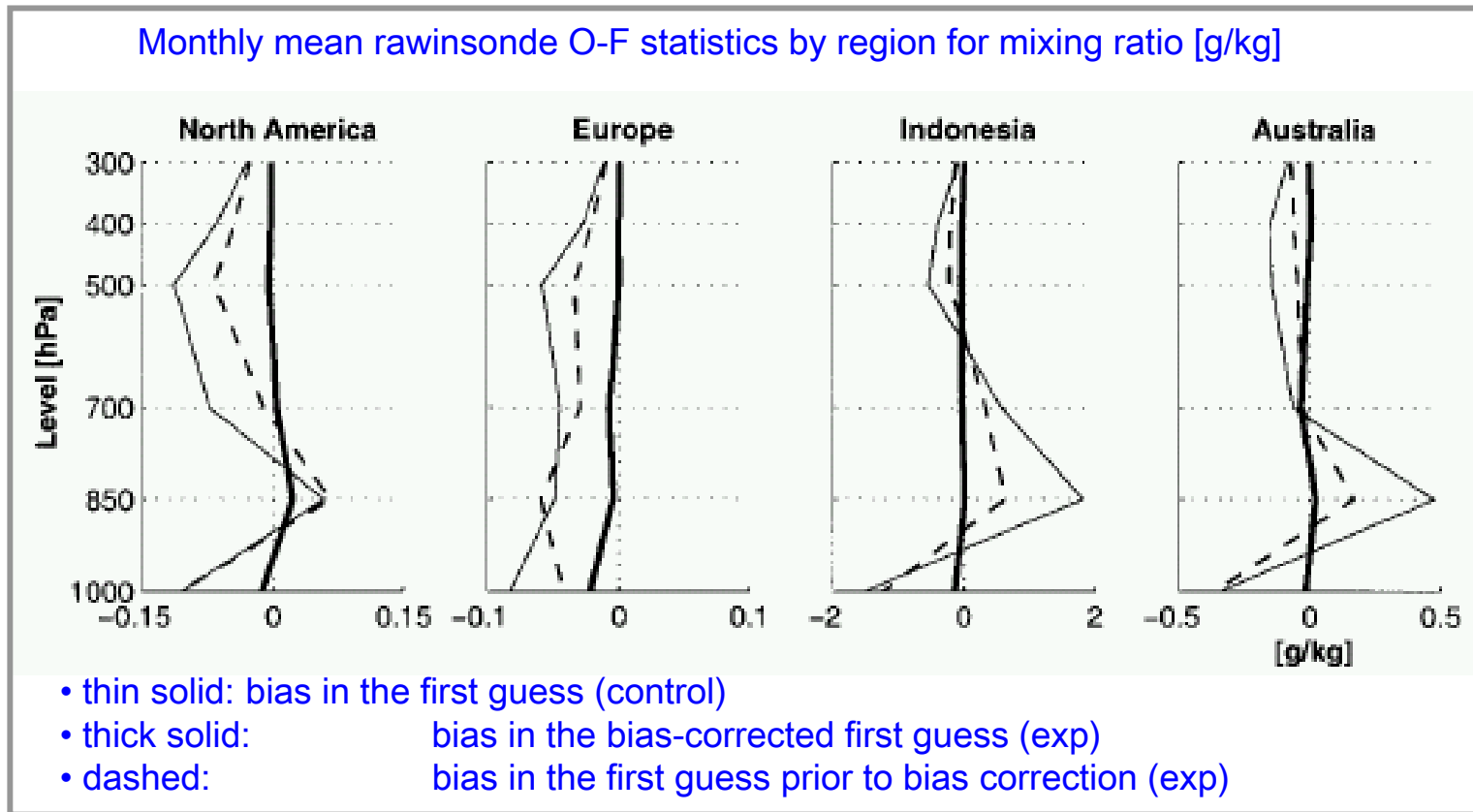
Moisture Bias Correction in GEOS-2

Model bias correction was implemented in the GEOS-2 data assimilation system at NASA Goddard Space Flight Center, based on rawinsonde observations only.



Reduction in First-guess Moisture Bias

Reduction of the first-guess bias suggests that the analyses have been improved: **Better analyses should lead to better forecasts.**



Simplification of the Algorithm (1)

The original algorithm is

$$\hat{\mathbf{b}}_k = \hat{\mathbf{b}}_{k-1} - \mathbf{L} \left[\mathbf{y}_k - \mathbf{H}(\mathbf{x}_k^f - \hat{\mathbf{b}}_{k-1}) \right]$$
$$\mathbf{x}_k^a = (\mathbf{x}_k^f - \hat{\mathbf{b}}_k) + \mathbf{K} \left[\mathbf{y}_k - \mathbf{H}(\mathbf{x}_k^f - \hat{\mathbf{b}}_k) \right]$$

which is expensive if all observations are used for the bias estimation.

If the bias varies slowly in time, we could use the previous bias estimate to correct the first guess, and reverse the order:

$$\mathbf{x}_k^a = (\mathbf{x}_k^f - \hat{\mathbf{b}}_{k-1}) + \mathbf{K} \left[\mathbf{y}_k - \mathbf{H}(\mathbf{x}_k^f - \hat{\mathbf{b}}_{k-1}) \right]$$
$$\hat{\mathbf{b}}_k = \hat{\mathbf{b}}_{k-1} - \mathbf{L} \left[\mathbf{y}_k - \mathbf{H}(\mathbf{x}_k^f - \hat{\mathbf{b}}_{k-1}) \right]$$

If we take $\mathbf{P}^b = \alpha \mathbf{P}^f$ and α is small, then $\mathbf{L} \approx \alpha \mathbf{K}$

Simplification of the Algorithm (2)

This leads to the very simple algorithm:

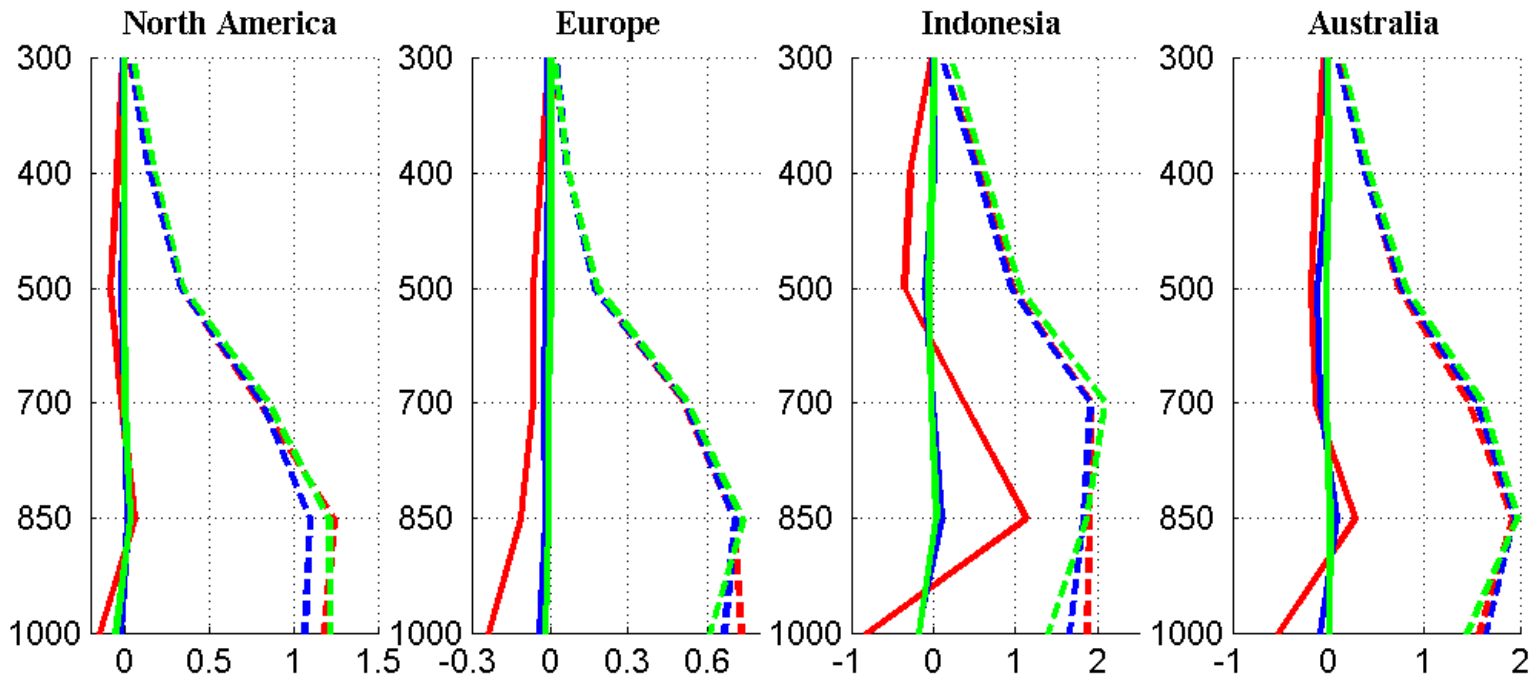
$\tilde{\mathbf{x}} = \mathbf{x}_k^f - \hat{\mathbf{b}}_{k-1}$	bias correction
$\mathbf{d}\mathbf{y} = \mathbf{y}_k - \mathbf{H}\tilde{\mathbf{x}}$ $\mathbf{d}\mathbf{x} = \mathbf{K}\mathbf{d}\mathbf{y}$ $\mathbf{x}_k^a = \tilde{\mathbf{x}} + \mathbf{d}\mathbf{x}$	} the usual analysis
$\hat{\mathbf{b}}_k = \hat{\mathbf{b}}_{k-1} - \alpha\mathbf{d}\mathbf{x}$	bias estimation

- the bias estimation and correction have been separated from the analysis
- it is essentially cost-free
- in a nonlinear analysis, one can put the entire algorithm in the inner loop
- one can choose to perform bias correction in terms of the model state representation, rather than the analysis control vector representation

Original Versus Simplified Scheme for Moisture

- Red: control (no bias correction)
- Green: original bias correction scheme
- Blue: simplified bias correction scheme

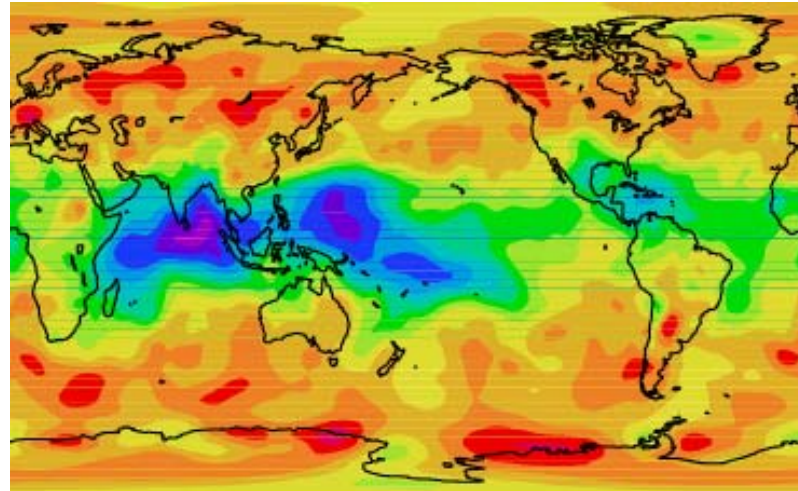
Jan1998 Observed-Forecast Means and Standard Deviations



Multivariate Bias Correction in fvDAS (1)

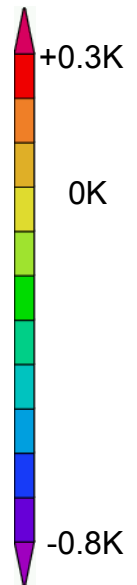
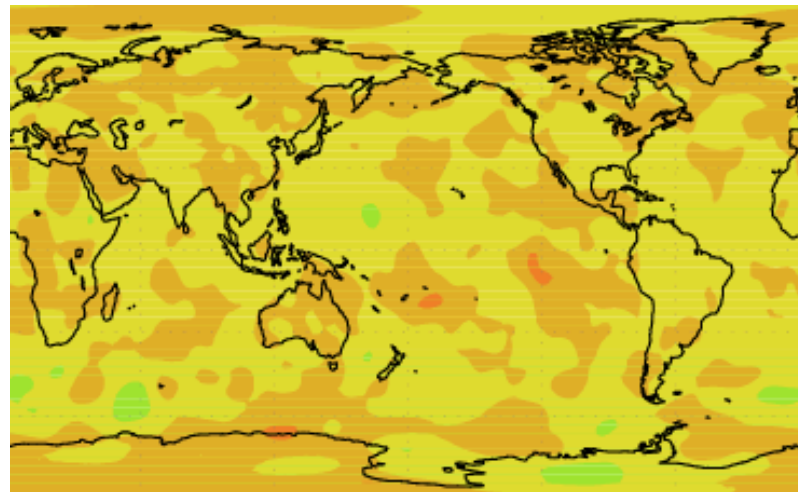
The model at 500hPa tends to be too cool in mid-latitudes; too warm in tropics

May 2011 mean temperature analysis increment, in fvDAS control assimilation, for model layer 8 (~500hPa)



Same, in experiment with multivariate bias correction

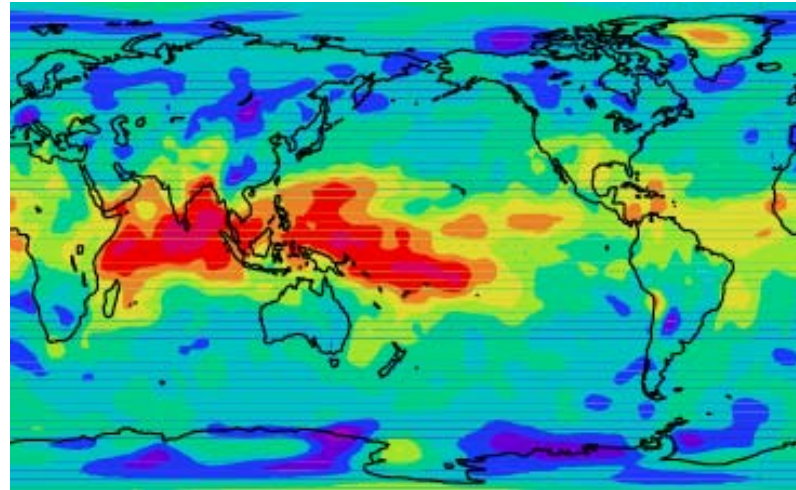
Bias correction reduces the mean analysis increments



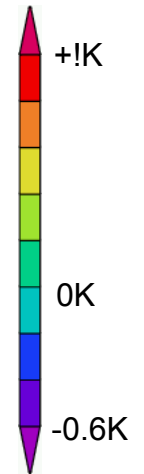
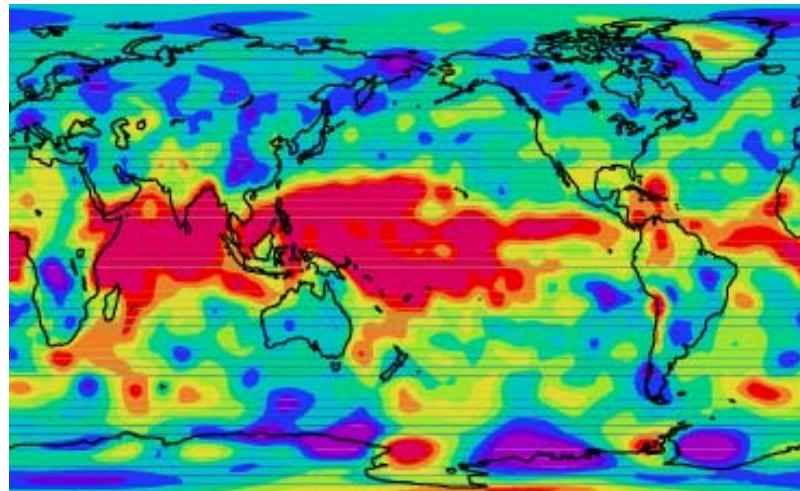
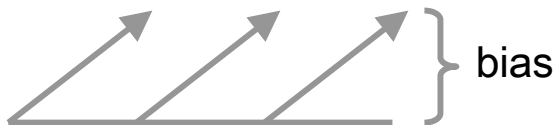
Multivariate Bias Correction in fvDAS (2)

But the temperature 6h forecast bias has increased!

May 2001 mean temperature bias, in control assimilation, for model layer 4 (~500hPa)



Same, in the experiment with multivariate bias correction



In this case, the bias correction clearly deteriorates the analyses.

What Could Possibly Go Wrong?

- We may have been inadvertently correcting for observation bias
- We may have introduced spurious effects via the background error covariance model

$$\tilde{\mathbf{x}} = \mathbf{x}_k^f - \hat{\mathbf{b}}_{k-1}$$

$$d\mathbf{y} = \mathbf{y}_k - \mathbf{H}\tilde{\mathbf{x}}$$

$$d\mathbf{x} = \mathbf{K}d\mathbf{y}$$

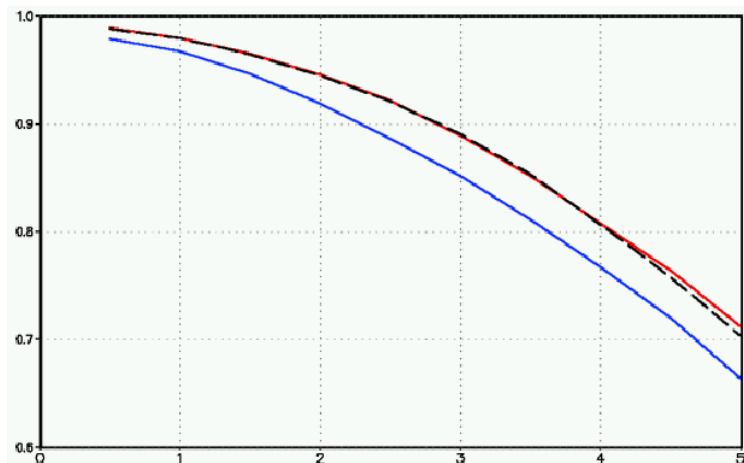
$$\mathbf{x}_k^a = \tilde{\mathbf{x}} + d\mathbf{x}$$

$$\hat{\mathbf{b}}_k = \hat{\mathbf{b}}_{k-1} - \alpha d\mathbf{x}$$

$$d\mathbf{x} \in \text{Range}(\mathbf{P})$$

- It may not be appropriate to apply persistent corrections to all model variables
- There may be a bug in the code..

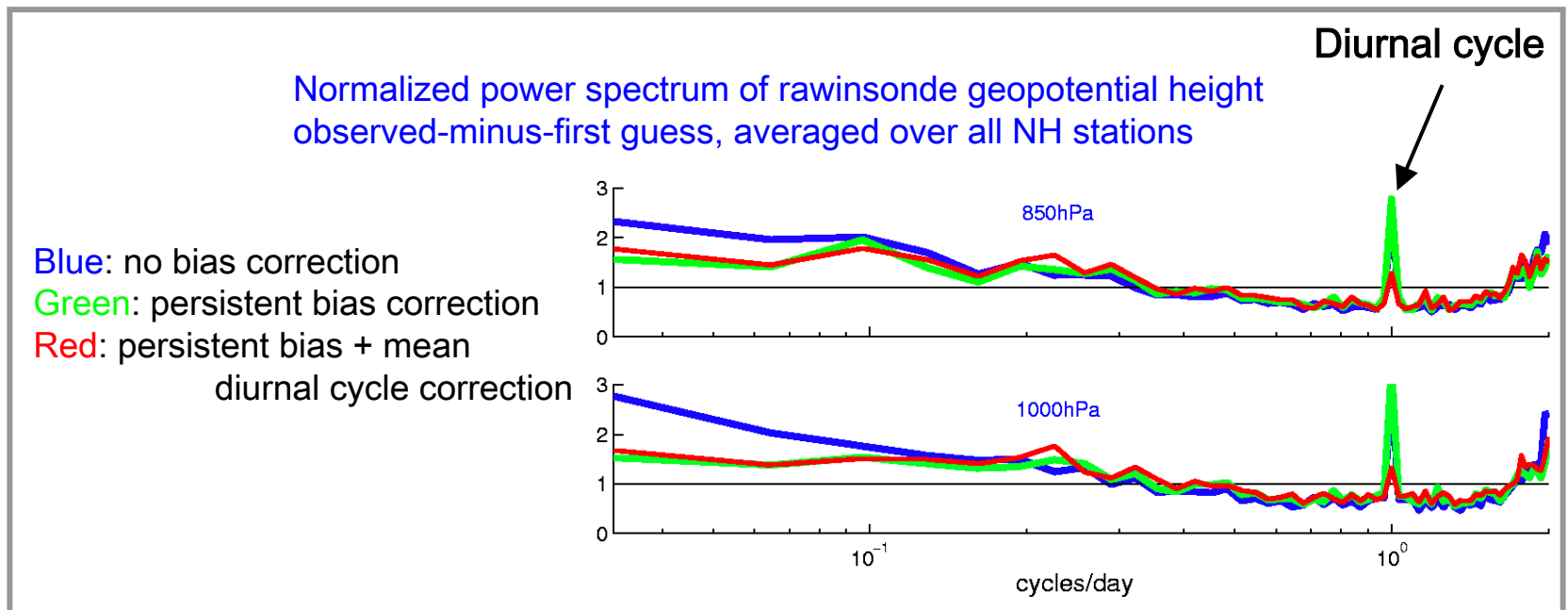
500hPa geopotential height anomaly correlations
27 cases – May/June 2001, Northern Hemisphere



Red: control (no bias correction)
Dashed: moisture bias correction
Blue: multivariate bias correction

Non-persistent Systematic Model Errors (1)

- Some model errors are clearly deterministic, yet not constant in time.
- Errors with known periodicity: Temperature bias in the model's diurnal cycle
- Visible in average power spectrum of station data time series



Non-persistent Systematic Model Errors (2)

Flow-dependent systematic errors:

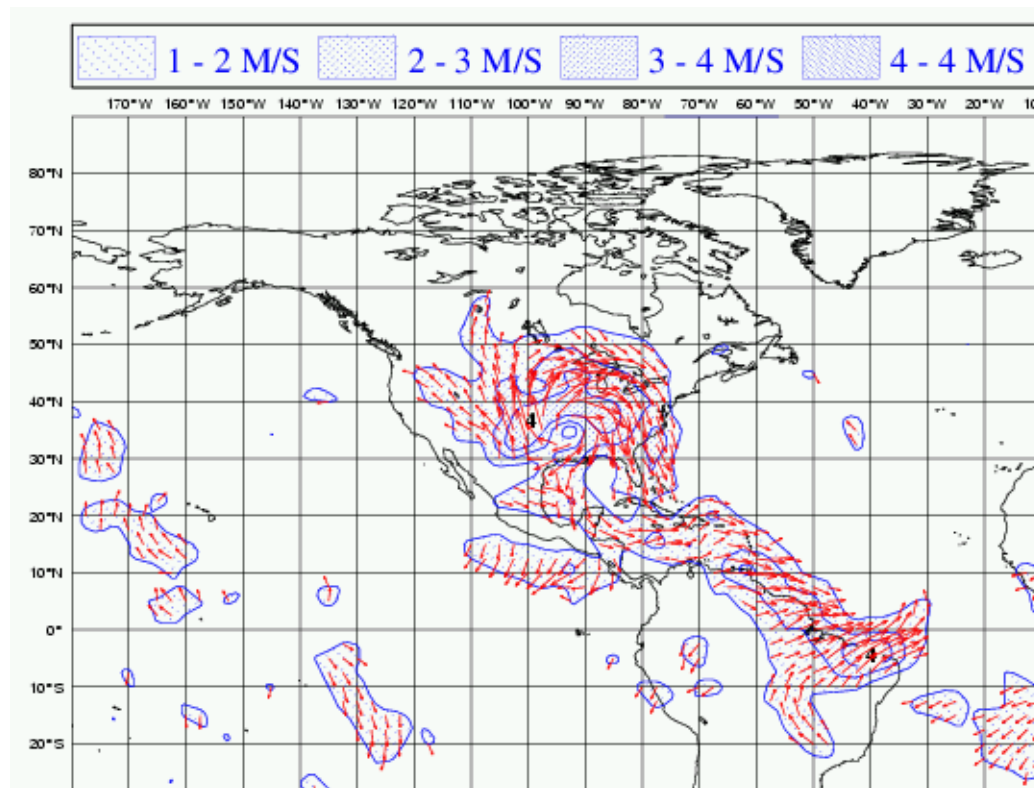
Related to dynamic developments that the model does not handle well – intermittent

From ECMWF MetOps:

Unusually large wind increments showing up in the monthly mean: excessive convective activity, divergence aloft

This problem has been addressed by improving the convective scheme in the model, and improving the humidity analysis

200 hPa Mean AN-FG Increment (m/s) at 12 UTC
Between 1/6/2003 and 30/6/2003



(Thanks to Antonio Garcia-Mendez)

Parameterization of Systematic Model Errors

We may be able to handle some (certainly not all) types of non-persistent systematic errors using statistical estimation.

Our basic assumption is that the first-guess error has a deterministic component (or bias) which is due to systematic model errors:

$$\mathbf{e}^f = \mathbf{b} + \tilde{\mathbf{e}}^f, \quad \langle \tilde{\mathbf{e}}^f \rangle = 0$$

Instead of assuming that the bias itself is persistent, it may be more natural to assume that it is a function of a set of persistent (or slowly varying) source parameters:

$$\mathbf{b} = \mathbf{b}(\boldsymbol{\beta}) \quad \boldsymbol{\beta} = [\beta_1 \ \beta_2 \ \dots \ \beta_N]^T$$

In principle, these could be related to the forcing of the model, or they could represent uncertain terms in the model formulation.

As long as we have a sufficient quantity of accurate observations that can be related to the unknown parameters, estimation theory provides tools to estimate those parameters.

A Measurement Model for the Bias Parameters

Our assumption is $\mathbf{e}^f = \mathbf{b}(\beta) + \tilde{\mathbf{e}}^f, \quad \langle \tilde{\mathbf{e}}^f \rangle = 0$

Then

$$\begin{aligned} \mathbf{y} - \mathbf{h}(\mathbf{x}^f) &\approx \mathbf{e}^o - \mathbf{H}\mathbf{e}^f \\ &= \mathbf{e}^o - \mathbf{H}\mathbf{b}(\beta) - \mathbf{H}\tilde{\mathbf{e}}^f \end{aligned}$$

If we define $\left\{ \begin{array}{l} \mathbf{d}\mathbf{y} = \mathbf{y} - \mathbf{h}(\mathbf{x}^f) \\ \tilde{\mathbf{e}} = \mathbf{e}^o - \mathbf{H}\tilde{\mathbf{e}}^f \\ \mathbf{g}(\beta) = -\mathbf{H}\mathbf{b}(\beta) \end{array} \right.$

then we have

$$\mathbf{d}\mathbf{y} = \mathbf{g}(\beta) + \tilde{\mathbf{e}} \quad \text{with} \quad \left\{ \begin{array}{l} \langle \tilde{\mathbf{e}} \rangle \approx 0 \\ \langle \tilde{\mathbf{e}}\tilde{\mathbf{e}}^T \rangle \approx \mathbf{H}\mathbf{P}^f\mathbf{H}^T + \mathbf{R} \end{array} \right.$$

β can be estimated, e.g., by variational methods!

Sequential Estimation of Linear Bias Parameters

Assume that $\mathbf{e}^f = \mathbf{b} + \tilde{\mathbf{e}}^f$, $\langle \tilde{\mathbf{e}}^f \rangle = 0$

where now $\mathbf{b} = \mathbf{b}(t_k) = \mathbf{B}_k \boldsymbol{\beta}$ and \mathbf{B}_k is a matrix with known coefficients.

Then the parameters $\boldsymbol{\beta}$ can be estimated from the data by

$$\hat{\boldsymbol{\beta}}_k = \hat{\boldsymbol{\beta}}_{k-1} - \mathbf{L}^\beta \left[\mathbf{y}_k - \mathbf{H}_k (\mathbf{x}_k^f - \mathbf{B}_k \hat{\boldsymbol{\beta}}_{k-1}) \right]$$

This estimator is optimal when

$$\mathbf{L}^\beta = \mathbf{P}^\beta \mathbf{B}^T \mathbf{H}^T [\mathbf{H} \mathbf{B} \mathbf{P}^\beta \mathbf{B}^T \mathbf{H}^T + \mathbf{H} \mathbf{P}^f \mathbf{H}^T + \mathbf{R}]^{-1}$$

with $\mathbf{P}^\beta = \langle (\hat{\boldsymbol{\beta}}_k - \boldsymbol{\beta})(\hat{\boldsymbol{\beta}}_k - \boldsymbol{\beta})^T \rangle$

The first-guess bias estimate $\hat{\mathbf{b}}_k = \mathbf{B}_k \hat{\boldsymbol{\beta}}_k$ can then be used to correct the background during data assimilation.

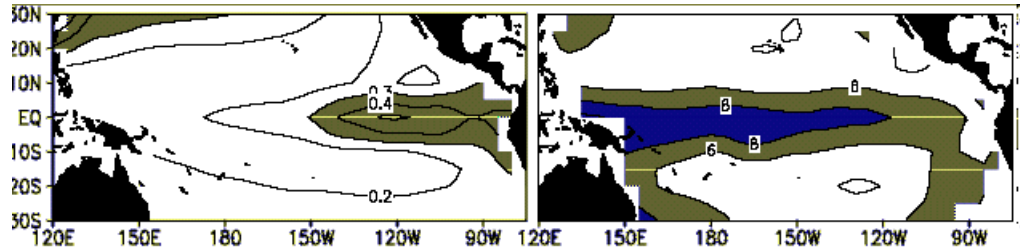
Bias Correction in a Tropical Ocean Model (1)

- Recent work with G. Chepurin and J. Carton (2003)
- Main goal is to identify large-scale, slowly varying features in model bias and to produce an analysis that is consistent with the observations
- Focus on tropical Pacific, 1970-2000, mixed layer temperature and thermocline depth
- Model: MOM2 (GFDL), $1x\frac{1}{2}x20L$ near equator to $1x1x20L$ mid-latitudes, sponge layer poleward of 62 degrees
- Data: temperature profiles from World Ocean Database + NOAA + TAO moorings, surface data from COADS
- Assimilation: simple OI with anisotropic background error correlations, and intermittent analysis update (IAU)

Bias Correction in a Tropical Ocean Model (2)

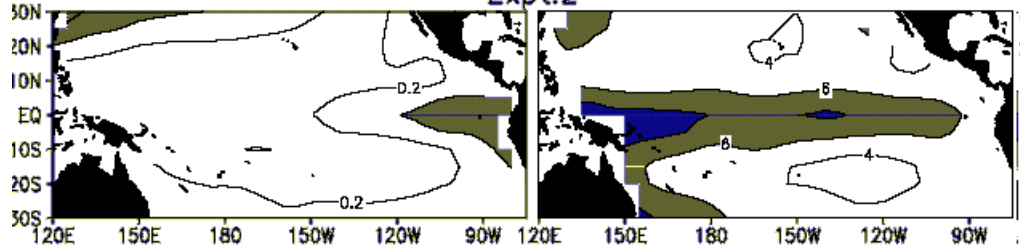
Rms observed-minus-forecast: Mixed layer temperature Thermocline depth

Control (no bias correction)



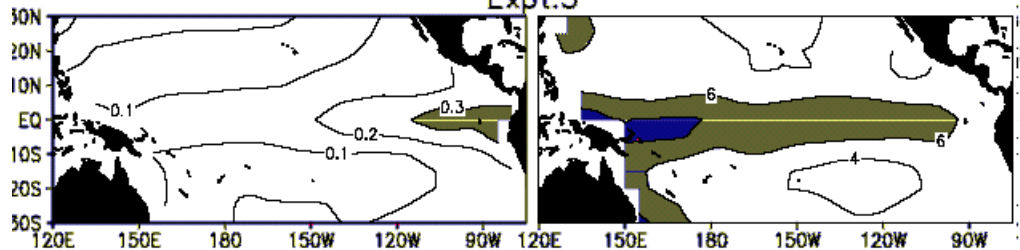
Time-mean bias correction

$$\beta_k^f = \bar{\beta}$$



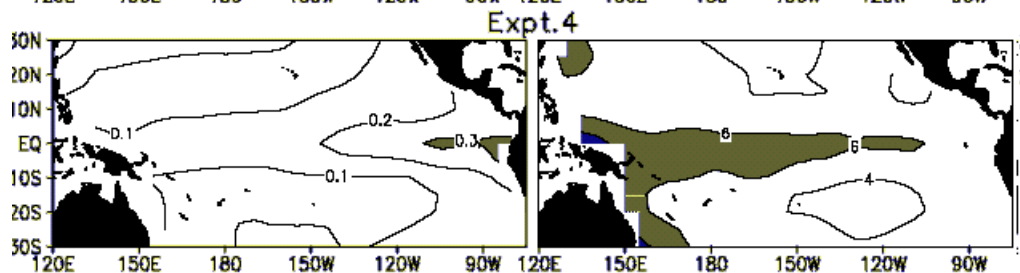
Time-mean plus annual cycle

$$\beta_k^f = \bar{\beta} + \beta_A e^{-i\omega_A t_k}$$



Time-mean plus annual cycle
plus 2 EOFs

$$\beta_k^f = \bar{\beta} + \beta_A e^{-i\omega_A t_k} + \sum_{n=1}^2 PC_n \cdot \tau_n(t_k)$$



Finally...

- Simple sequential estimation methods can be used to account for some deterministic components of first-guess error. These are statistical techniques that treat the model as a black box
- Similar methods can also be developed in the variational context (Derber 1992, Griffith and Nichols 1996, presentation by Trémolet)
- Ultimately one would like to identify and correct the model terms that cause the systematic errors (e.g. Bell, Martin, and Nichols 2001)
- We have assumed throughout that the observations are unbiased, but of course this is not to be taken for granted. Similar techniques can be used to estimate observation biases
- Separating model bias from observation bias requires hypotheses on the nature of each and lots of data (e.g. McNally and Watts, current work on AIRS)
- Computational tools and methods for handling the estimation problem are available – the data are forthcoming...