

Computational Optimizations and Emulation of EDA Perturbed Members and Statistics

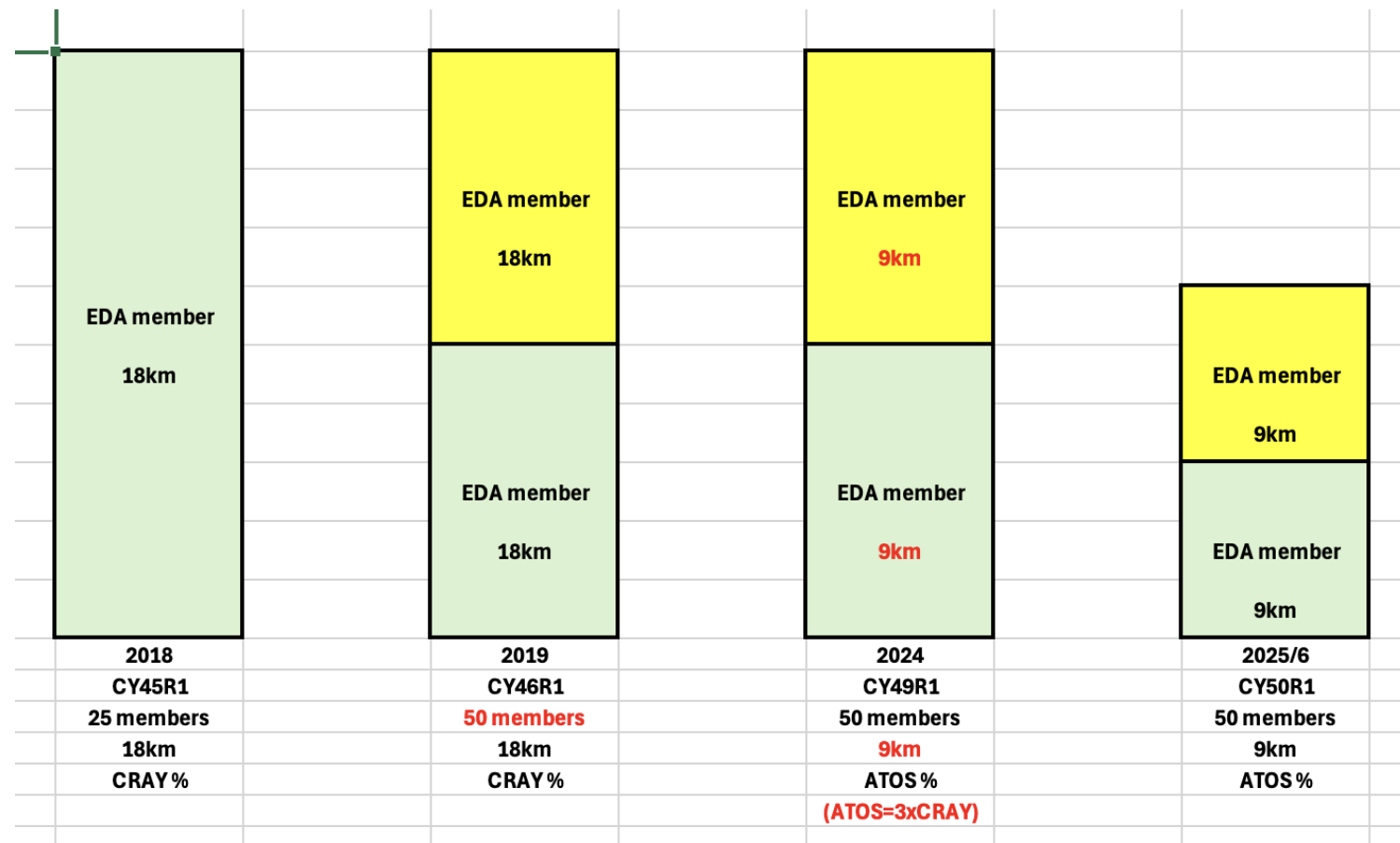
**ECMWF Workshop on Data Assimilation
Bonn 9-10 April 2025**

Elas Holm and Wei Pan

with Jorge Bandejas, Massimo Bonavita, Marcin Chrust, Alan Geer, Paddy Gillies, Simon Lang, Peter Lean and more contributing ideas and optimizations

Recent Efficiency Gains in EDA (x5) and How Used to Improve EDA

- Recent improvements in the computational efficiency of the EDA from technical and algorithmic developments (factor 5 since 2019) has been used to **increase the consistency and accuracy of the EDA using same or reduced fraction of the ECMWF supercomputers.**
- The two main applications of EDA are to **provide background errors B for the analysis (deterministic and ensemble)** and **contribute perturbations for ENS initial conditions.** The EDA developments have at each stage **improved deterministic and ensemble forecast scores.**



Recent Efficiency Gains in EDA (x5) and How Used to Improve EDA

- **2019 46R1:** 50% technical efficiency gain enables EDA members 25-->50 at same cost to match 50 ENS members number (both at 18km) ==>
 - improved high-resolution analysis (better background error B),
 - Improved EDA-based perturbations for the IC's of 50-member ENS,
 - Exchangeability of members improves validity of RD tests with few members ENS.
- **2024 49R1:** 30% algorithmic efficiency gain from soft re-centring enables EDA 18km-->9km to match ENS resolution (meanwhile increased to 9km), while using same proportion of new ATOS clusters as the old CRAY clusters. Combined with scientific improvement of model uncertainty parameterization (SPPT-->SPP) ==>
 - improved spread, improved high-resolution analysis (better background error B),
 - improved ENS initial conditions.
- **2025/6 50R1:** 30-40% technical and algorithmic efficiency gain (single precision trajectory, memory optimize 4D-Var, optimized inner loop resolution) enables to run EDA on 30-40% less resources (memory x wallclock). Combined with scientific improvement of spread calculation ==>
 - improved high-resolution analysis (better background error B),
 - improved ENS initial conditions.

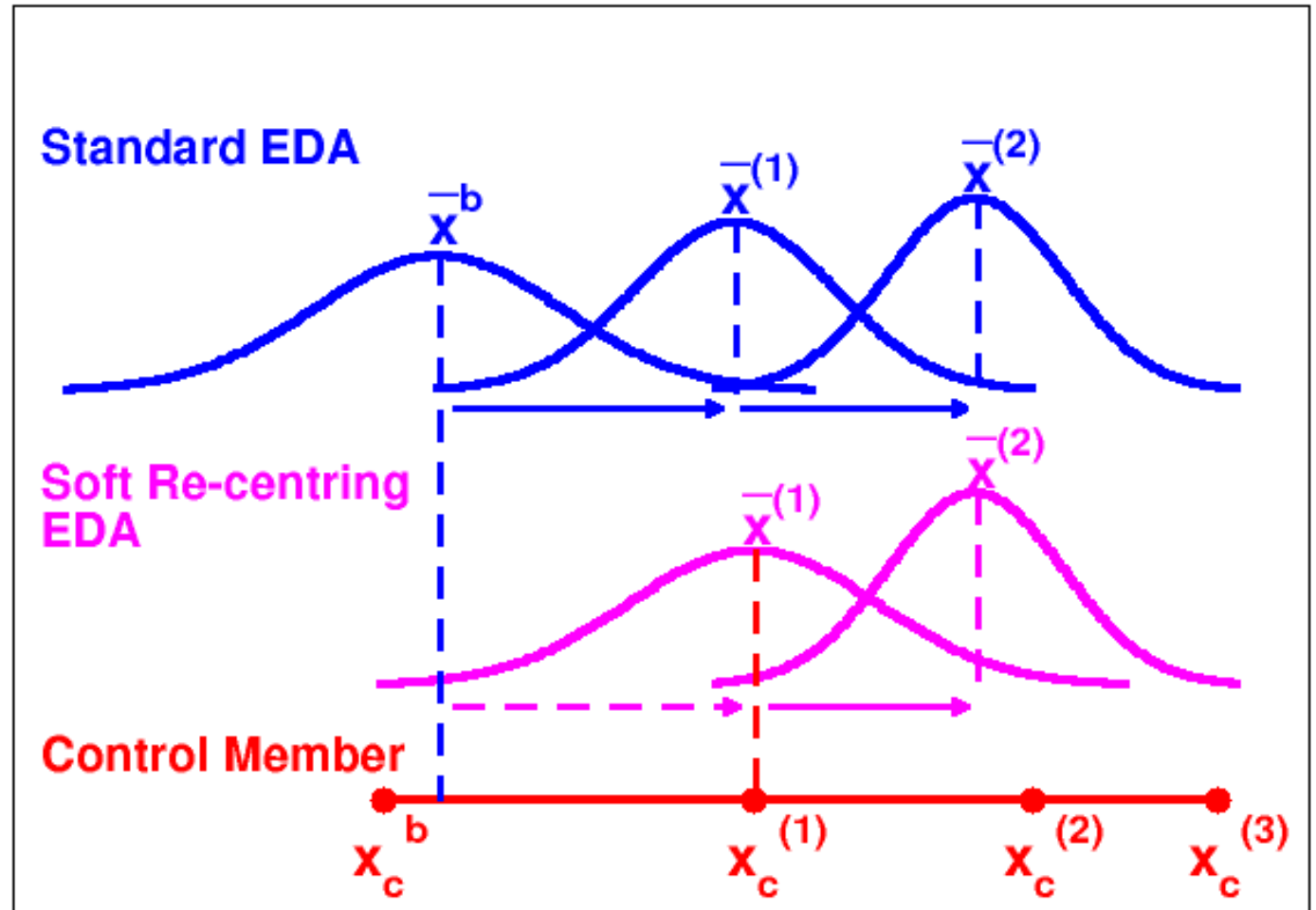
49R1 Soft Re-Centred EDA => One Outer Loop Equals Two!

The 49R1 EDA resolution increase to TCo1279 and introduction of SPP was needed for more realistic spread.

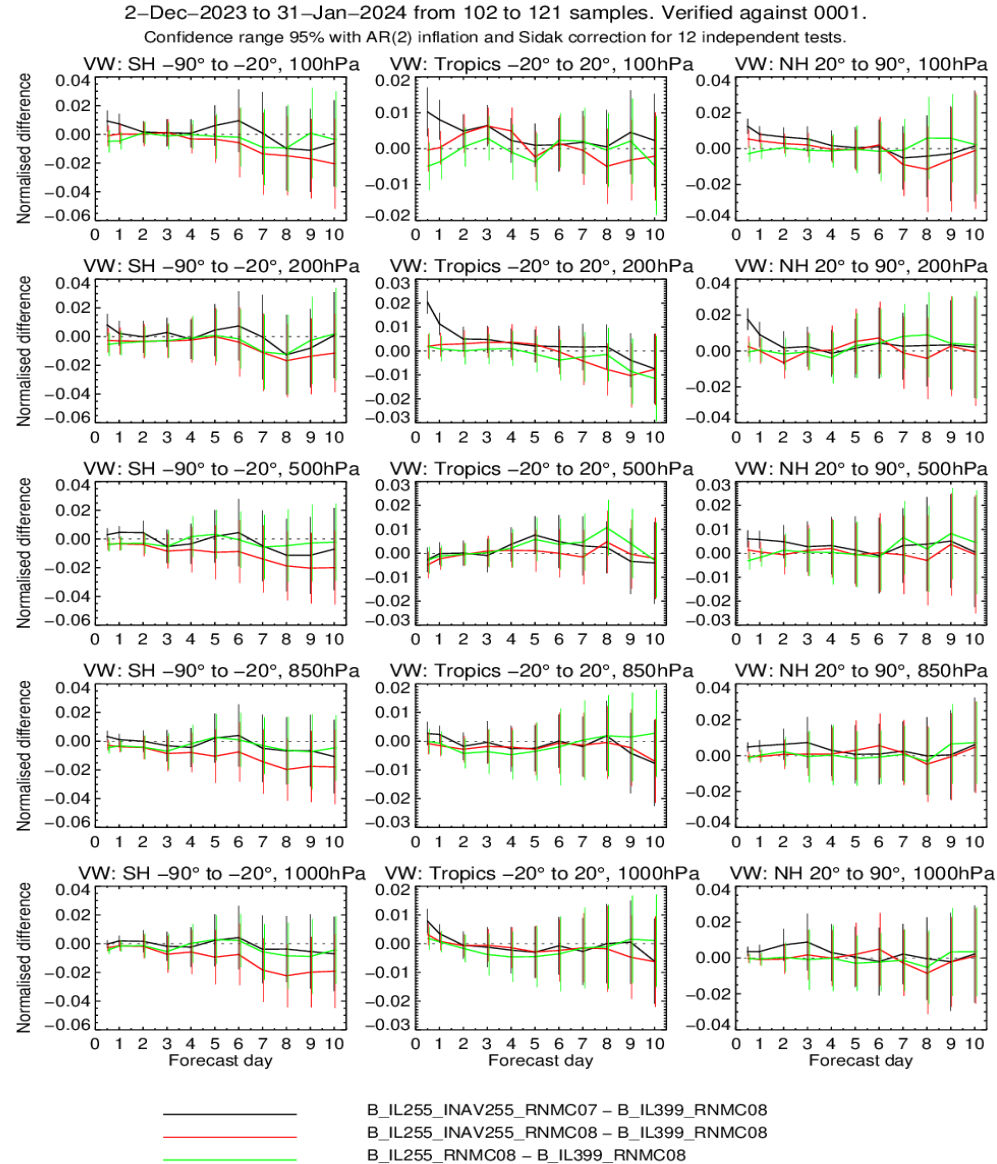
Required computational efficiency gain

➔ Soft Re-centred EDA:

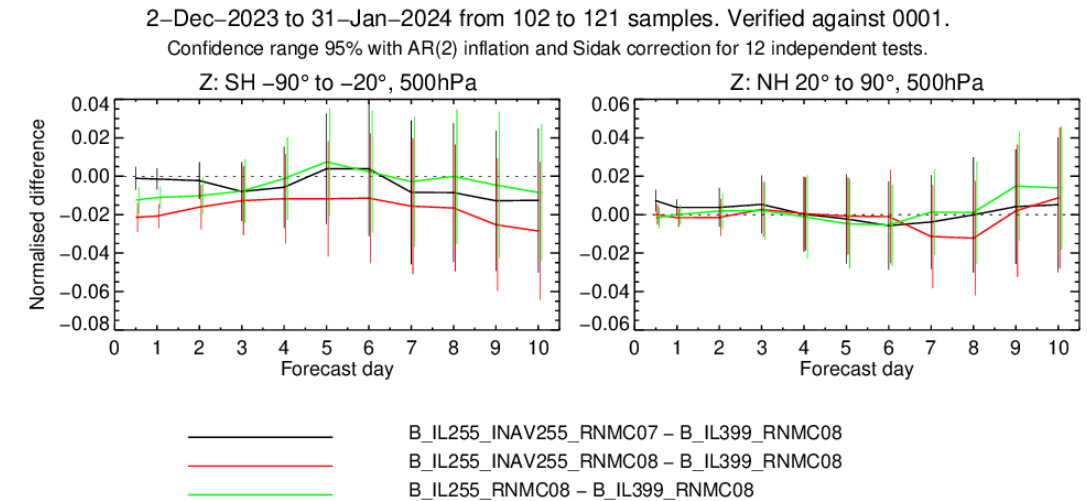
- Technically AND quality wise equivalent to doing two outer loops.
- Saves one EDA minimization, 30%, for same quality, enabling higher resolution EDA.
- Control runs more/higher resolution minimization for higher quality.



50R1 TCo1279 4D-Var using B from 10-mem TCo1279/TL255-319-511 EDA's:



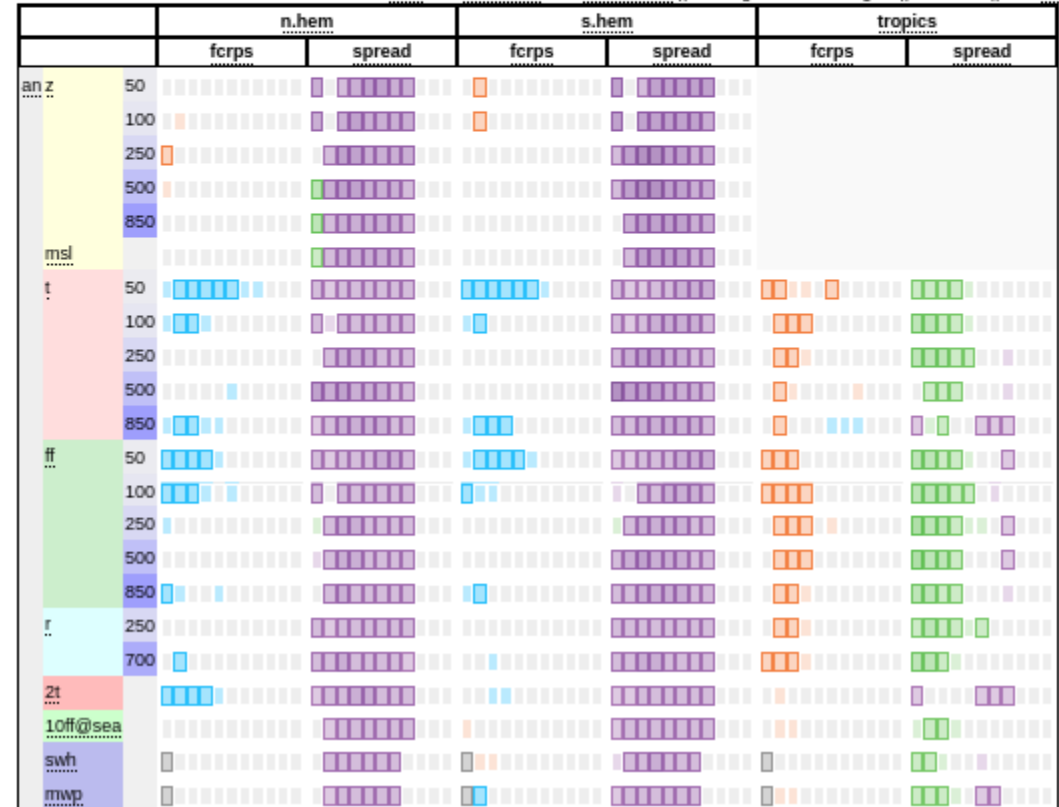
- 4D-Var scores wrt OPER neutral to advantage for TL255 inner loop with T255 spread (red line).
- Only TL255 inner loop not as good (green line).
- REDMNC=0.7 worse (black line), so current 0.8 fine.
- Note: Show OPER scores, OWN rmse scores biased in short range due to spread change of EDA/B.



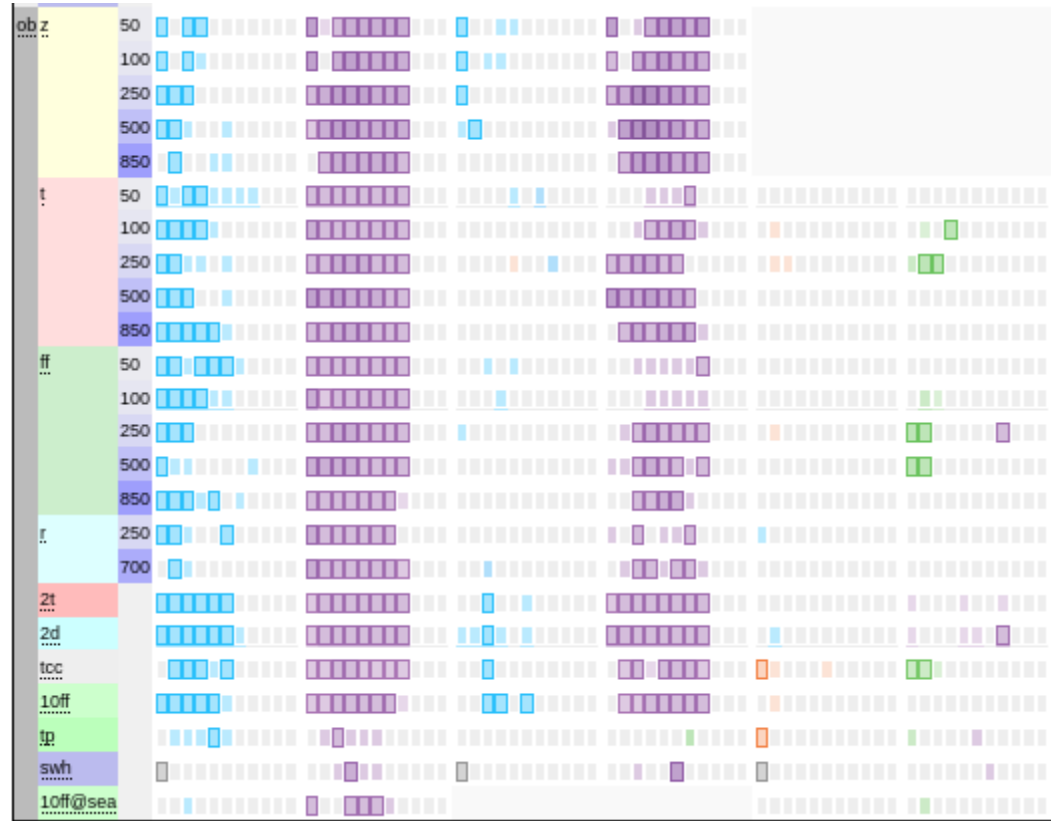
50R1 TCo1279 10-mem ENS/EDA with 50R1 Re-centring on e-suite+standard SV: TL255 IL+T255 Spread vs Reference TL399 IL+T159 spread (Dec-Jan 2023-24)

IL255_INAV255_RMC0.8_vs_IL399_INAV159_RMC0.8 scorecard

dates=[2023-12-02 00:00:00,2023-12-02 12:00:00,2023-12-03 00:00:00,...,2024-01-31 00:00:00,2024-01-31 12:00:00]
steps=[0, 24, 48, 72, 96, 120, 144, 168, 192, 216, 240]
vstreams=['qrdx_an', 'qrdx_ob']
classs=rd
streams=['enfo', 'waef']
expvers=(cntrl:ik97, exper:ik9a)
reftypes=['an', 'ob']
☒ fcrps ☒ spread
☒ n.hem ☒ s.hem ☒ tropics ☐ europe ☐ n.atl ☐ n.pac (☒ all)
shaded boxes for confidence boundaries: ☐ 95% ☐ 50%/95% ☒ 95%/99.7% || ☐ significance triangles || ☐ bars || ☐ sam



FCRPS wrt OBS/OPER improved/neutral in Extratropics.
Tropics worse short range wrt OPER, neutral wrt OBS
(same as in T255 IL only)



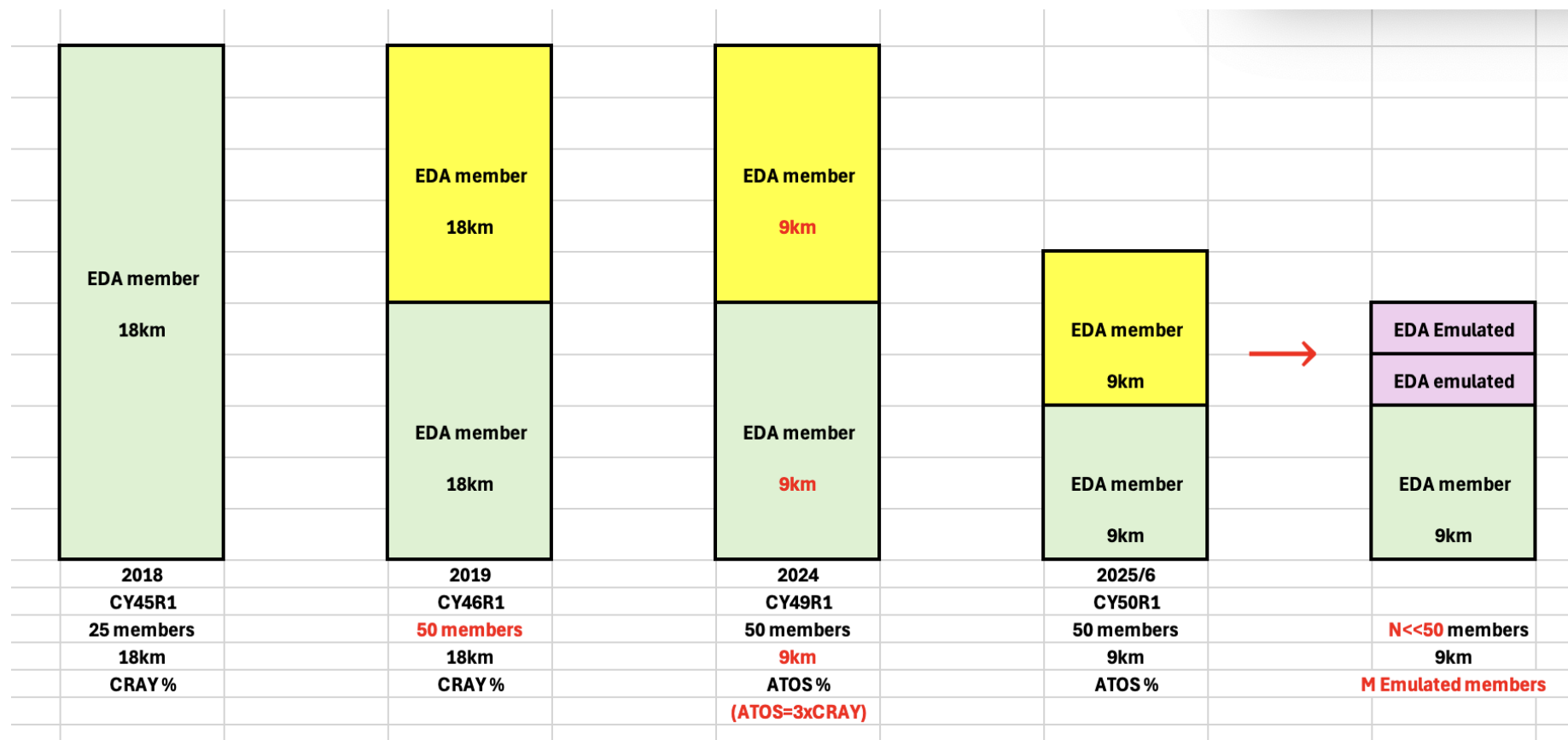
OPER

EUROPEAN CENTRE FOR MEDIUM-RANGE WEATHER FORECASTS

OBS

Future Efficiency Gains of EDA: Emulate Subset Members with Fast ML?

- Can we keep only a few full-resolution physical EDA members, while emulating the rest with ML and keep equal or better background errors and ENS initial conditions that improve deterministic and ensemble forecast scores?



4DVar Data Assimilation

- ◇ Bayes theorem + Gaussianity assumption:

$$\underbrace{\pi(\mathbf{x}|\mathbf{y}_{\text{obs}})}_{\text{posterior}} \propto \underbrace{\exp(-0.5(\mathbf{x} - \bar{\mathbf{x}})^T \mathbf{B}^{-1}(\mathbf{x} - \bar{\mathbf{x}}))}_{\text{prior}} \underbrace{\pi(\mathbf{y}_{\text{obs}}|\mathbf{x})}_{\text{likelihood}}. \quad (1)$$

- ◇ An *analysis*

$$\hat{\mathbf{x}} = \arg \max_{\mathbf{x}} \pi(\mathbf{x} | \mathbf{y}_{\text{obs}}) \quad (2)$$

is the optimiser of

$$L(\mathbf{x}) = \frac{1}{2}(\mathbf{x} - \bar{\mathbf{x}})^T \mathbf{B}^{-1}(\mathbf{x} - \bar{\mathbf{x}}) + \text{"log of likelihood"} \quad (3)$$

- ◇ We use *ensemble statistics* to estimate $\bar{\mathbf{x}}$ and \mathbf{B} :

$$\hat{\mathbf{x}}^1, \dots, \hat{\mathbf{x}}^{50}.$$

Cost saving using emulation

Why – Optimisation is expensive! Need many evaluations of TL/AD operators.
Can we make things cheaper using alternative methods?

Cost saving using emulation

- Why** – Optimisation is expensive! Need many evaluations of TL/AD operators.
Can we make things cheaper using alternative methods?
- How** – Build a model to *generate* physically consistent perturbations cheaply.

Cost saving using emulation

Why – Optimisation is expensive! Need many evaluations of TL/AD operators.
Can we make things cheaper using alternative methods?

How – Build a model to *generate* physically consistent perturbations cheaply.

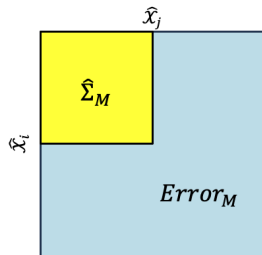
What

- **Statistics emulator** for the 'diagonal' part of \mathbf{B} is being tested for its impact on operational ensemble 4DVar.
- **Perturbations generator** is in development for generating analysis perturbations.

Statistics emulator

The 'diagonal' part of \mathbf{B} is denoted by Σ , whose entries are pointwise bg error stddev, estimated using EDA sample variance

$$\Sigma^2 = \frac{1}{N^2} \sum_{i=1}^N \sum_{j>i}^N (\hat{\mathbf{x}}^i - \hat{\mathbf{x}}^j)^{\otimes 2}, \quad (N = 50 \text{ currently}). \quad (4)$$



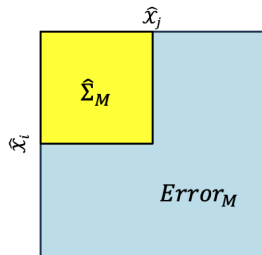
Statistics emulator

The 'diagonal' part of \mathbf{B} is denoted by Σ , whose entries are pointwise bg error stddev, estimated using EDA sample variance

$$\Sigma^2 = \frac{1}{N^2} \sum_{i=1}^N \sum_{j>i}^N (\hat{\mathbf{x}}^i - \hat{\mathbf{x}}^j)^{\otimes 2}, \quad (N = 50 \text{ currently}). \quad (4)$$

We take a subset of EDA members, $\hat{\mathbf{x}}^1, \dots, \hat{\mathbf{x}}^M$, $M < N$, calculate

$$\hat{\Sigma}_M^2 = \frac{1}{M^2} \sum_i^M \sum_{j>i}^M (\hat{\mathbf{x}}^i - \hat{\mathbf{x}}^j)^{\otimes 2}, \quad (\text{currently } M = 5) \quad (5)$$



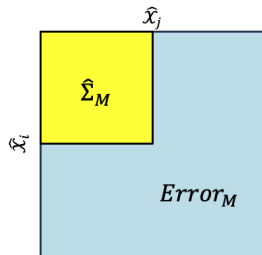
Statistics emulator

The 'diagonal' part of \mathbf{B} is denoted by Σ , whose entries are pointwise bg error stddev, estimated using EDA sample variance

$$\Sigma^2 = \frac{1}{N^2} \sum_{i=1}^N \sum_{j>i}^N (\hat{\mathbf{x}}^i - \hat{\mathbf{x}}^j)^{\otimes 2}, \quad (N = 50 \text{ currently}). \quad (4)$$

We take a subset of EDA members, $\hat{\mathbf{x}}^1, \dots, \hat{\mathbf{x}}^M$, $M < N$, calculate

$$\hat{\Sigma}_M^2 = \frac{1}{M^2} \sum_{i=1}^M \sum_{j>i}^M (\hat{\mathbf{x}}^i - \hat{\mathbf{x}}^j)^{\otimes 2}, \quad (\text{currently } M = 5) \quad (5)$$



then have a trained parametric model provide a guess of Σ based on $\hat{\Sigma}_M$

$$\hat{\Sigma}_M \mapsto \Sigma.$$

Note, by C.L.T.: $\hat{\Sigma}_M^2 - \Sigma^2 \sim \frac{1}{\sqrt{M}} \text{sqrt}(\mu_4 - \Sigma^4) \mathcal{N}(0, 1)$

Statistics emulator

Our emulation model is a parametric conditional probability distribution,

$$\hat{\Sigma}_M \mapsto p_{\theta}(\tilde{\Sigma} \mid \hat{\Sigma}_M) \xrightarrow{\text{sample}} \bar{\Sigma} \approx \Sigma. \quad (6)$$

Statistics emulator

Our emulation model is a parametric conditional probability distribution,

$$\hat{\Sigma}_M \mapsto p_{\theta}(\tilde{\Sigma} | \hat{\Sigma}_M) \xrightarrow{\text{sample}} \tilde{\Sigma} \approx \Sigma. \quad (6)$$

We implement $p_{\theta}(\tilde{\Sigma} | \hat{\Sigma}_M)$ as a nonlinear regression,

$$g_{\theta_1} \left(f_{\theta_2}(\hat{\Sigma}_M), e \right) \xrightarrow{\text{sample}} \tilde{\Sigma}, \quad e \sim q_{\theta_3}(\mathbf{z} | \hat{\Sigma}_M). \quad (7)$$

Statistics emulator

Our emulation model is a parametric conditional probability distribution,

$$\hat{\Sigma}_M \mapsto p_{\theta}(\tilde{\Sigma} | \hat{\Sigma}_M) \xrightarrow{\text{sample}} \tilde{\Sigma} \approx \Sigma. \quad (6)$$

We implement $p_{\theta}(\tilde{\Sigma} | \hat{\Sigma}_M)$ as a nonlinear regression,

$$g_{\theta_1}(f_{\theta_2}(\hat{\Sigma}_M), e) \xrightarrow{\text{sample}} \tilde{\Sigma}, \quad e \sim q_{\theta_3}(\mathbf{z} | \hat{\Sigma}_M). \quad (7)$$

We build six emulators, one for each model variable (vo, dv, Insp, q, O3 and t).

Each model has ~ 199114 trainable parameters.

Statistics emulator

Our emulation model is a parametric conditional probability distribution,

$$\hat{\Sigma}_M \mapsto p_{\theta}(\tilde{\Sigma} | \hat{\Sigma}_M) \xrightarrow{\text{sample}} \tilde{\Sigma} \approx \Sigma. \quad (6)$$

We implement $p_{\theta}(\tilde{\Sigma} | \hat{\Sigma}_M)$ as a nonlinear regression,

$$g_{\theta_1}(f_{\theta_2}(\hat{\Sigma}_M), e) \xrightarrow{\text{sample}} \tilde{\Sigma}, \quad e \sim q_{\theta_3}(\mathbf{z} | \hat{\Sigma}_M). \quad (7)$$

We build six emulators, one for each model variable (vo, dv, Insp, q, O3 and t).

Each model has ~ 199114 trainable parameters.

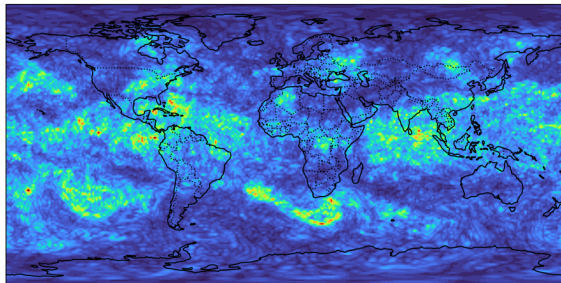
Impact on analysis: consider $\tilde{\mathbf{B}} := \mathbf{B} + \epsilon \delta \mathbf{B}$, (for 3DVar) we have

$$\|d \delta \mathbf{x}\| \leq c \|\delta \mathbf{B}\| \quad (8)$$

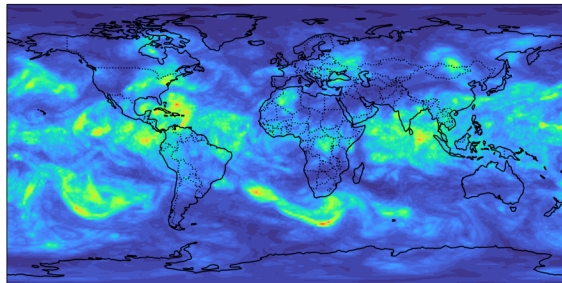
where $d \delta \mathbf{x}$ denotes *change in analysis increment*.

Experiment results – vo mlev74 20220601 1800

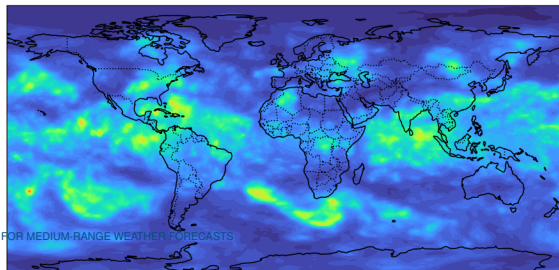
5-EDA-es



es

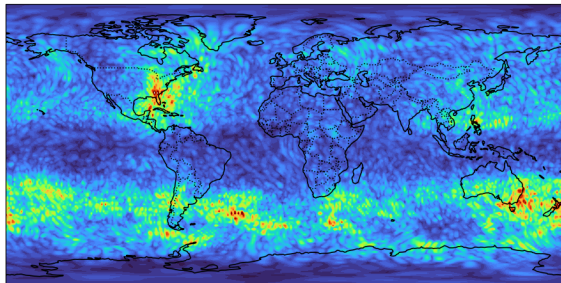


ML-es

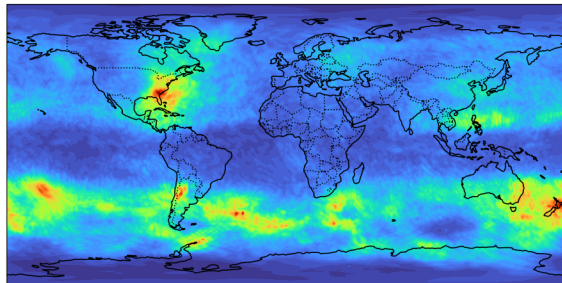


Experiment results – ucdv mlev02 20220601 1800

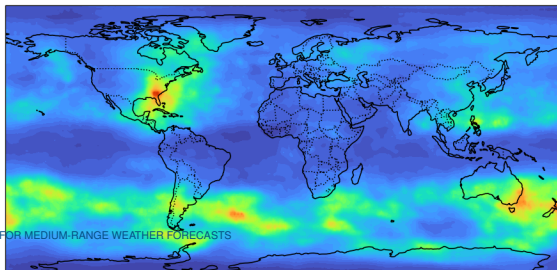
5-EDA-es



es



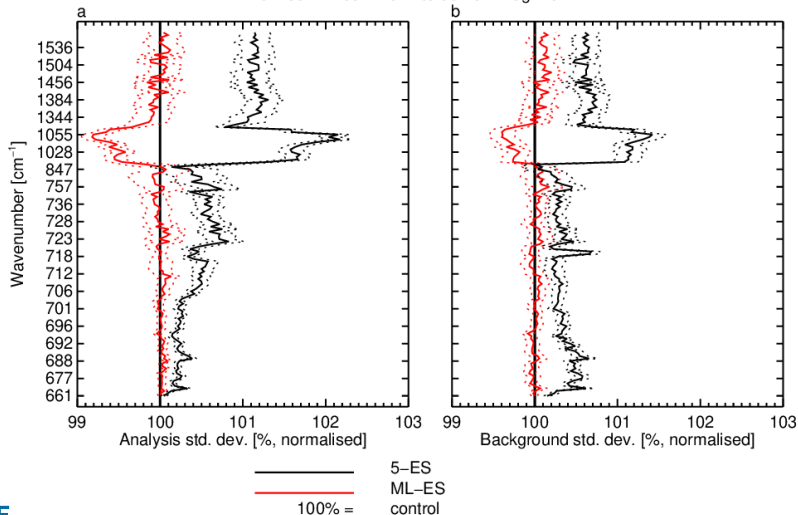
ML-es



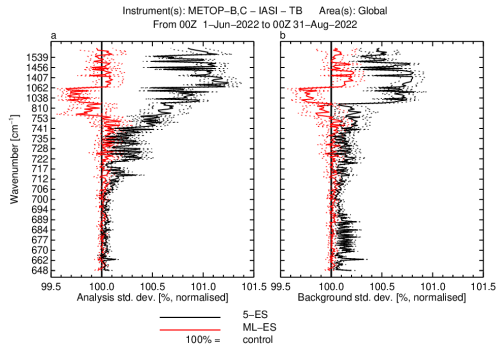
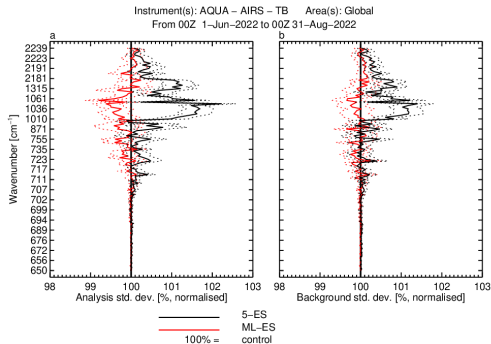
Experiment results – (iver) ref exp 50r1 control

Instrument(s): NOAA-20; NPP – CRIS – TB Area(s): Global

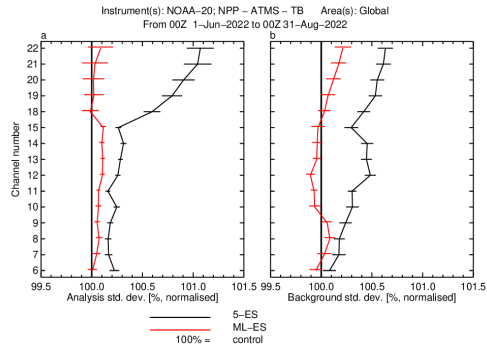
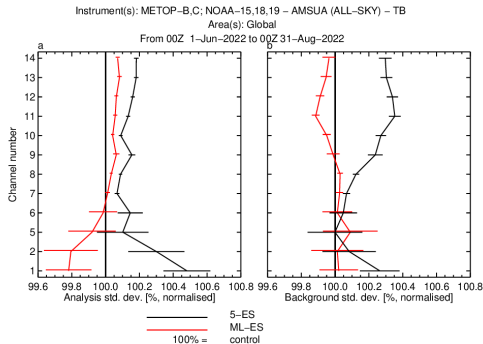
From 00Z 1-Jun-2022 to 00Z 31-Aug-2022



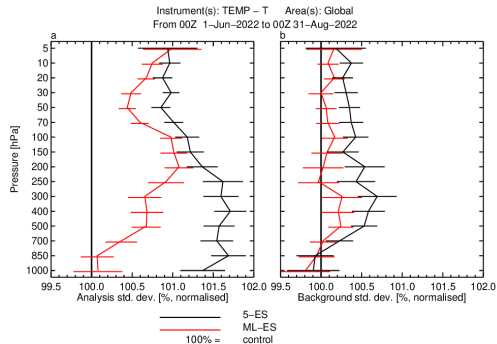
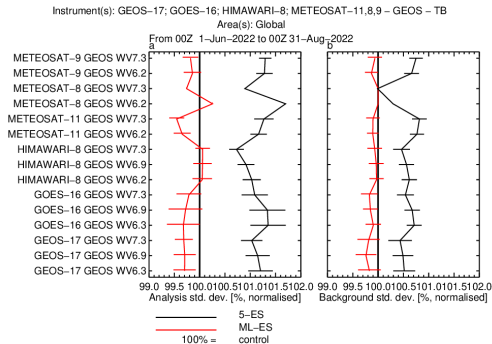
Experiment results – (iver) ref exp 50r1 control



Experiment results – (iver) ref exp 50r1 control



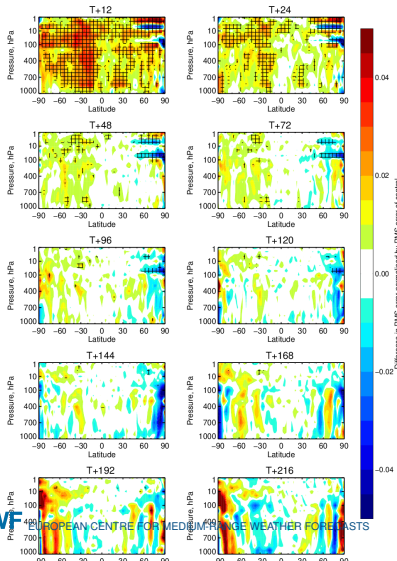
Experiment results – (iver) ref exp 50r1 control



Experiment results – (iver) ref exp 50r1 control - T fc

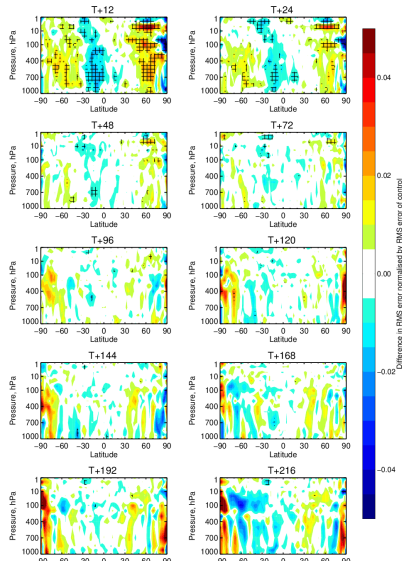
Change in RMS error in T (5-ES-control)

1-Jun-2022 to 31-Aug-2022 from 164 to 183 samples. Verified against 0001.
Cross-hatching indicates 99% confidence with Sidak correction for 20 independent tests.



Change in RMS error in T (ML-ES-control)

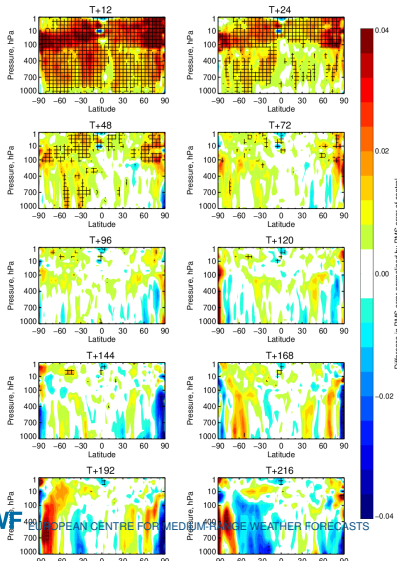
1-Jun-2022 to 31-Aug-2022 from 164 to 183 samples. Verified against 0001.
Cross-hatching indicates 99% confidence with Sidak correction for 20 independent tests.



Experiment results – (iver) ref exp 50r1 control - wind fc

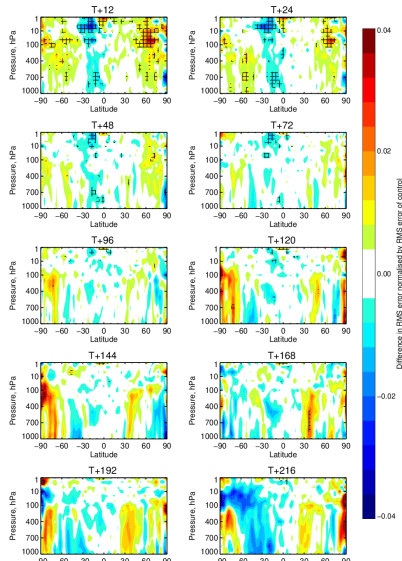
Change in RMS error in VW (5-ES-control)

1-Jun-2022 to 31-Aug-2022 from 164 to 183 samples. Verified against 0001.
Cross-hatching indicates 99% confidence with Sidak correction for 20 independent tests.



Change in RMS error in VW (ML-ES-control)

1-Jun-2022 to 31-Aug-2022 from 164 to 183 samples. Verified against 0001.
Cross-hatching indicates 99% confidence with Sidak correction for 20 independent tests.



Perturbations generator

The perturbations generator is a parametric conditional probability distribution

$$\nu_{\theta}(\mathbf{x}', t \mid \mathbf{x}'_1, \dots, \mathbf{x}'_M) \xrightarrow{\text{sample}} \mathbf{x}'^a. \quad (9)$$

which we view as the solution to a specific instance of the following problem.

Problem

Let $\mathfrak{X}(\mathcal{M})$ denote the space of vector fields on a smooth manifold \mathcal{M} . Let $\mathcal{P}(\mathfrak{X}(\mathcal{M}))$ denote the space of probability measures on $\mathfrak{X}(\mathcal{M})$. Let D be a metric on $\mathcal{P}(\mathfrak{X}(\mathcal{M}))$.

Given $\mu, \tilde{\mu} \in \mathcal{P}(\mathfrak{X}(\mathcal{M}))$ with $\tilde{\mu} \ll \mu$, find $\nu^* \in \mathcal{P}(\mathfrak{X}(\mathcal{M}))$ such that

$$\nu^* = \arg \min_{\substack{\nu \in \mathcal{P}(\mathfrak{X}(\mathcal{M})) \\ \nu \ll \mu}} D(\mu, \alpha \tilde{\mu} + (1 - \alpha)\nu), \quad (10)$$

for some $\alpha \in (0, 1)$.

Perturbations generator

ETKF for Lorenz96 standard setup: $F = 8$, $N_x = 40$

ETKF analysis perturbs ensemble:

$$x_i^{'a}(t), \quad i = 1, \dots, 20$$

Subset-ensemble
+ ML generated perturbs:

$$x_i^{'a}(t), \quad i = 1, \dots, 10$$

$$\text{and } \tilde{x}_j^{'a}(t), \quad j = 11, \dots, 20$$

Compare analysis
covariance matrices $P_{\text{full}}^a(t)$,
 $P_{\text{sub}}^a(t)$ and $P_{\text{sub+ML}}^a(t)$

Errors between analysis covariance matrices



Summary and outlook

- ◇ First results (statistics emulator) sets a benchmark for further improvements.
- ◇ More tests and evaluations are in the pipeline.
- ◇ Ensemble correlations.
- ◇ Continued development of the perturbations generator.