# Alberto Carrassi

## Merging DA and ML at various degrees
*Examples from applications to*
*Arctic sea ice and ocean biogeochemistry*

With:

**Ivo Pasmans, Giovanni De Cillis, Ieuan Higgs, Tobias Finn, Yumeng Chen, Marc Bocquet, Julien Brajard, Matteo Broccoli, Laurent Bertino, Ross Bannister, Stefano Ciavatta, Jozef Skakala, Simon Driscoll**
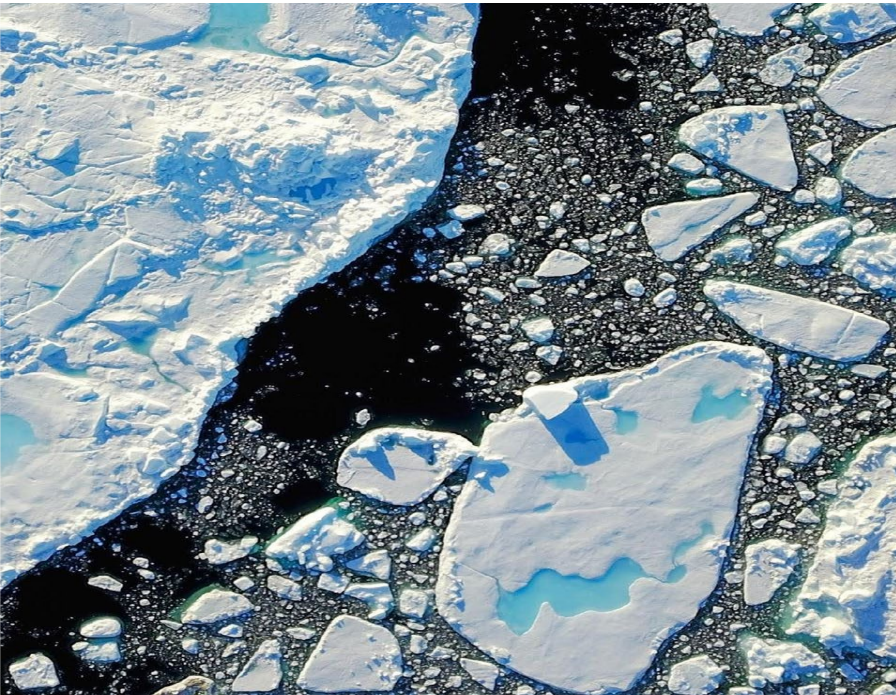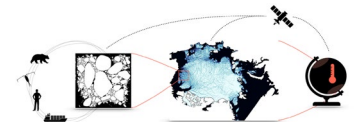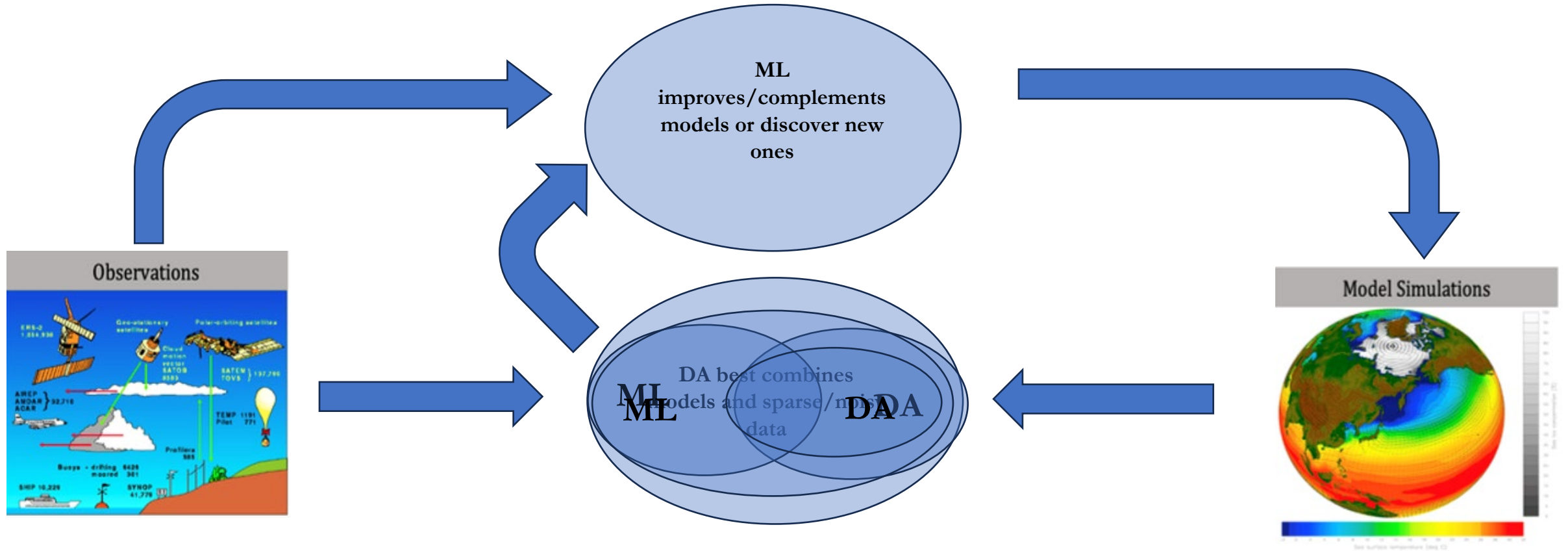
ALMA MATER STUDIORUM
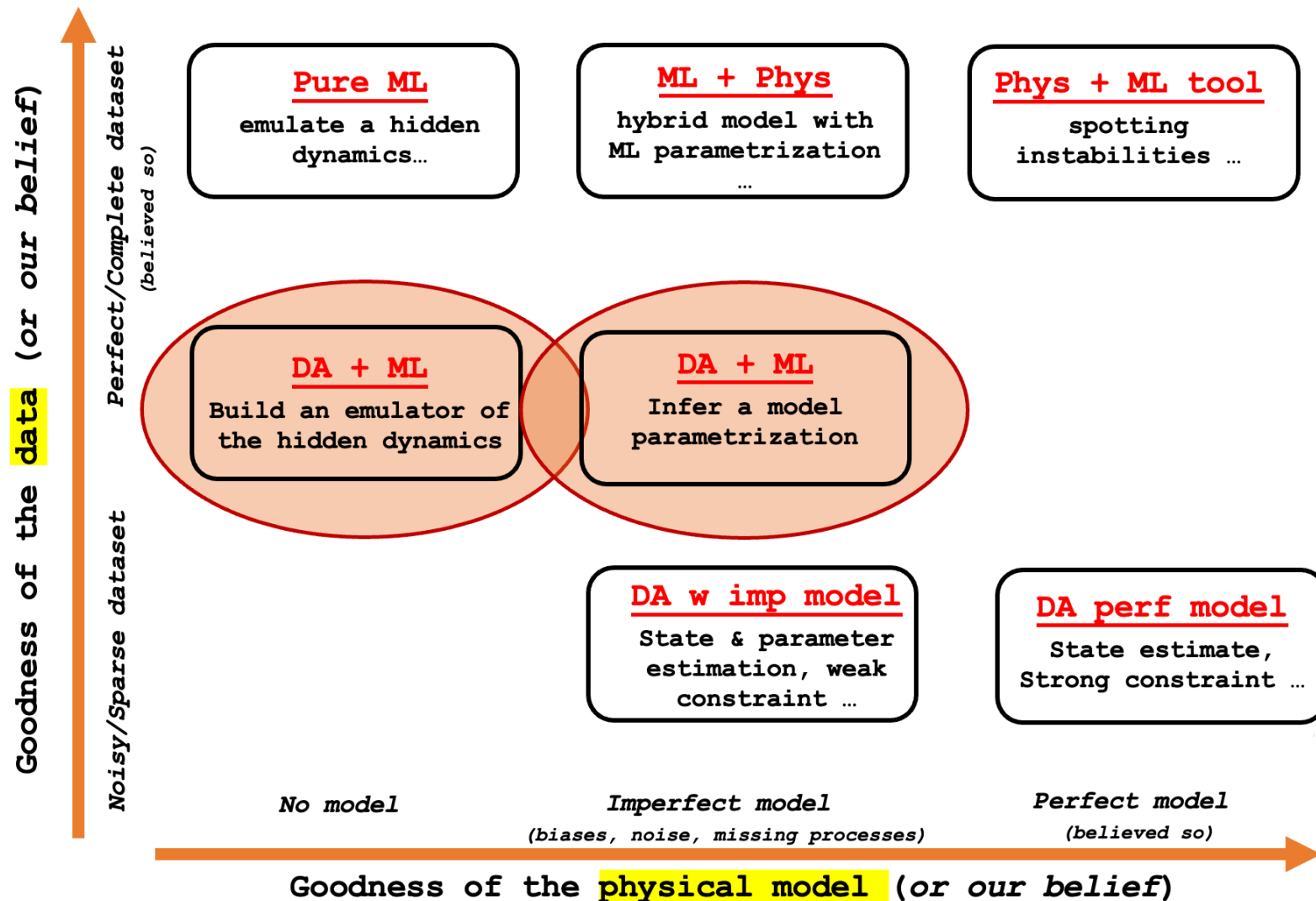UNIVERSITÀ DI BOLOGNA

University of Reading

The Scale-Aware Sea Ice Project

MELTED

Schmidt Sciences

# DA and ML: the very many ways of making them working together



ML improves/complements models or discover new ones

DA best combines models and sparse/noisy data

ML

DA

Observations

Model Simulations

Merging DA and ML to cope with their respective weakness (e.g., DA in the latent model and po space, nonlinear DA using NN) and prediction

## Combined DA-ML to infer unresolved scales parametrizations

**The objective is to produce a hybrid (physical/data-driven) model**

$$\mathbf{x}(t + \delta t) = \mathcal{M}^{\varphi}[\mathbf{x}(t)] + \mathcal{M}^{\text{UN}}[\mathbf{x}(t)],$$
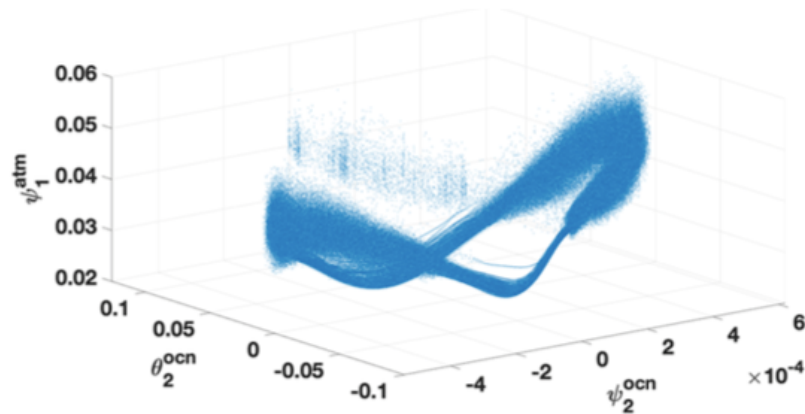
where:

- $\mathbf{x}(t)$ is the state of the dynamical system
- $\mathcal{M}^{\varphi}$ is the physical model (assumed to be known a priori)
- $\mathcal{M}^{\text{UN}}$ is the unresolved component of the dynamics to be inferred from data
- $\delta t$ is the integration time step

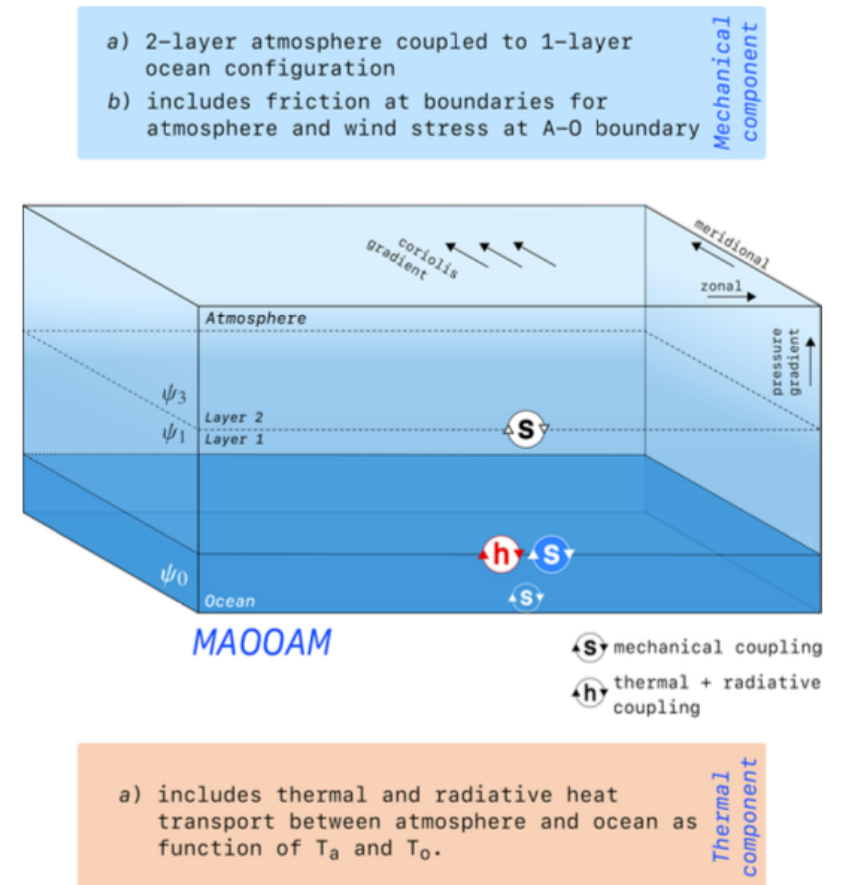$\mathcal{M}^{\text{UN}}$ is approximated by a **data-driven model** represented under the form of a neural network whose parameters are $\boldsymbol{\theta}$: $\mathcal{M}_{\boldsymbol{\theta}}[\mathbf{x}(t)]$

## Experiments with the coupled atmosphere-ocean MAOOAM

▶ **MAOOAM**: Modular arbitrary-order ocean-atmosphere model (Le Cruz *et al*, 2016)

▶ A two-layer QG atmosphere coupled, thermally and mechanically, to a QG shallow-water ocean layer in the $\beta$-plane.

▶ MAOOAM is resolved in spectral space, for streamfunction and potential temperature, with adjustable resolution.



▶ We implement a **strongly coupled EnKF** (Tondeur *et al* 2020).



a) 2-layer atmosphere coupled to 1-layer ocean configuration
b) includes friction at boundaries for atmosphere and wind stress at A-O boundary

*Mechanical component*

*MAOOAM*

mechanical coupling

thermal + radiative coupling

a) includes thermal and radiative heat transport between atmosphere and ocean as function of $T_a$ and $T_o$.

*Thermal component*

## Proposed approach

Simplified description of the algorithm:

**1** Estimating the state $\mathbf{x}^{\mathrm{a}}_{1:K}$. At each time $t_k$, we calculate a forecast $\mathbf{x}^{\mathrm{f}}$:

$$\mathbf{x}^{\mathrm{f}}_{k+1} = \mathbf{x}^{\mathrm{f}}(t_k + \Delta t) = (\mathcal{M}^{\varphi})^{N_c}(\mathbf{x}^{\mathrm{a}}_k)$$

An observation $\mathbf{y}_{k+1}$ is assimilated with **strongly coupled EnKF** to produce an analysis $\mathbf{x}^{\mathrm{a}}_{k+1}$

**2** Determining an estimation of the unknown part of the model. We assume that:

- $\mathbf{x}(t + \Delta t) \approx (\mathcal{M}^{\varphi})^{N_c}(\mathbf{x}(t)) + N_c \cdot \mathcal{M}^{\mathrm{UN}}[\mathbf{x}(t)]$
- $\mathbf{x}(t) \approx \mathbf{x}^{\mathrm{a}}(t)$

We consider that $\mathcal{M}^{\mathrm{UN}}(\mathbf{x}_k) \approx \mathbf{z}_{k+1} = 1/N_c \cdot \left(\mathbf{x}^{\mathrm{a}}_{k+1} - \mathbf{x}^{\mathrm{f}}_{k+1}\right) \implies$ The "target" (*i.e.* the model error) is estimated using the *analysis increments* (Carrassi and Vannitsem, 2011).

**3** Training a neural network $\mathcal{M}_{\boldsymbol{\theta}}$ by minimising the loss $L(\boldsymbol{\theta}) = \sum_{k=0}^{K-1} ||\mathcal{M}_{\boldsymbol{\theta}}(\mathbf{x}^{\mathrm{a}}_k) - \mathbf{z}_{k+1}||^2$

**4** Using the hybrid model $\mathcal{M}^{\varphi} + \mathcal{M}_{\boldsymbol{\theta}}$ to produce new simulations (*e.g.* to make forecasts).

## Experiments with MAOOAM

**①** **Truth**: $n_a = 20$ and $n_o = 8$ modes for atmosphere and ocean. **Total dimension** $N_x = 56$.

**②** **Truncated**: $n_a = 10$ and $n_o = 8$ modes for atmosphere and ocean. **Total dimension** $N_x = 36$.

▶ The truncated model is **missing 20 high-order atmospheric variables**

▶ There is not locality in spectral space so the NN is made of 3 layers multi-layer perceptrons

### RMSE-f of hybrid and truncated MAOOAM models

| RMSE-f(lead time $\tau$) | $\psi_{o,2}$(2 years) | $\theta_{o,2}$(2 years) | $\psi_{a,1}$(1 day) |
|---|---|---|---|
| Truncated | 0.23 | 0.21 | 0.36 |
| **Coupled DA-ML hybrid** | **0.10** | **0.06** | **0.28** |

- The hybrid models have superior skill than the truncated model.

- The improvement is larger for the ocean that is fully resolved $\implies$ **Enhanced representation of the atmosphere-ocean coupling processes thanks to coupled DA**.

- The atmosphere is improved less: the hybrid is not very good in representing the fast processes.
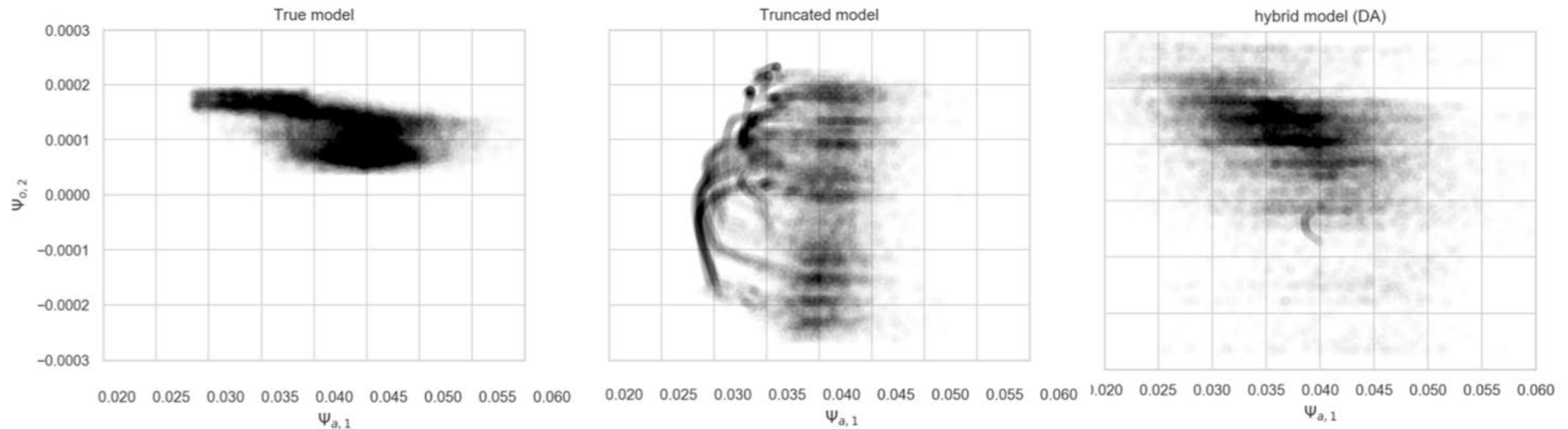
## Numerical experiments: atmosphere-ocean model MAOOAM

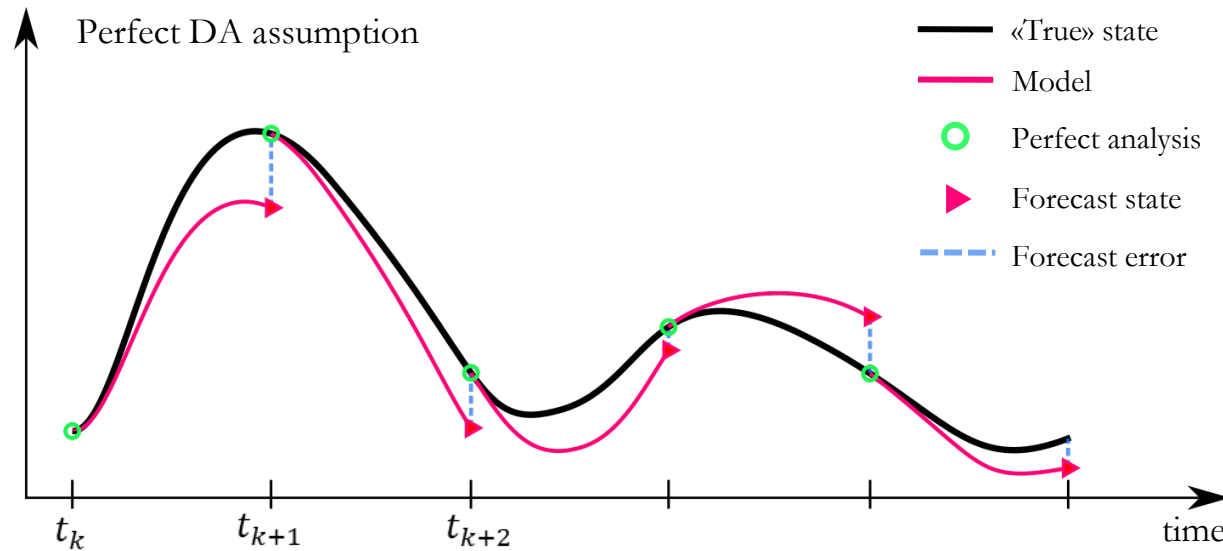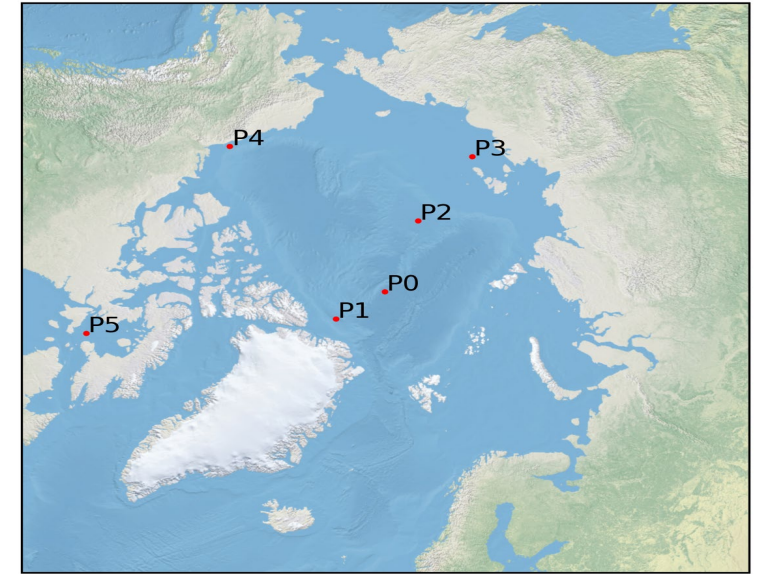### Reconstruction of the model attractor



▶ The truncated model visits areas of the phase space that are not admitted in the real dynamics.

▶ Discrepancies are reduced by the hybrid models.

# Hybrid sea ice modelling: ML-ICEPACK



Perfect DA assumption

- «True» state
- Model
- ○ Perfect analysis
- ▶ Forecast state
- ╌ Forecast error

- 6 locations simulated
- 15 years atmospheric forcing (1993 - 2007)

➢ 30 different configurations were generated sampling the parameters independently from the distributions in the table.

➢ One training dataset per configuration (30 in total).

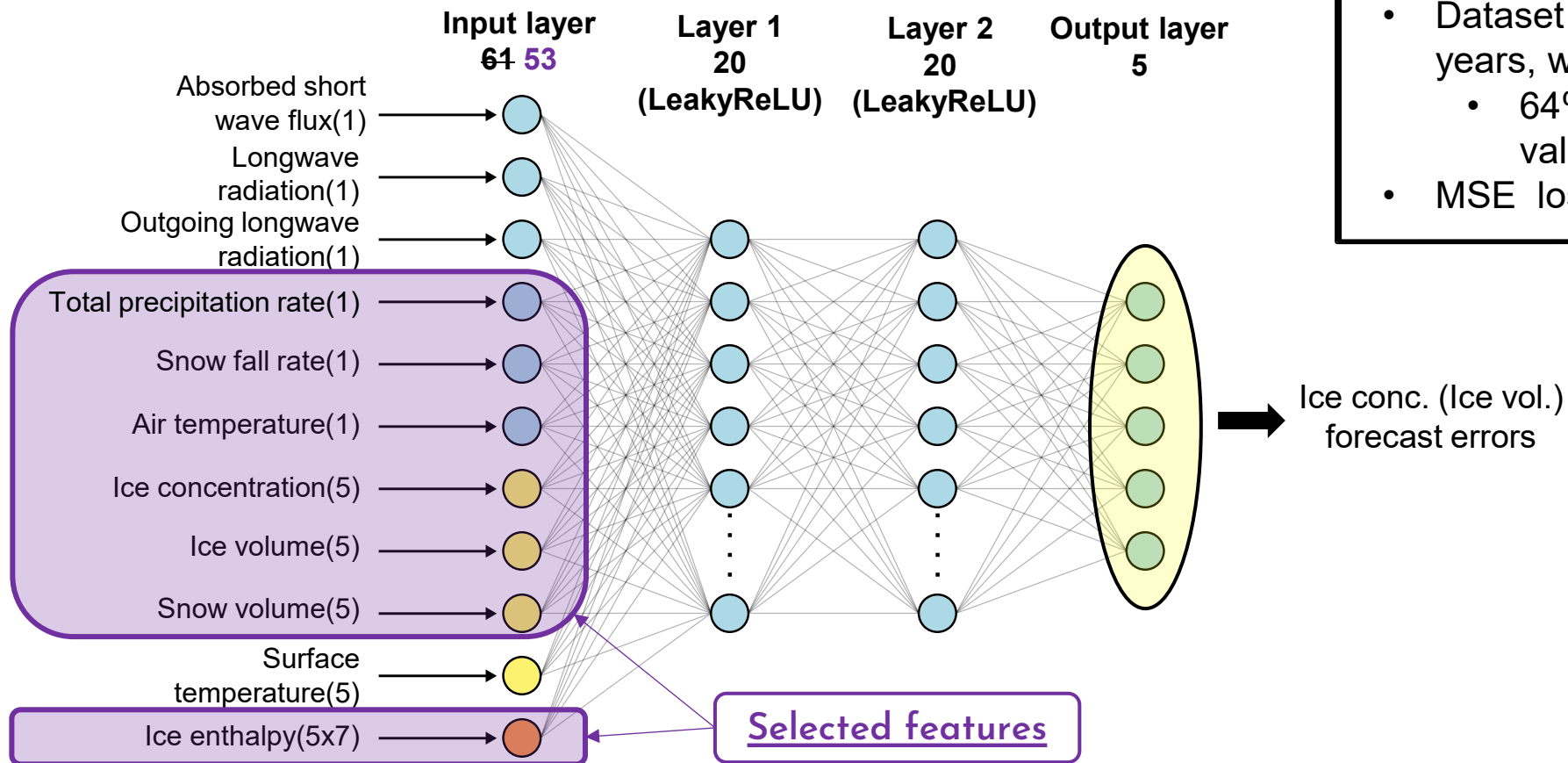➢ Atmospheric forcing perturbations introduced in the execution of the hybrid models

## Parametric error in the snow thermodynamic and radiative properties

| Parameter symbol | Parameter description | Default (true) value | Distribution | Min value | Max value | Mode |
|---|---|---|---|---|---|---|
| ksno | Thermal conductivity of snow (W m$^{-1}$K$^{-1}$) | 0.3 | Uniform | 0.03 | 0.65 | |
| rsnw_mlt | Max. melting snow grain size (µm) | 1500 | Triangular | 250 | 3000 | 1500 |

See Urrego-Blanco et al. 2016

*Estimating and correcting model error in ICEPACK*

ML MODELS

FEATURE SELECTION: RECURSIVE FEATURE ELIMINATION

Input layer
61 53

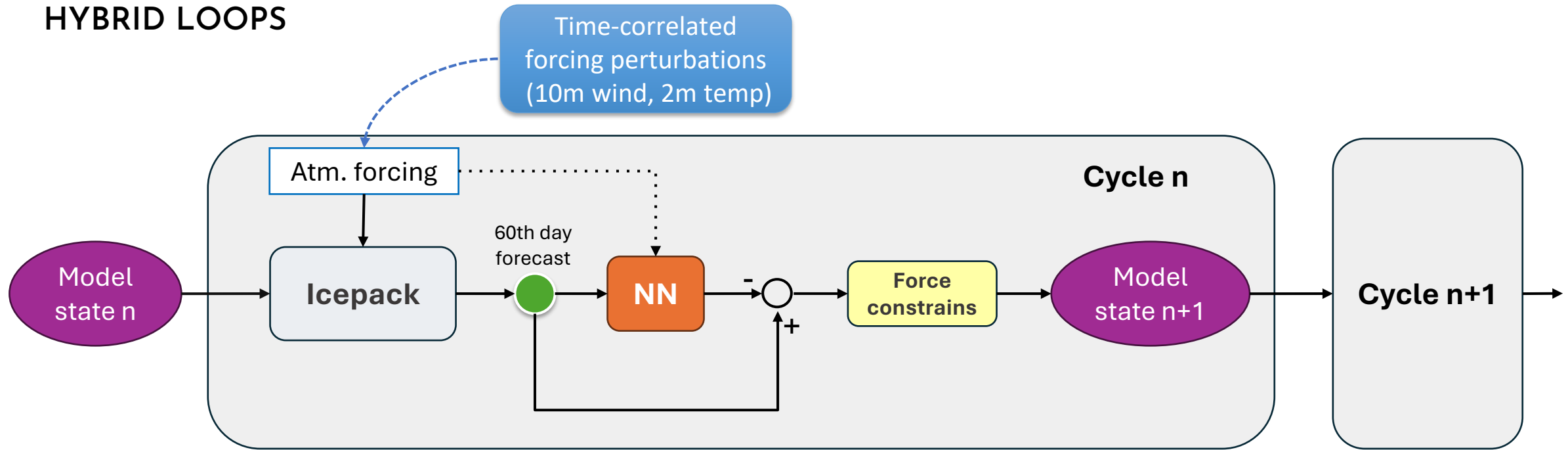Layer 1
20
(LeakyReLU)

Layer 2
20
(LeakyReLU)

Output layer
5

Absorbed short wave flux(1)
Longwave radiation(1)
Outgoing longwave radiation(1)
Total precipitation rate(1)
Snow fall rate(1)
Air temperature(1)
Ice concentration(5)
Ice volume(5)
Snow volume(5)
Surface temperature(5)
Ice enthalpy(5x7)

Selected features

Ice conc. (Ice vol.) forecast errors

- 2 MLP for each member (for ice conc and ice volume with the same architecture)
- Dataset size: 3132 instances (10 years, weekly frequency, 6 locations)
  - 64% training (2005), 16% validation (501), 20% test (626)
- MSE loss

*De Cillis et al., 2025 (in preparation)*

**HYBRID LOOPS**



➢ Although trained for parametric error, the hybrid model is used also in the presence of error in the external forcing

➢ In the experiments we compare our NN-based corrected model with a climatological correction WCLIM
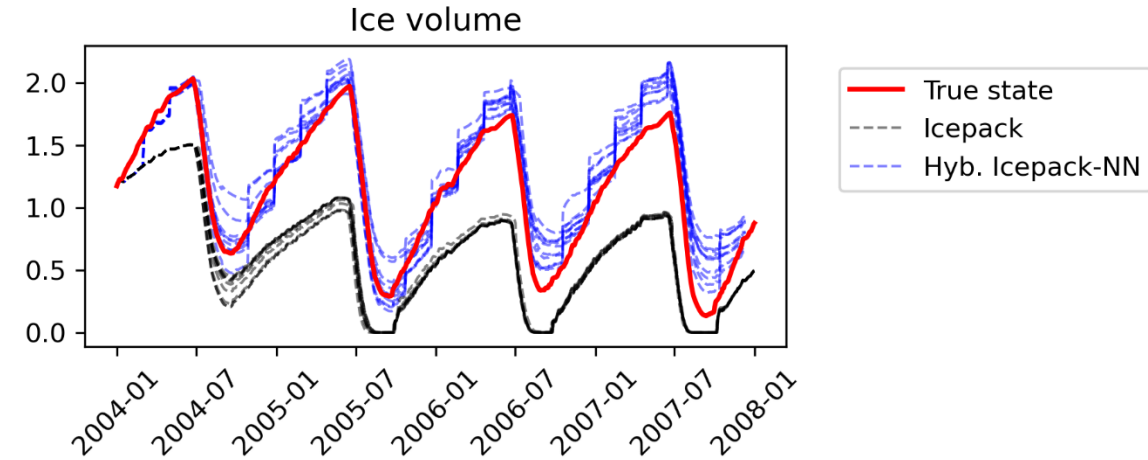
Ice volume evolution on on-line test run with <u>atmospheric forcing perturbations</u>
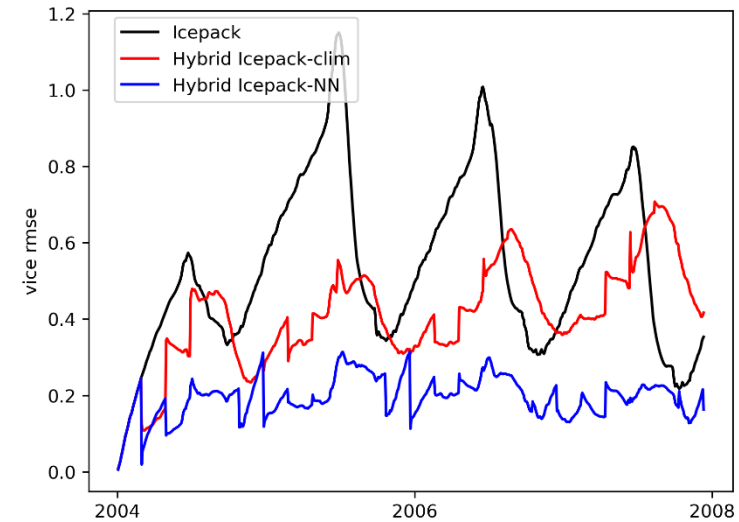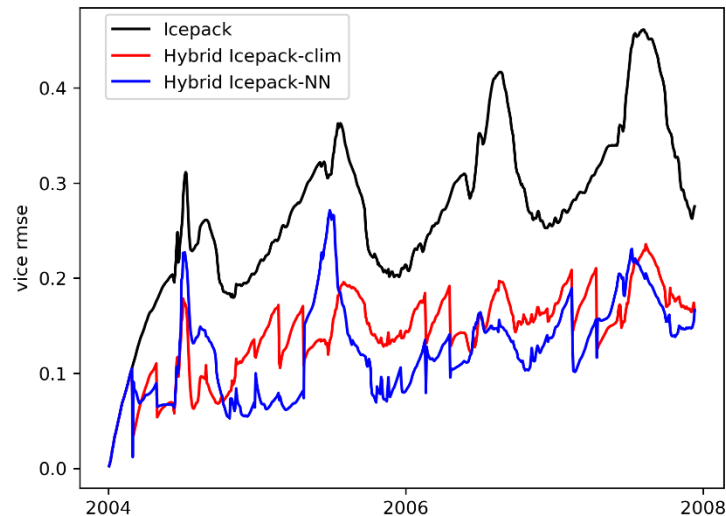
**The  original model overestimates**          At P0 (the pole)          **The original model underestimates**
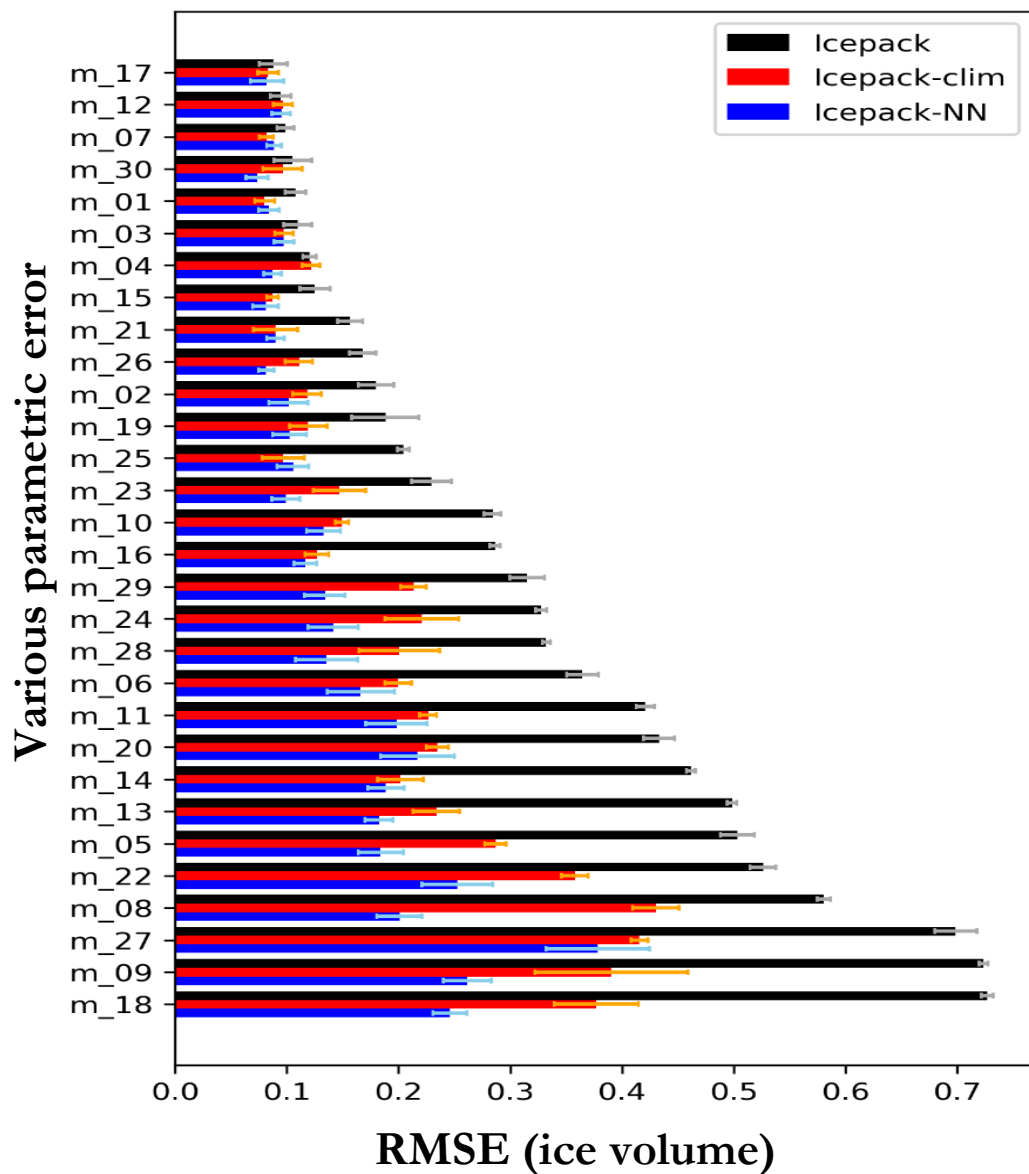


RMSE (average over locations and ensemble members with perturbed forcing)
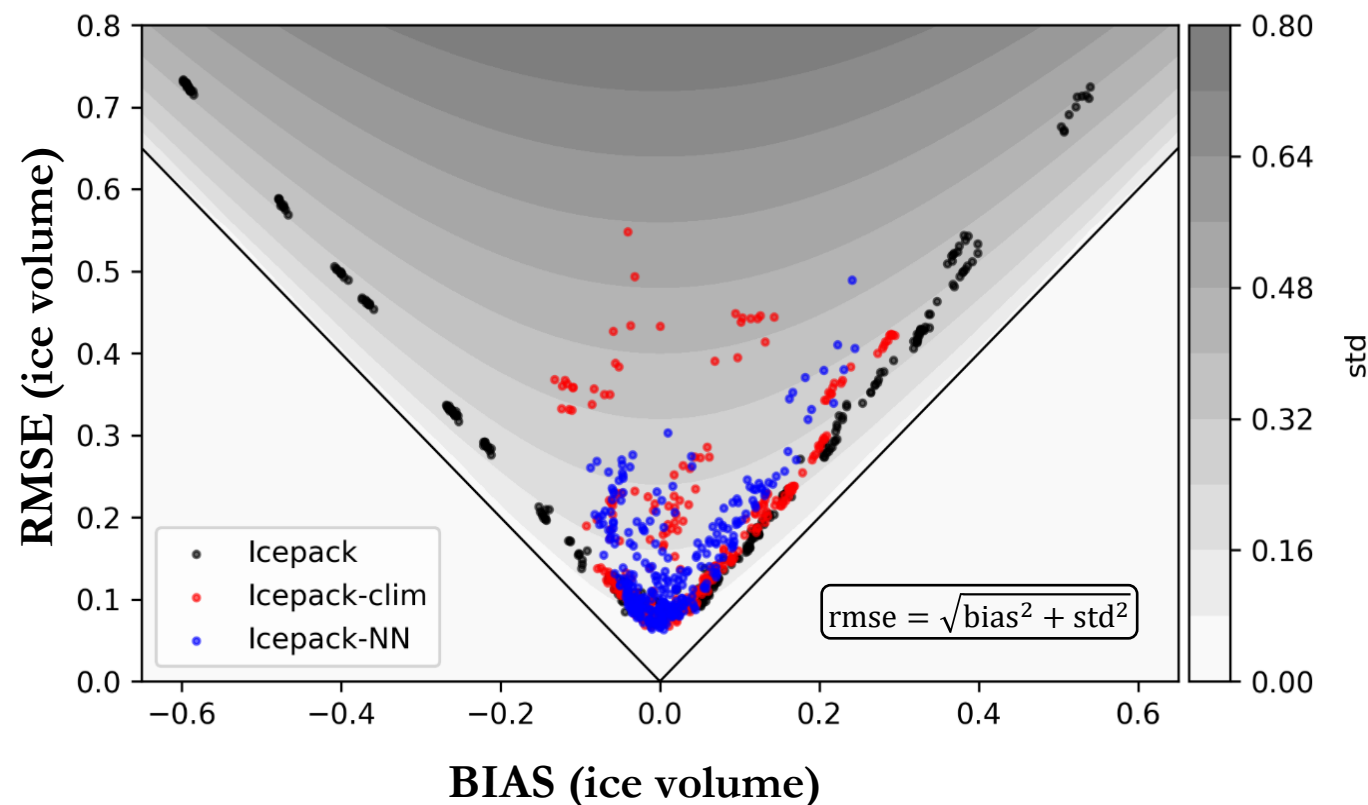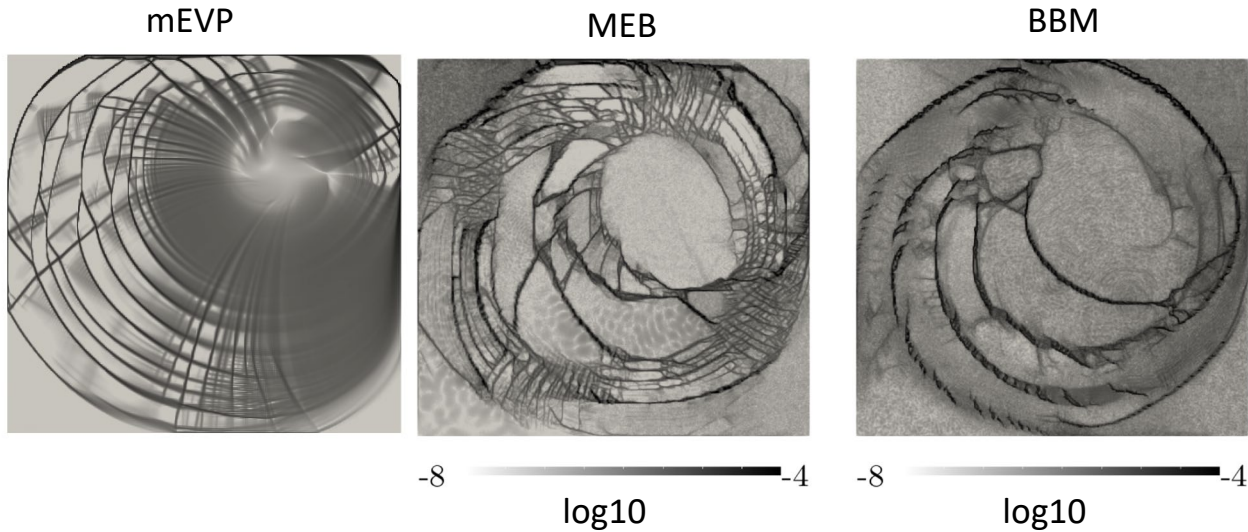
*De Cillis et al., 2025 (in preparation)*

➢ NN performs well despite atmospheric forcing perturbations
➢ Reduced bias and RMSE compared to the uncorrected (free run) and climatology
➢ The NN works better when there are large errors to correct
➢ Larger residual errors are observed for models with positive bias

**Features in sea ice models are non-Gaussian and constraint to submanifolds**

mEVP MEB BBM



-8    -4      -8    -4
log10       log10

- A new sea ice model is being developed in SASIP.
- It will incorporate the new Maxwell-elasto-brittle (MEB) and Bingham-brittle (BBM) rheology. **Contrary to the elasto-viscous (mEVP) rheology, these rheologies have a memory**.
- Brittle rheologies produce more realistic structures and respect scaling laws observed in the sea ice fields.

*Shear deformation after 2 days for different rheologies and orders of basis polynomials (Richter et al., 2023)*

➤ **Errors in sea ice models are non-Gaussian and constraint to submanifolds**. This violates EnKF assumptions.

➤ Variational autoencoder (VAE) provides encoder $p_\theta(z|x)$ and decoder $q_\varphi(x|z)$ linking an arbitrary distribution $p(x)$ in state space to the standard normal $N(z; 0, \mathbf{I})$ in latent space.

- Variational autoencoder (VAE) tries to link an arbitrary $p_{\mathcal{X}}(\mathbf{x})$ distribution to the standard normal $\mathcal{N}(\mathbf{z}; \mathbf{0}, \mathbf{I})$ via an encoder and decoder:

- Encoder:

$$\mathcal{N}(\mathbf{z}; \mathbf{0}, \mathbf{I}) = \int \boxed{q_\phi(\mathbf{z}|\mathbf{x})} p_{\mathcal{X}}(\mathbf{x}) \, \mathbf{dx}$$

Decoder:

$$p_{\mathcal{X}}(\mathbf{x}) = \int \boxed{p_\theta(\mathbf{x}|\mathbf{z})} \mathcal{N}(\mathbf{z}; \mathbf{0}, \mathbf{I}) \, \mathbf{dz}$$
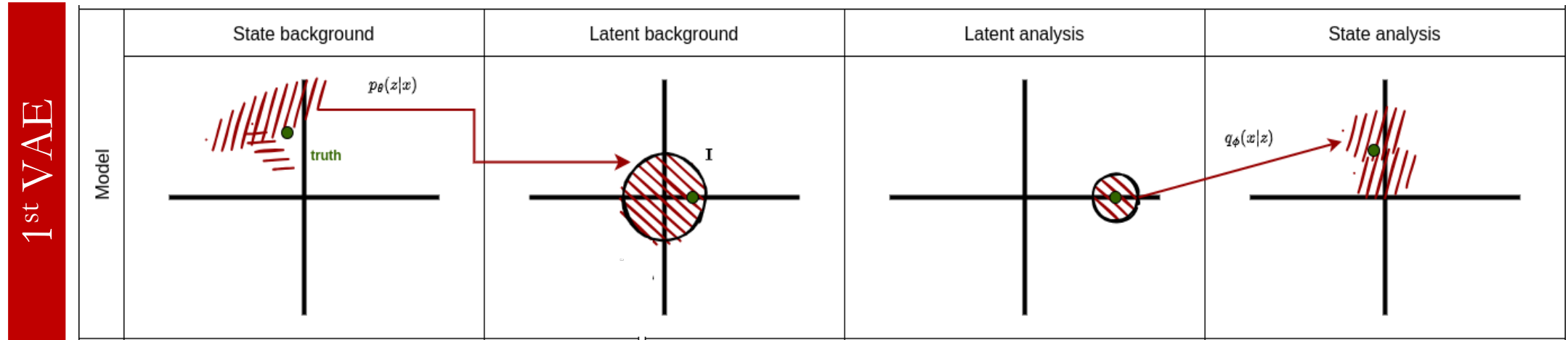
and $\theta, \phi$ are the weights in a neural network.

- $\theta, \phi$ can be found be maximizing the ELBO

$$\mathcal{L}(\theta, \phi, \mathbf{x}) = \int q_\phi(\mathbf{z}|\mathbf{x}) \log p_\theta(\mathbf{x}|\mathbf{z}) \, \mathbf{dz} - \int q_\phi(\mathbf{z}|\mathbf{x}) \log \frac{q_\phi(\mathbf{z}|\mathbf{x})}{\mathcal{N}(\mathbf{z}; \mathbf{0}, \mathbf{I})} \, \mathbf{dz}$$

- **Hypothesis:** better adherence to Gaussian assumptions can be achieved by correcting the ensemble in latent space instead of physical space.
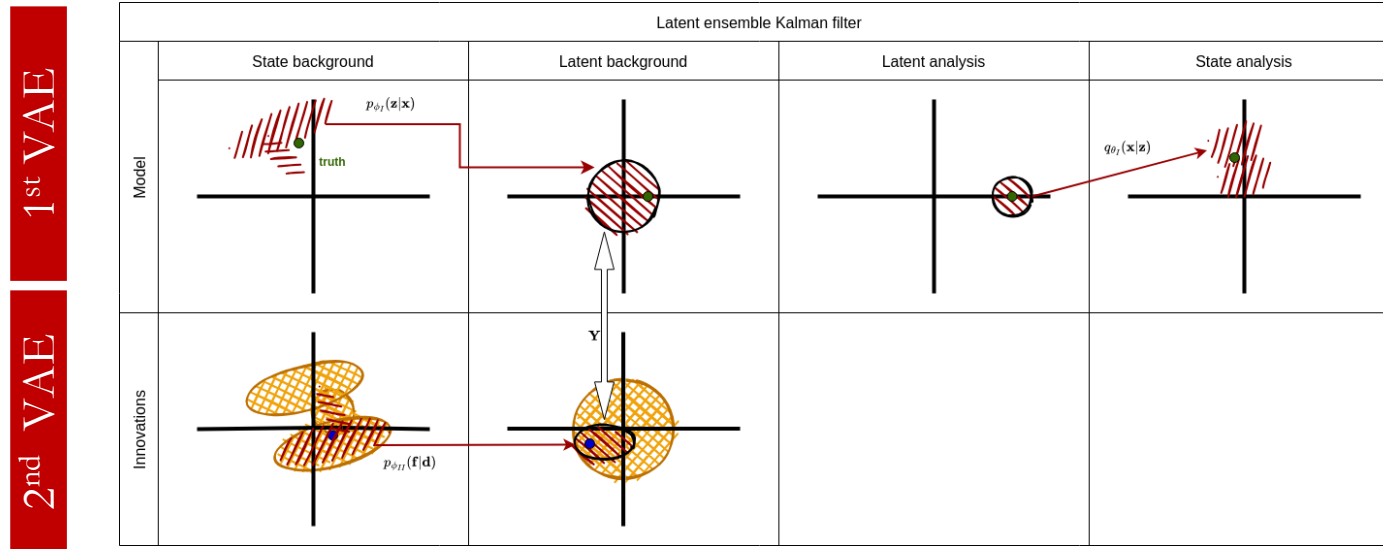
- (Re)train 1st *VAE*.
- Use the **encoder** to sample one ensemble member in latent space for each ensemble member in state space replacing $\mathbf{X}^f$ with $\mathbf{Z}^f$.
- Apply *ETKF* to the ensemble in latent space.
- Use the **decoder** to sample one ensemble member in physical state space for each ensemble member in the latent space replacing $\mathbf{Z}^a$ with $\mathbf{X}^a$.
- Run the ensemble forward in time using a physical model.

| Forward propagation | DA method | States | Observations | Gaussian | Non-Gaussian |
|---|---|---|---|---|---|
| Physical model | ETKF | Latent space 1st VAE | Physical space | $\mathbf{Z}^f, \mathbf{Z}^a$ | $\mathbf{X}^f, \mathbf{X}^a, \mathbf{Y}, \mathbf{D}$ |

# …a *double VAE-ETKF* configuration
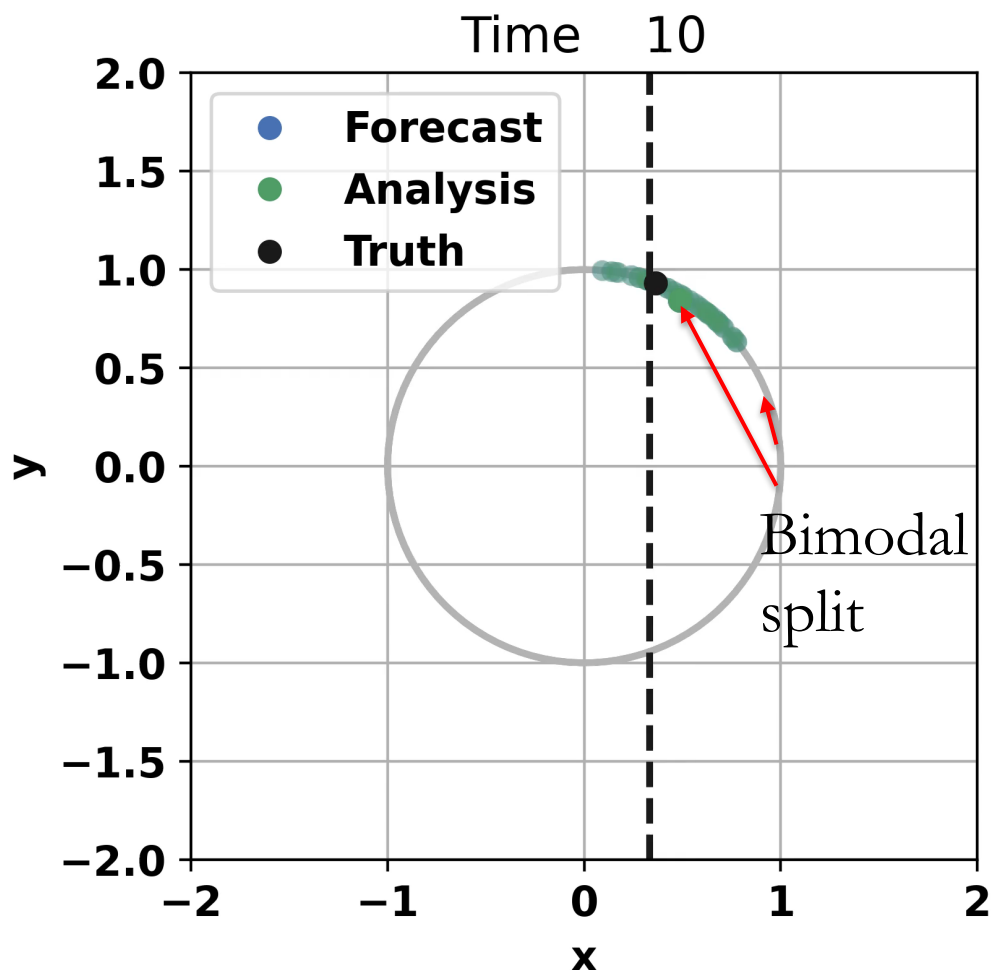


Pasmans et al., 2025

- Create the ensemble in the 1st latent space as described on the previous slide.
- Create artificial innovations using bootstrap. Train 2nd VAE on these innovations.
- Create innovations that match observed values by adding observational errors. Use 2nd encoder to sample expanded ensemble **D** in latent space from this.
- Use **2nd encoder** to sample one innovation in the 2nd latent space from each unperturbed innovation in the physical space. Use this to construct **Y**.
- Carry out ETKF and follow steps on the previous slide.

| Forward propagation | DA method | States | Observations | Gaussian | Non-Gaussian |
|---|---|---|---|---|---|
| Physical model | ETKF | Latent space 1st VAE | Latent space 2nd VAE | $\mathbf{Z}^f, \mathbf{Z}^a, \mathbf{Y}, \mathbf{D}$ | $\mathbf{X}^f, \mathbf{X}^a$ |

# Dynamical model

- VAE-DA is tested using a rotating point.
- Difficult for classical *ETKF:*
  **discontinuity for angle=0**
  **solution restricted to manifold.**
- Optional change in radius controlled by $\boldsymbol{\alpha_r}$.

- *x*-coordinate assimilated every 10 timesteps.
- Observational error is realization Gaussian with a standard deviation of 0.1.

$$\frac{\mathrm{d}\theta}{\mathrm{d}t} = \alpha_\theta \theta(t) \text{ with } 0 \le \theta < 2\pi$$

$$\frac{\mathrm{d}r}{\mathrm{d}t} = \alpha_r \frac{2\pi}{50} \cos(\frac{2\pi t}{50})$$



*(light green) 64-member ensemble together with the (dark green) ensemble mean, (black dot) truth and (dashed black line) imperfect observation x-coordinate.*
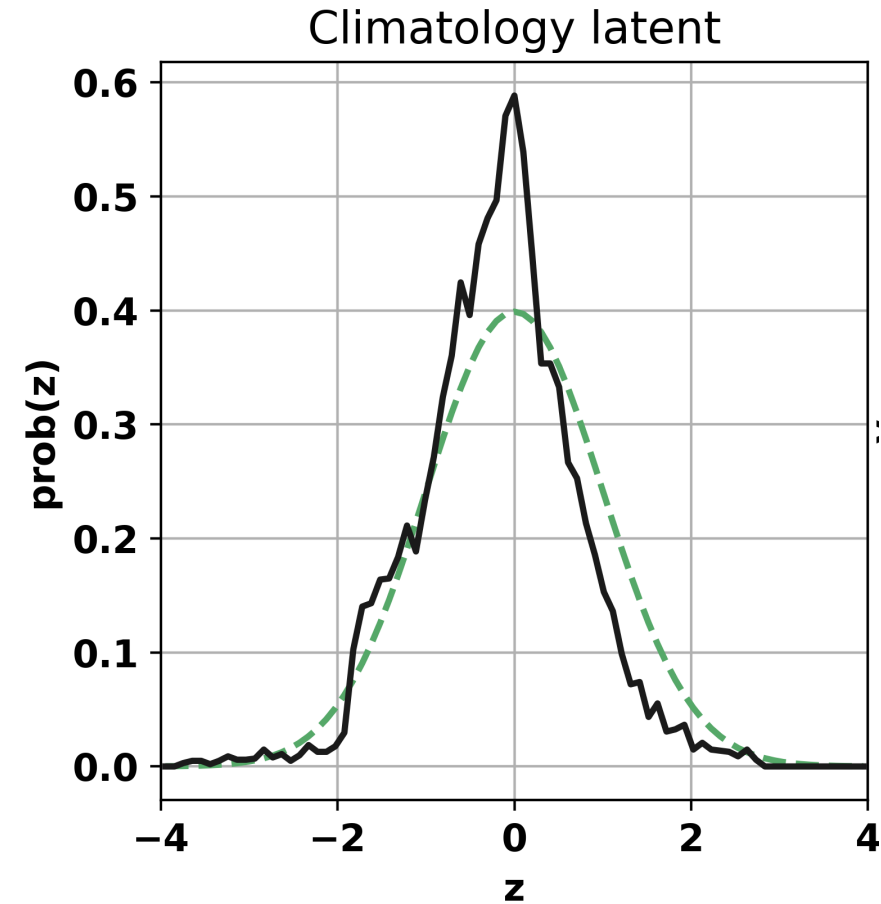
- VAE trained on every 10th sample of 10,000 long model run.
- VAE succeeds in capturing the submanifold.

Reconstruction state space

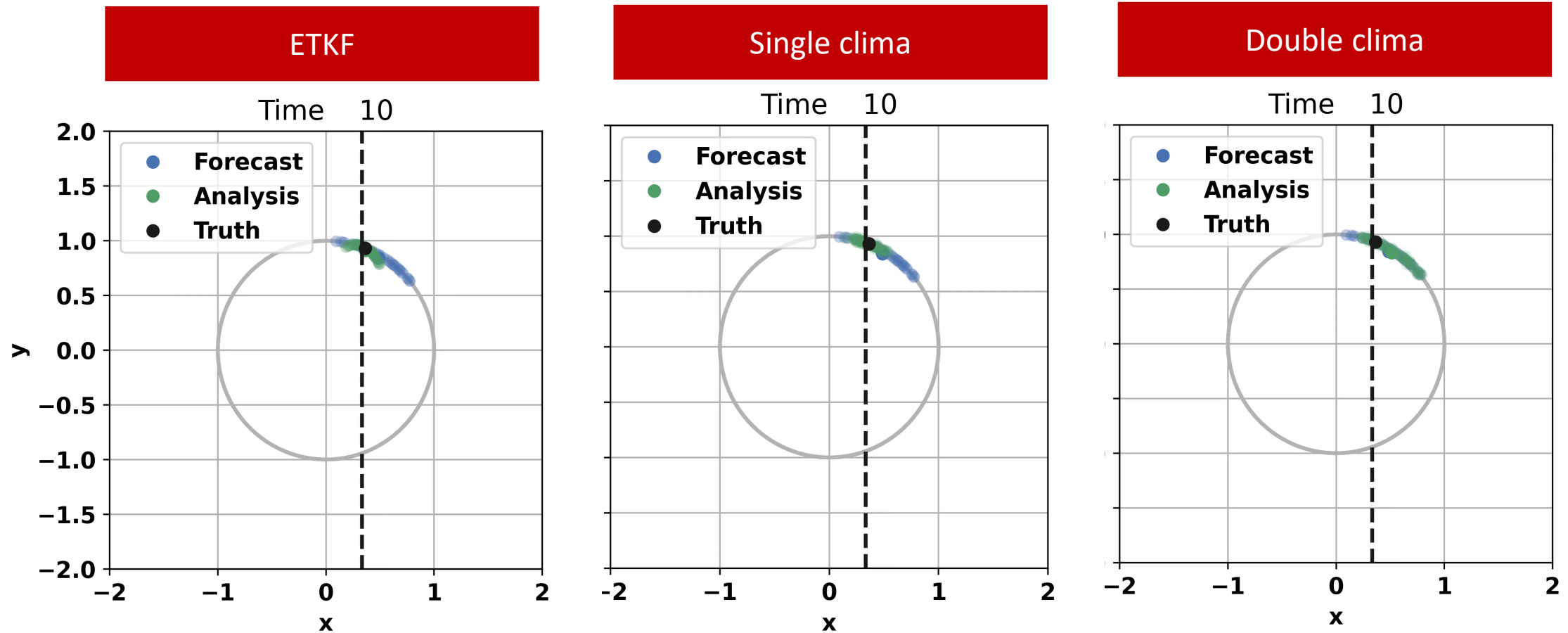Distribution latent space



Climatology state

Climatology latent

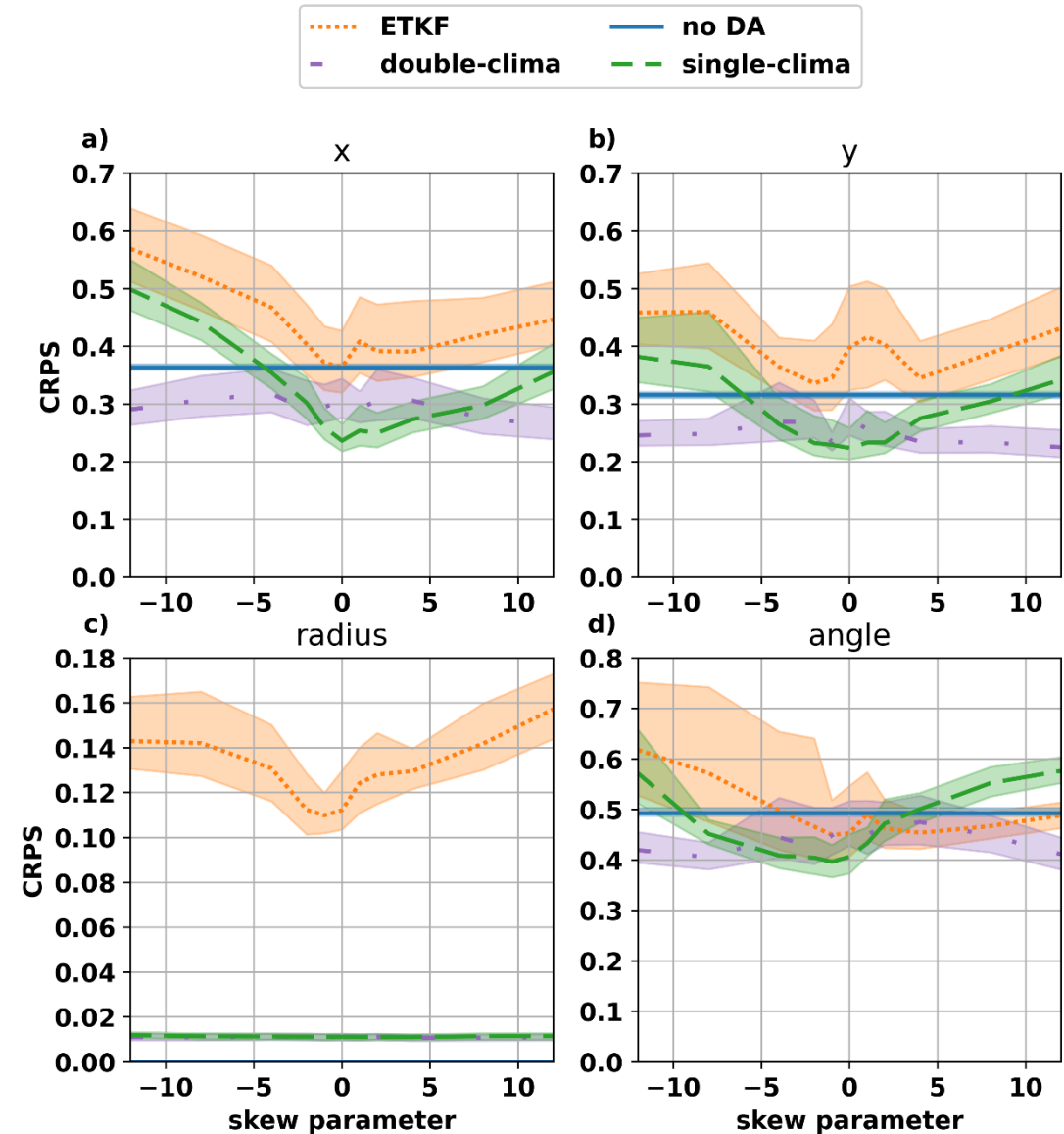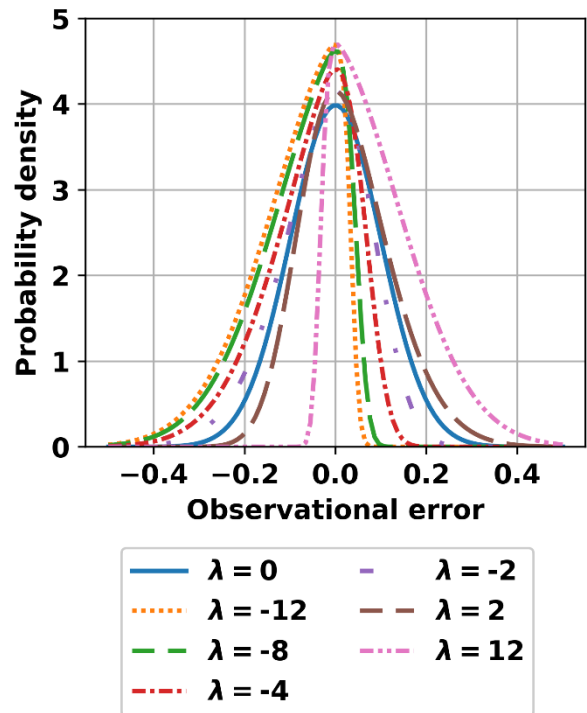# Experiment I: stationary climate

**Improvement *VAE-DA* over *ETKF* is caused by the decoder <u>constraining the ensemble members onto the circle</u>.**



*64-member ensemble just (light blue) before and (light green) after DA together with the (dark) ensemble mean, (black dot) truth and (dashed black line) imperfect observation x-coordinate in experiment I.*
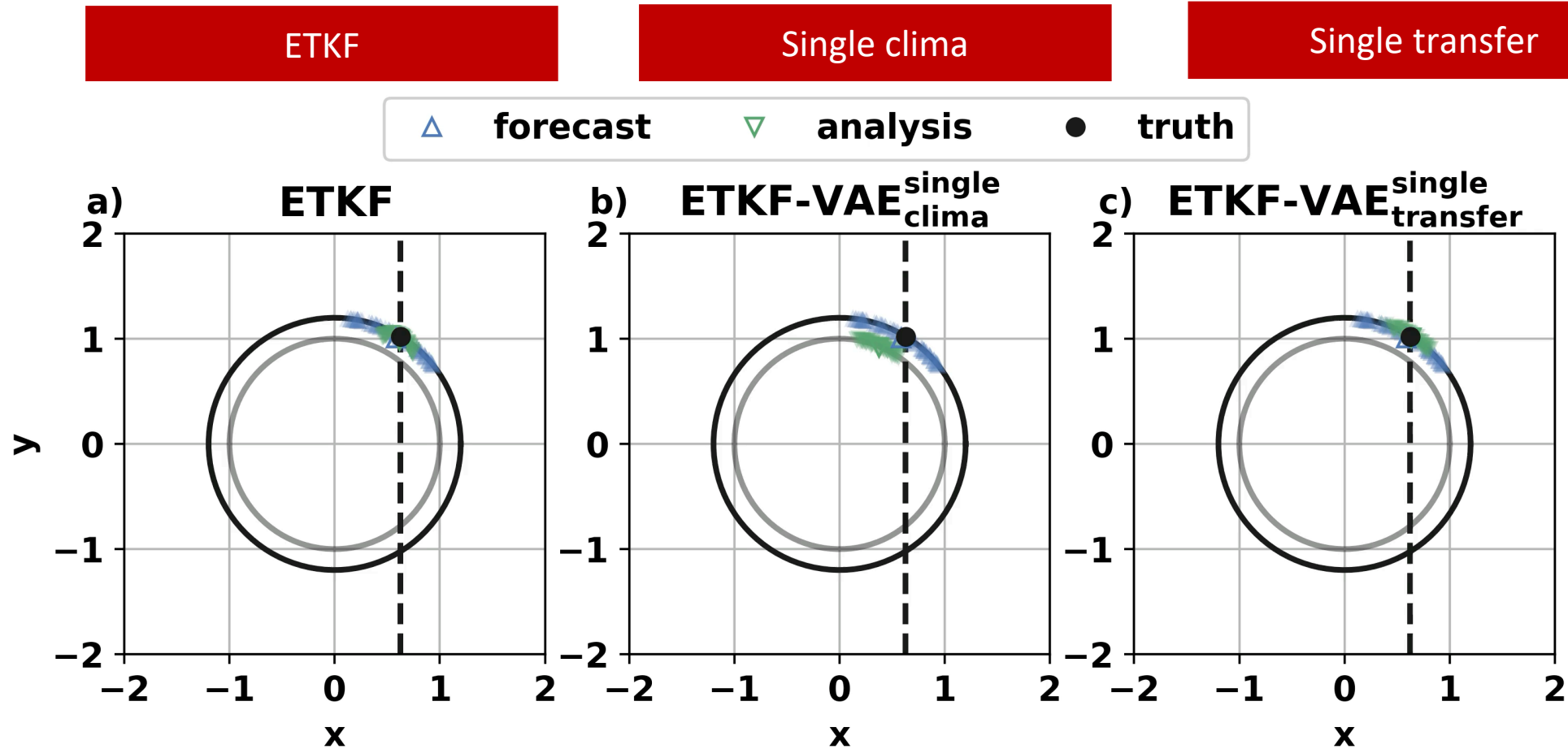
- ➢ Experiments with non-Gaussian observation error

- ➢ The performance is measured using the **Continuous Rank Probability Score** (CRPS)

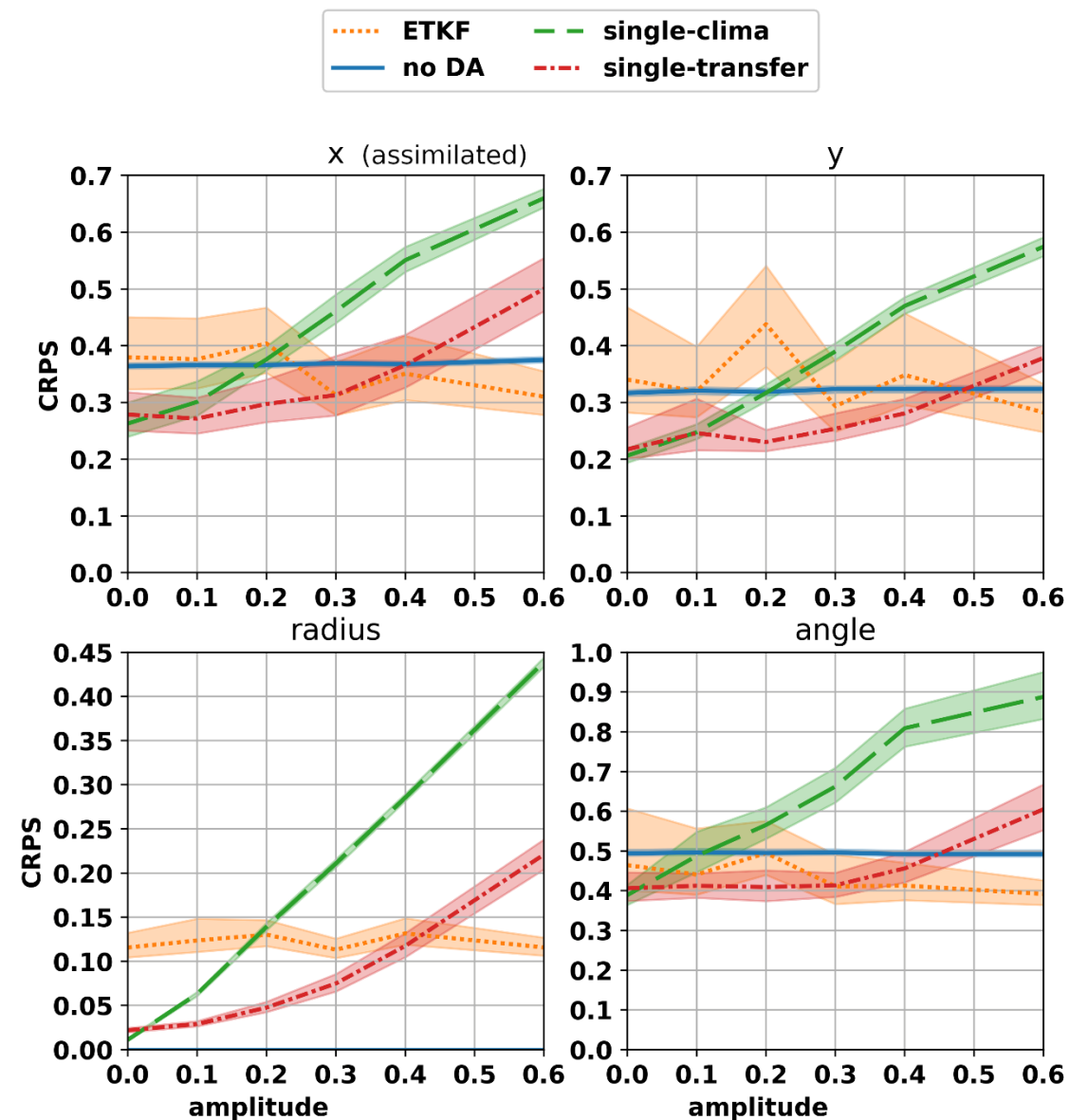- ➢ Double-clima outperforms when the deviation from Gaussianity becomes larger

- $\alpha_r = 0.2$ so truth oscillates around the unit circle.
- *VAE-DA* configurations with online updating (*transfer*) capture this, those without (*clima*) map the ensemble back to unit circle.

ETKF

Single clima

Single transfer

△ **forecast**   ▽ **analysis**   ● **truth**



64-member ensemble just (light blue) before and (light green) after DA together with the (dark) ensemble mean, (black dot) truth and (dashed black line) imperfect observation x-coordinate in experiment II.
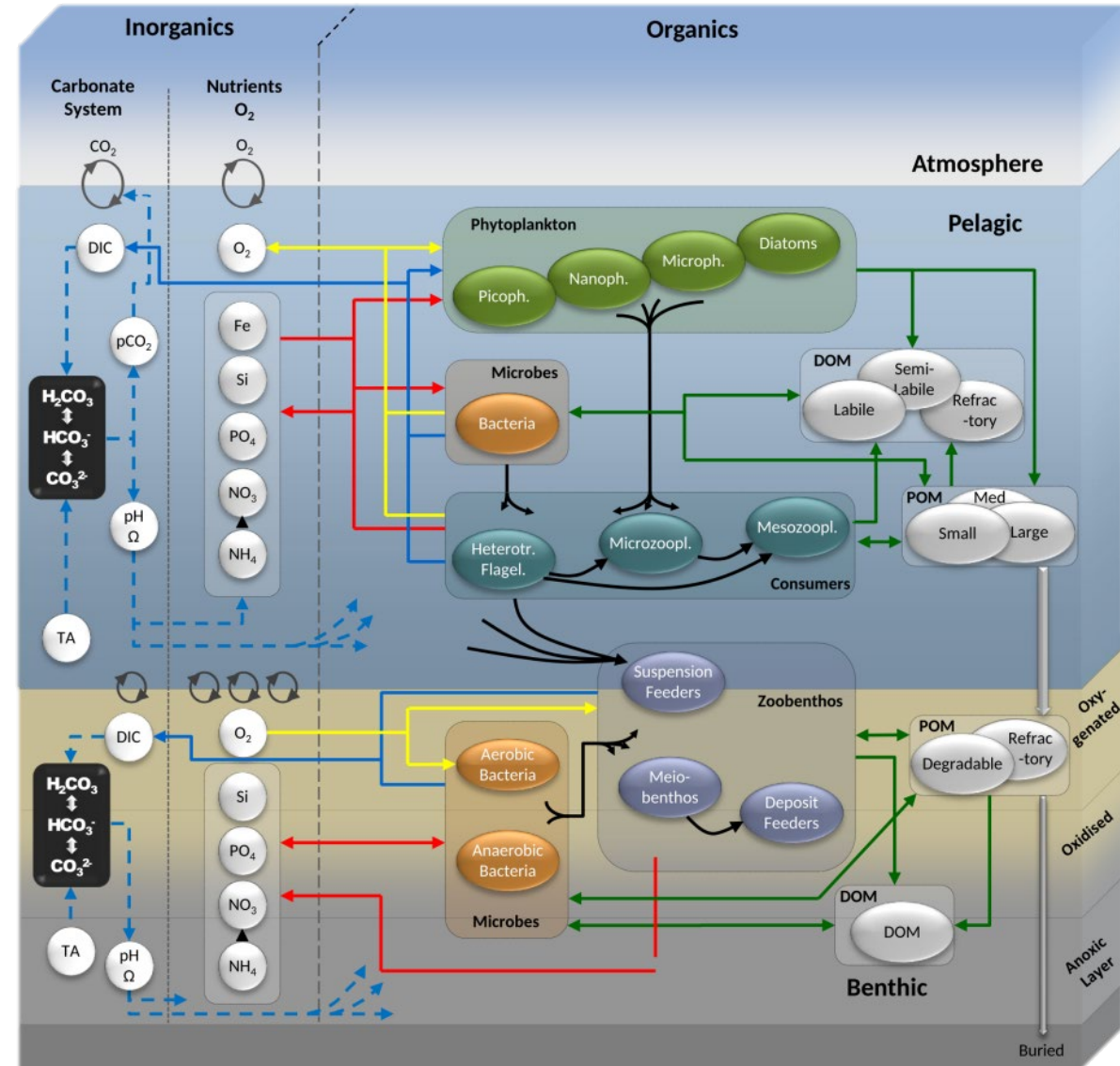
> - Experiments with varying-in-time circle
>
> - The performance is measured using the **Continuous Rank Probability Score** (CRPS)
>
> - Single-transfer outperforms when the oscillation around a steady circle have larger amplitude

Pasmans et al., 2025

# Domain considerations

- Modelling this is a complex task

- Large system of interactions coupled to a physical model

- Huge uncertainties in model parameters & feedback mechanisms

- Observation Space <<<< Model Space (only observe total chlorophyll, derived from ocean color at surface)

# Hybrid DA-ML for ocean biogeochemistry applications

➢ We need to to describe accurately the (potentially non-linear) relationship between the few observations and the rest of the unseen domain

➢ Currently, a univariate EnKF + a stochiometric balancing scheme is adopted operationally (e.g. at MetOfficeUK)

➢ We aim at using machine learning to unveil the observed to unobserved relationship, and use it effectively to update the unobserved space
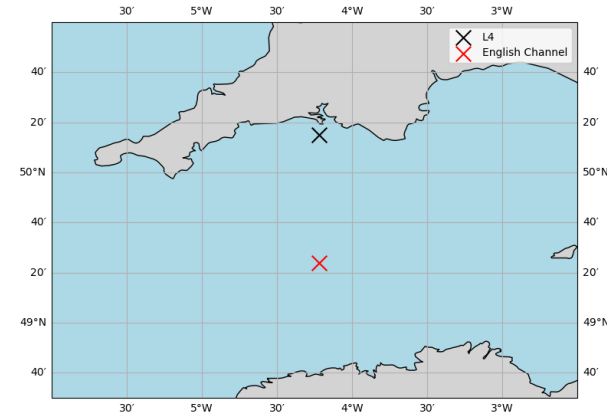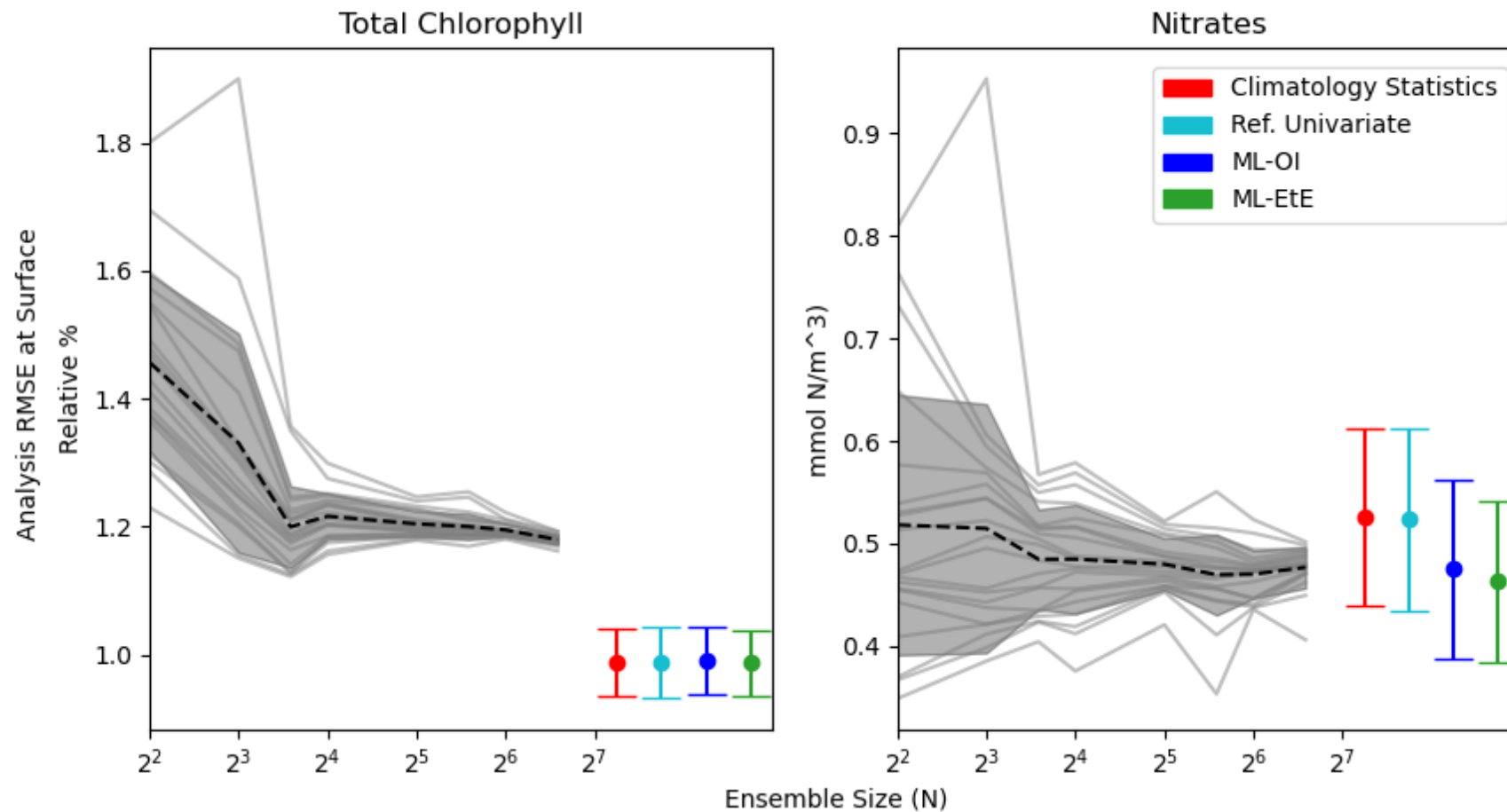
Learning the cross-covariance: ML-OI

$$x_v^a = x_v^f + \frac{P_{VC}^f}{P_{CC}^f + R} \cdot (y - x_C^f)$$
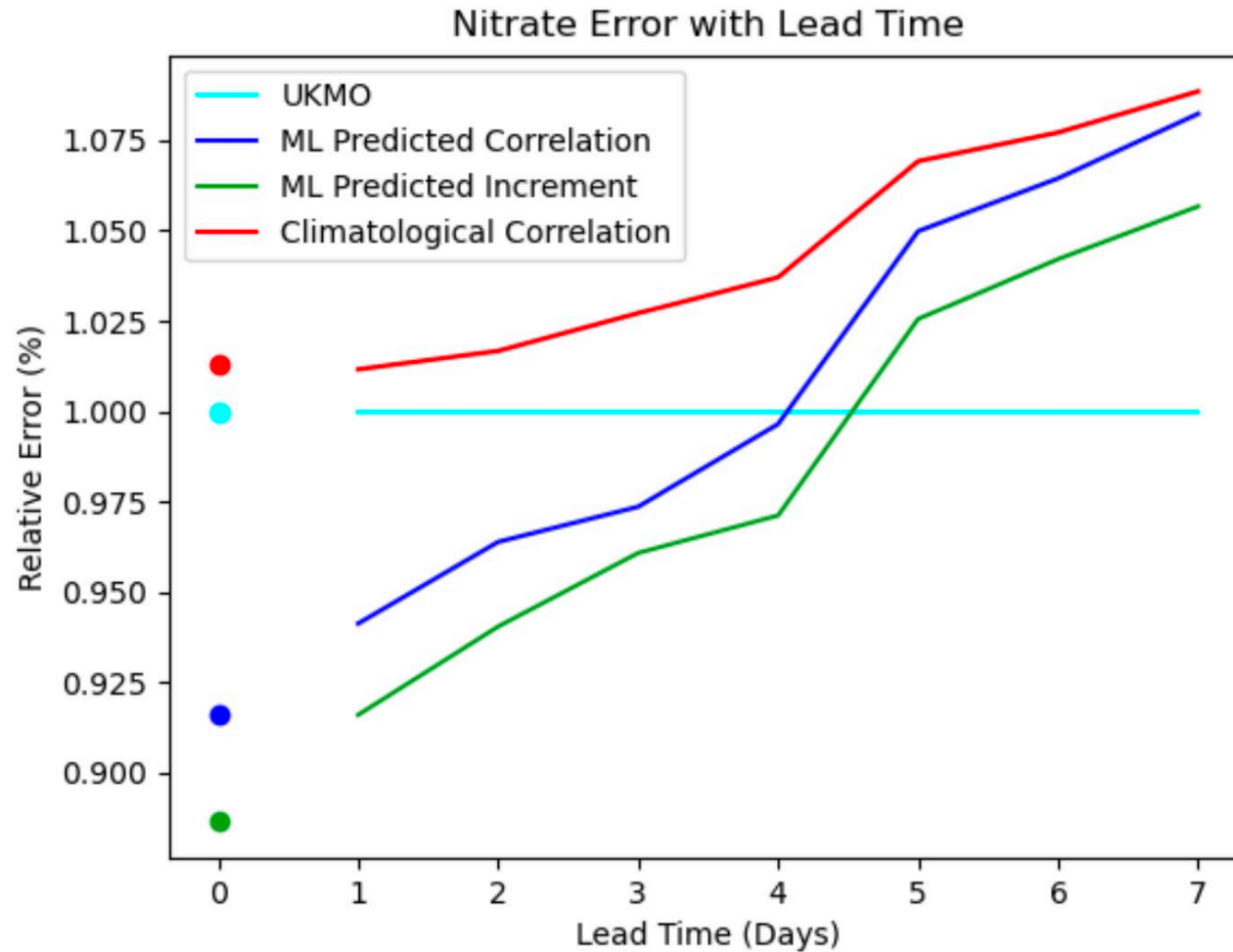
Learning the full analysis update: ML-EtE

Higgs et al., 2025

Performance at analysis

Higgs et al., 2025

Forecasting performance



Higgs et al., 2025

# Recap & looking ahead

✅ **DA and ML can be combined in different ways: here we showed several examples**

✅ **ML can substitute some (possibly more nonlinear/non-Gaussian) steps in the DA process: learning the cross-covariance; DA with VAE.**

✅ **ML can be used to learn and correct for model error based on analysis from DA**

✅ **A lot of what done needs to be studied on larger dimension...**

⚙️ Several ongoing directions: scale-up applications, coupled chaotic dynamics, full end-to-end....