Learning surrogate models and data assimilation processes for advanced geophysical dynamics forecasting

Marc Bocquet,*

Alban Farchi, †,* Tobias Finn, * Charlotte Durand, * Sibo Cheng, * Wenbo Yu, * Yumeng Chen, $^{\parallel}$ Ivo Pasmans, $^{\parallel}$ Alberto Carrassi, §

* CEREA, ENPC, EDF R&D, Institut Polytechnique de Paris, Île-de-France, France

Department of Meteorology and National Centre for Earth Observation, University of Reading, United-Kingdom

§ Department of Physics and Astronomy, University of Bologna, Italy

† ECMWF, Reading, United Kingdom,

[‡] NERSC, Bergen, Norway,











Outline

- Learning data assimilation from artificial intelligence: first results
 - Sequential data assimilation for chaotic dynamics
 - Learning data assimilation
 - Sensitivity analysis and first results
 - Investigation and interpretation
- Learning data assimilation from artificial intelligence: advanced results
 - Scalability and locality
 - Learning the analysis operator from observations only
 - Learning nonlinearity
- Conclusions

Outline

- Learning data assimilation from artificial intelligence: first results
 - Sequential data assimilation for chaotic dynamics
 - Learning data assimilation
 - Sensitivity analysis and first results
 - Investigation and interpretation
- Learning data assimilation from artificial intelligence: advanced results
 - Scalability and locality
 - Learning the analysis operator from observations only
 - Learning nonlinearity
- Conclusions

Sequential data assimilation for chaotic dynamics

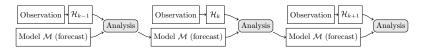
▶ Here, data assimilation (DA) methods are formulated from

$$\mathbf{x}_{k+1} = \mathcal{M}(\mathbf{x}_k),\tag{1a}$$

$$\mathbf{y}_k = \mathcal{H}_k(\mathbf{x}_k) + \boldsymbol{\varepsilon}_k, \qquad \boldsymbol{\varepsilon}_k \sim N(\mathbf{0}, \mathbf{R}_k),$$
 (1b)

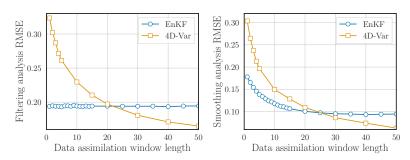
where \mathcal{M} is the *autonomous* evolution model, \mathbf{x}_k is the state vector at time τ_k , \mathbf{y}_k is the observation vector, \mathcal{H}_k is the observation operator, ε_k is the observation error, assumed to be additive, unbiased, white in time, and Gaussian of covariance matrix \mathbf{R}_k .

▶ DA for geofluids has to be *sequential* in time because (i) observations need to be assimilated *as they arrive* to update the state estimation, (ii) applied to *chaotic dynamics*, typical errors have an exponential growth.



The edge of ensemble filtering methods

- ➤ The variational methods (3D–Var, 4D–Var): can handle nonlinearity of the operators, asynchronous observations, but cannot handle the errors of the day.
- ▶ The ensemble filtering methods (EnKFs): can only handle weak nonlinearity of the operators, cannot handle asynchronous observations, can handle the errors of the day through the ensemble.
- ▶ Testing the EnKF ($N_{\rm e}=20$), 4D–Var, and IEnKS ($N_{\rm e}=20$) variants with the chaotic 40–variable Lorenz 96 model [Bocquet et al. 2013; Asch et al. 2016]:



▶ In mild nonlinear regime, the EnKF significantly outperforms the (basic) 4D–Var with moderately large DA windows because it captures the *errors of the day*.

Our focus: learning the analysis

▶ Let us assume that \mathcal{M} is known, that the Jacobian of \mathcal{H}_k is \mathbf{H}_k , and that we wish to learn an incremental analysis operator a_{θ} , typically a neural network parametrised by θ .

▶ If $\mathbf{E}_k^{\mathrm{a}}, \mathbf{E}_k^{\mathrm{f}} \in \mathbb{R}^{N_{\mathrm{x}} \times N_{\mathrm{e}}}$ are the analysis and forecast ensemble matrices at time τ_k , a_{θ} is defined via the (ensemble) update:

$$\mathbf{E}_{k}^{\mathrm{a}} = \mathbf{E}_{k}^{\mathrm{f}} + a_{\theta} \left(\mathbf{E}_{k}^{\mathrm{f}}, \mathbf{H}_{k}^{\mathsf{T}} \mathbf{R}_{k}^{-1} \delta_{k} \right), \tag{2a}$$

where δ_k , the innovation at time τ_k , is defined by

$$\boldsymbol{\delta}_{k} \stackrel{\Delta}{=} \mathbf{y}_{k} - \mathcal{H}_{k} \left(\bar{\mathbf{x}}_{k}^{\mathrm{f}} \right), \quad \bar{\mathbf{x}}_{k}^{\mathrm{f}} \stackrel{\Delta}{=} \frac{1}{N_{\mathrm{e}}} \sum_{i=1}^{N_{\mathrm{e}}} \mathbf{x}_{k}^{\mathrm{f},i}.$$
 (2b)

 \longrightarrow Notice our trick: $a_{\theta}\left(\mathbf{E}_{k}^{\mathrm{f}}, \delta_{k}\right) \longrightarrow a_{\theta}\left(\mathbf{E}_{k}^{\mathrm{f}}, \mathbf{H}_{k}^{\mathsf{T}} \mathbf{R}_{k}^{-1} \delta_{k}\right)$, i.e., uplift of observation information in state space.

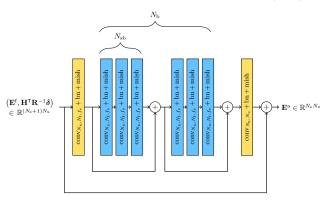
▶ The DA forecast step propagates the analysis ensemble, member-wise:

$$\mathbf{E}_{k+1}^{\mathrm{f}} = \mathcal{M}\left(\mathbf{E}_{k}^{\mathrm{a}}\right). \tag{3}$$

▶ The a_{θ} -based sequential DA will be called DAN in the following.

Neural network architecture

 \blacktriangleright We choose a_{θ} to have a simple residual convolutional neural network (CNN) architecture.



Architecture of the residual convolutional network, where $N_{\rm b}=2$, $N_{\rm sb}=3$. ${\rm conv}_{N_1,N_2,f}$ is a generic one-dimensional convolutional layer of dimension N_1 , with N_2 filters of kernel size f.

Training scheme -1/2

- ▶ Full end-to-end, models and DA [Allen et al. 2025; Alexe et al. 2024; Lean et al. 2025]: not our objective here!
- ▶ Literature (focused on sequential data assimilation):
 - ▶ Learning the analysis of sequential DA is not new [Härter et al. 2012; Cintra et al. 2018], though barely explored.
 - ▶ Learning key components of the analysis in the (En)KF [H. Hoang et al. 1994; S. Hoang et al. 1998] possibly leveraging auto-differentiable structure [Haarnoja et al. 2016; Chen et al. 2022; Luk et al. 2024] was also investigated.
 - Only two key papers so far focused on a non-parametrised analysis using backpropagation through the DA cycles: [McCabe et al. 2021; Boudier et al. 2023].
- ▶ Our training loss (supervised learning):1

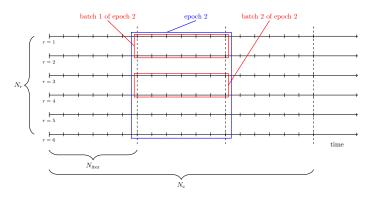
We consider $N_{\rm r}$, $N_{\rm c}$ cycle-long ensemble DA runs, based on $N_{\rm r}$ independent concurrent trajectories of the dynamics $\mathbf{x}_k^{{\rm t},r}$ and as many sequences of observation vectors \mathbf{y}_k^r .

The analysis ensemble is $\mathbf{x}_k^{\mathrm{a},i,r} \in \mathbb{R}^{N_{\mathrm{x}} \times N_{\mathrm{e}}}$. The loss function is defined by

$$\mathcal{L}(\boldsymbol{\theta}) = \sum_{r=1}^{N_{\rm r}} \sum_{k=1}^{N_{\rm c}} \left\| \mathbf{x}_k^{\mathrm{t},r} - \bar{\mathbf{x}}_k^{\mathrm{a},r}(\boldsymbol{\theta}) \right\|^2, \quad \bar{\mathbf{x}}_k^{\mathrm{a},r} \stackrel{\Delta}{=} \frac{1}{N_{\rm e}} \sum_{i=1}^{N_{\rm e}} \mathbf{x}_k^{\mathrm{a},i,r}. \tag{4}$$

¹[Bocquet et al. 2024]

Training scheme -2/2

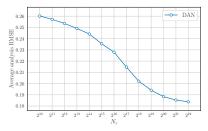


Structure of the dataset organised as a function of time, trajectory sample, batches and epochs.

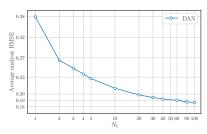
- ▶ Like [McCabe et al. 2021; Boudier et al. 2023], We use truncated backpropagation through time TBPTT [Tang et al. 2018; Aicher et al. 2020], with a truncation at $N_{
 m iter} \ll N_{
 m c}$.
- ► For numerical efficiency, we choose to generate the samples *online*, as the training progresses, i.e. an *infinite training dataset!*

Hyperparameter sensitivity analysis

▶ Sensitivity analysis on key hyperparameters such as the number of trajectories $N_{\rm r}$ in the dataset, and the architecture parameters ($N_{\rm f}$, $N_{\rm b}$, $N_{\rm sb}$) using the standard Lorenz 96 DA configuration ($\mathcal{H}=\mathbf{I}_{\rm x}$, $\mathbf{R}=\mathbf{I}_{\rm x}$).



Filtering performance vs the number of trajectories $N_{
m r}$

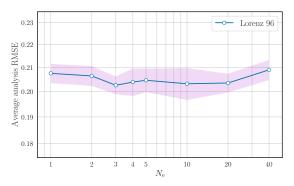


Filtering performance vs the number of filters $N_{
m f}$

- ► The learned DA scheme *yields EnKF-like accuracy!*
- ▶ Compromise between a_{θ} 's size and its accuracy: $N_{\rm r}=2^{18}$, $N_{\rm f}=40$, $N_{\rm b}=5$, $N_{\rm sb}=5$.

Sensitivity to the ensemble size

▶ First key observation: The performance of a_{θ} barely depends on the ensemble size $N_{\rm e}$. Hence localisation is irrelevant and unnecessary.



- Second key observation: a_{θ} does not require inflation and is incredibly robust to noise (as we shall see it applies its own inflation).
- Explanation from the optimisation standpoint: feature collapse of a_{θ} with respect to $N_{\rm e}$ in the training. Potential better solution when $N_{\rm e}>1$, but a_{θ} with $N_{\rm e}=1$ is as accurate as the EnKF!

Consequences and further checks

lackbox Hence, from now on, we will focus on the mode: $N_{
m e}=1$.

▶ Recall

$$\mathbf{E}_{k}^{\mathrm{a}} = \mathbf{E}_{k}^{\mathrm{f}} + a_{\theta} \left(\mathbf{E}_{k}^{\mathrm{f}}, \mathbf{H}_{k}^{\mathsf{T}} \mathbf{R}_{k}^{-1} \delta_{k} \right). \tag{5}$$

 \blacktriangleright Performance of a_{θ} compared to baselines such as optimally tuned 3D-Var, the learned optimal linear filter, optimally tuned EnKF:

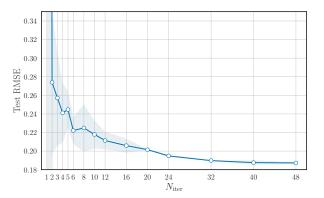
DA method	well-tuned classical	DL-based	aRMSE
EnKF-N, $N_{\rm e}=20$	yes		0.191
EnKF-N, $N_{\rm e}=40$	yes		0.179
3D-Var	yes		0.40
a_{θ} , $N_{\rm e} = 1$, $N_{\rm f} = 40$		yes	0.191
a_{θ} , $N_{\rm e} = 1$, $N_{\rm f} = 100$		yes	0.185
linear $a_{m{ heta}}$, $N_{ m e}=1$, $N_{ m f}=40$		yes	0.384
simplified \hat{a}_{θ} , $N_{\rm e}=1$, $N_{\rm f}=40$		yes	0.382

where (simplified Ansatz)

$$\mathbf{E}_{k}^{\mathbf{a}} = \mathbf{E}_{k}^{\mathbf{f}} + \hat{a}_{\boldsymbol{\theta}} \left(\mathbf{H}_{k}^{\mathsf{T}} \mathbf{R}_{k}^{-1} \boldsymbol{\delta}_{k} \right). \tag{6}$$

Sensitivity to the training (truncation) depth $N_{ m iter}$

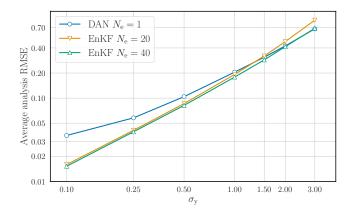
- lacktriangle Training through $N_{
 m iter}=1$ cycle cannot learn about the direct impact of the dynamics on DA.
- ▶ Training through $N_{\rm iter}$ chained cycles is expected to be crucial to the accuracy and robustness of the learned a_{θ} [Bocquet et al. 2025].



▶ Training depth does matter! As expected, $N_{\rm iter} \geq 2$ cycles are required to see significant benefits.

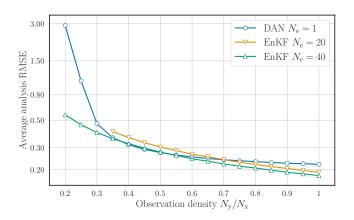
Sensitivity to observation error magnitude

- ▶ Next, we carry out a series of experiments that are not central to our message here but further ground the *viability of such learned* a_{θ} (assuming here $\mathbf{R}_k \stackrel{\Delta}{=} \mathbf{I}_x$).
- ▶ Impact of the *observation noise magnitude* on the data assimilation tests:



Sensitivity to observation sparsity

▶ Impact of the *sparsity of the observation dataset* on the data assimilation tests:



 $ightharpoonup a_{m{ heta}}$ trained with time-dependent, random, observation numbers and positions.

Operator expansion of the analysis

- ▶ We look for a classical Kalman update that would be a good match to a_{θ} seen as a mathematical map, at least for small analysis increments.
- ▶ To that end, we define the time-dependent normalised scalar anomalies

$$b_k = \frac{1}{\sqrt{N_{\mathbf{x}}}} \|a_{\boldsymbol{\theta}}(\mathbf{x}_k, \mathbf{0})\|, \tag{7}$$

along with the associated mean bias b and the standard deviation s of b_k in time.

- \longrightarrow We obtain $b \simeq 5 \times 10^{-3}$ and $s \simeq 10^{-3}$, which are indeed very small compared to the typical aRMSE of an either DAN or EnKF run, i.e., 0.20.
- \blacktriangleright Next, expanding with respect to the innovation, the following functional form for a_{θ} is assumed:

$$a_{\theta}(\mathbf{x}, \mathbf{H}^{\mathsf{T}} \mathbf{R}^{-1} \boldsymbol{\delta}) \approx \mathbf{K}(\mathbf{x}) \cdot \boldsymbol{\delta},$$
 (8)

owing to the fact that no state update is needed when the innovation vanishes, and only keeping the leading order term in δ .

Identifying the operators in the expansion

▶ Innovations $\{\delta_j\}_{j=1,...,N_{\rm p}}$ are sampled from $\delta_j \sim N(\mathbf{0},\mathbf{R})$.

This yields a set of corresponding incremental updates $\left\{\mathbf{a}_j = a_{\theta}(\mathbf{x}, \mathbf{H}^{\mathsf{T}} \mathbf{R}^{-1} \delta_j)\right\}_{j=1,\dots,N_{\mathrm{p}}}$. $\mathbf{K}(\mathbf{x})$ is then estimated with the least squares problem

$$\mathcal{L}_{\mathbf{x}}(\mathbf{K}) = \sum_{j=1}^{N_{p}} \left\| \mathbf{a}_{j} - \bar{\mathbf{a}} - \mathbf{K}(\mathbf{x}) \cdot \left(\boldsymbol{\delta}_{j} - \bar{\boldsymbol{\delta}} \right) \right\|^{2}, \tag{9}$$

where $ar{\mathbf{a}} = N_{\mathrm{p}}^{-1} \sum_{j=1}^{N_{\mathrm{p}}} \mathbf{a}_j$ and $ar{\pmb{\delta}} = N_{\mathrm{p}}^{-1} \sum_{j=1}^{N_{\mathrm{p}}} \pmb{\delta}_j$.

▶ Within the best linear unbiased estimator framework, K is related to \mathbf{P}^{a} through $K = \mathbf{P}^{\mathrm{a}}\mathbf{H}^{\mathsf{T}}\mathbf{R}^{-1}$ so that from Eq. (8),

$$a_{\theta}(\mathbf{x}, \mathbf{H}^{\mathsf{T}} \mathbf{R}^{-1} \boldsymbol{\delta}) \approx \mathbf{P}^{\mathbf{a}} \mathbf{H}^{\mathsf{T}} \mathbf{R}^{-1} \boldsymbol{\delta},$$
 (10)

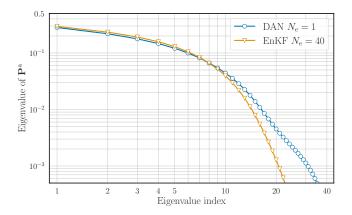
which suggests that an expansion in the second variable $\pmb{\zeta} \in \mathbb{R}^{N_{ ext{x}}}$ of $a_{\pmb{ heta}}$ yields

$$a_{\theta}(\mathbf{x}, \zeta) \approx \mathbf{P}^{\mathbf{a}}(\mathbf{x}) \cdot \zeta.$$
 (11)

Hence, we can obtain a numerical estimation of an equivalent $\mathbf{P}^{\mathrm{a}}(\mathbf{x})$.

What is learned? Supporting numerical results

▶ The surrogate ${\bf P}^a$, denoted ${\bf P}^a_{\rm DAN}$ and estimated from Eq. (11), is compared to that of a concurrent well-tuned EnKF with $N_{\rm e}=40$, whose analysis error covariance matrix is ${\bf P}^a_{\rm EnKF}$.



▶ Time-averaged eigenspectra of $\mathbf{P}_{\mathrm{DAN}}^{\mathrm{a}}$ and $\mathbf{P}_{\mathrm{EnKF}}^{\mathrm{a}}$. They are remarkably close to each other for the first 10 modes. Beyond these modes the a_{θ} operator is likely to selectively apply some multiplicative inflation, as one would expect from such stable DA runs.

Main interpretation

- ▶ Conclusion 1: a_{θ} depends on the innovation but also directly on $\mathbf{x}_k^{\mathrm{f}}$ when $N_{\mathrm{e}}=1$, as opposed to the incremental update of the EnKF: a_{θ} extracts important information from $\mathbf{x}_k^{\mathrm{f}}$.
- ▶ Conclusion 2: a_{θ} manages to assess a $\mathbf{P}_{\mathrm{DAN}}^{\mathrm{a}}$ with $N_{\mathrm{e}}=1$ which is very close to $\mathbf{P}_{\mathrm{EnKF}}^{\mathrm{a}}$ with $N_{\mathrm{e}}=40$, for the dominant axes, and applies multiplicative inflation on the less unstable modes.² We conclude that a_{θ} directly learns about the dynamics features. Hence, for a_{θ} , critical pieces of information on $\mathbf{P}_{k}^{\mathrm{a}}$ are encoded, and thus exploitable, in $\mathbf{x}_{k}^{\mathrm{f}}$ alone.
- → Supported by results from Sacco et al. 2024; Sakov 2025.
- Explanation, conclusion 3: Furthermore, if the DA run (the forecast and analysis cycle) is considered as an ergodic dynamical system of its own, 3 the *multiplicative ergodic theorem* guarantees the existence of a mapping between $\mathbf{x}_k^{\mathbf{f}}$ and $\mathbf{P}_k^{\mathbf{a}}$ that a_{θ} is able to guess. We believe that a generalised variant of the multiplicative ergodic theorem for non-autonomous random dynamics should be applicable. 4

²[Bocquet et al. 2015]

³[Carrassi et al. 2008]

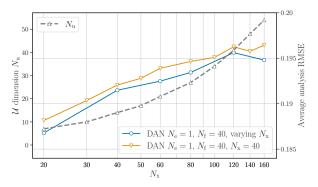
⁴[Arnold 1998; Flandoli et al. 2021]

Outline

- Learning data assimilation from artificial intelligence: first results
 - Sequential data assimilation for chaotic dynamics
 - Learning data assimilation
 - Sensitivity analysis and first results
 - Investigation and interpretation
- Learning data assimilation from artificial intelligence: advanced results
 - Scalability and locality
 - Learning the analysis operator from observations only
 - Learning nonlinearity
- Conclusions

Locality and scalability – 1/2

▶ a_{θ} is now trained without changing the architecture and the hyperparameters ($N_{\rm f}=40$), but with a changing state space dimension $N_{\rm x}\in[20,160]$. Almost as good as well tuned EnKFs with changing dimension $N_{\rm x}$ and $N_{\rm e}=N_{\rm x}$!



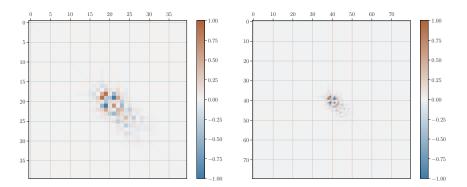
- \longrightarrow We conjecture that a_{θ} extracts *local* pieces of information from $\mathbf{x}_{k}^{\mathrm{f}}$.
- ▶ a_{θ} , learned from Lorenz 96 with $N_{\rm x}=40$ is now tested on Lorenz 96 models with $N_{\rm x}$ ranging from 20 to 160 (same weights and biases!). The performance is still on par with retraining! We called this a *transdimensional transfer*.

Locality and scalability – 2/2

▶ These local patterns (for a_{θ} , not \mathcal{M}) can be pictured from the sensitivity:

$$\mathbf{S} = \left\langle \mathbf{C} : \left[\nabla_{\mathbf{x}} \nabla_{\zeta} a_{\theta}(\mathbf{x}, \zeta)_{|\zeta=0} \right] \right\rangle_{\mathbf{x} \in \mathcal{T}} = \left\langle \mathbf{C} : \left[\nabla_{\mathbf{x}} \mathbf{P}^{\mathbf{a}}(\mathbf{x}) \right] \right\rangle_{\mathbf{x} \in \mathcal{T}}, \tag{12}$$

where $\mathcal T$ is a long L96 trajectory, and $\mathbf C$ is a tensor that leverages translational invariance of the L96 model: $[\mathbf C]_{ij}^{nmk} = \frac{1}{N_v} \delta_{n,i+k} \delta_{m,j+k}$.



Semi-supervised learning

- lackbox What if we do not have access to the truth $\mathbf{x}_k^{\mathrm{t}}$ but to the observations only \mathbf{y}_k ?
- Assume (i) \mathcal{H}_k is linear, (ii) $\mathbf{y}_k \perp \mathbf{y}_{k+1}$, and the estimator $\mathbf{z}_{k+1}^{\boldsymbol{\theta}}$ only depends on $(\mathbf{x}_k, \mathbf{y}_k)$.
- ▶ We define the semi-supervised loss function as

$$\mathcal{L}(\theta) = \sum_{k=1}^{N_c} \left\| \mathbf{y}_k - \mathbf{H}_k \mathbf{z}_k^{\theta} \right\|^2 = \sum_{k=1}^{N_c} \mathcal{L}_k(\theta).$$
 (13)

But we have from the above assumptions:

$$\mathbb{E}_{\mathbf{y}}\left[\mathcal{L}_{k}(\boldsymbol{\theta})\right] = \mathbb{E}_{\mathbf{y}}\left[\left\|\mathbf{y}_{k} - \mathbf{H}_{k}\mathbf{x}_{k}^{t}\right\|^{2}\right] + \mathbb{E}_{\mathbf{y}}\left[\left\|\mathbf{H}_{k}\left(\mathbf{x}_{k}^{t} - \mathbf{z}_{k}^{\boldsymbol{\theta}}\right)\right\|^{2}\right]$$
(14a)

$$= \operatorname{Cst} + \mathbb{E}_{\mathbf{y}} \left[\left\| \mathbf{H}_{k} \left(\mathbf{x}_{k}^{t} - \mathbf{z}_{k}^{\theta} \right) \right\|^{2} \right]$$
 (14b)

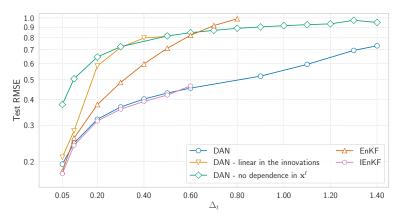
▶ Hence, generalising [McCabe et al. 2021] to non-trivial \mathbf{H}_k , we can learn $\mathbf{z}_k^{\boldsymbol{\theta}}$ from the observation only, with further assumptions on $\left\{\mathbf{H}_k\right\}_{k=1,\ldots,K}$. For instance, we can choose $\mathbf{z}_k^{\boldsymbol{\theta}}$ such that:⁵

$$\mathcal{L}_{k}(\boldsymbol{\theta}) = \left\| \mathbf{y}_{k+1} - \mathbf{H}_{k+1} \mathcal{M} \left\{ \mathbf{x}_{k}^{f} + a_{\boldsymbol{\theta}} \left(\mathbf{x}_{k}^{f}, \mathbf{H}_{k}^{\mathsf{T}} \mathbf{R}_{k}^{-1} \left(\mathbf{y}_{k} - \mathbf{H}_{k} \mathbf{x}_{k}^{f} \right) \right) \right\} \right\|^{2}.$$
 (15)

⁵[Bocquet et al. 2025]

Data assimilation networks in stronger nonlinear conditions

- ▶ Testing DANs as the update time-step Δ_t is increased, with the L96 model.
- ► Comparison with well tuned *EnKF*, *IEnKF*, and well-tuned static background DA methods.



- ▶ Performing at least as well as the IEnKF, without an ensemble, without any nonlinear iterative solver, without inflation and localisation!
- ▶ In addition to the MET map, DANs also implicitly *learn non-Gaussian priors*.

Outline

- Learning data assimilation from artificial intelligence: first result
 - Sequential data assimilation for chaotic dynamics
 - Learning data assimilation
 - Sensitivity analysis and first results
 - Investigation and interpretation
- Learning data assimilation from artificial intelligence: advanced results
 - Scalability and locality
 - Learning the analysis operator from observations only
 - Learning nonlinearity
- Conclusions

Conclusions

- ▶ One can learn robust DA methods, without an ensemble, without inflation, that are as accurate as the very best baseline methods (EnKF), including in stronger nonlinear regimes (IEnKF).
- ▶ We have carried similar numerical experiments with the *Kuramoto–Sivashinski* model and a *single-layer QG model on the sphere*, with similar conclusions.
- ▶ The performance is achieved through (i) implicitly learning the map $\mathbf{x}^f \mapsto \mathbf{P}^f(\mathbf{x}^f)$, and (ii) through implicitly learning non-Gaussian priors!
- ▶ This suggests that end-to-end approaches such as GraphDOP that only have a snapshot of the physical system, can still implicitly rely on dynamical and non-Gaussian priors!
- ▶ Will such *multiplicative ergodic theorem* still be valid in more anisotropic, non-autonomous, forced, multivariate, heterogeneously observed systems?
- ▶ In any case, this promotes a rethinking of the popular sequential DA schemes for chaotic dynamics.

Talk mainly based on:

- M. Bocquet, A. Farchi, T. S. Finn, C. Durand, S. Cheng, Y. Chen, I. Pasmans, and A. Carrassi. "Accurate deep learning-based filtering for chaotic dynamics by identifying instabilities without an ensemble". In: Chaos 29 (2024), p. 091104.
- M. Bocquet, T. S. Finn, S. Cheng, W. Yu, and A. Farchi. "On the performance of data assimilation neural networks in nonlinear conditions". In: (2025). In preparation.

References I

(2013), pp. 803-818.

- C. Aicher, N. J. Foti, and E. B. Fox. "Adaptively Truncating Backpropagation Through Time to Control Gradient Bias". In: Proceedings of The 35th Uncertainty in Artificial Intelligence Conference. Ed. by Ryan P. Adams and Vibhav Gogate. Vol. 115. Proceedings of Machine Learning Research. PMLR, 22–25 Jul 2020, pp. 799–808.
- [2] M. Alexe et al. GraphDOP: Towards skilful data-driven medium-range weather forecasts learnt and initialised directly from observations. 2024. arXiv: 2412.15687 [physics.ao-ph].
- [3] A. Allen et al. "End-to-end data-driven weather prediction". In: Nature (2025).
- [4] L. Arnold. Random Dynamical Systems. Springer Berlin, Heidelberg, 1998, p. 586.
- [5] M. Asch, M. Bocquet, and M. Nodet. Data Assimilation: Methods, Algorithms, and Applications. Fundamentals of Algorithms. SIAM, Philadelphia, 2016, p. 324.
- [6] M. Bocquet, A. Farchi, T. S. Finn, C. Durand, S. Cheng, Y. Chen, I. Pasmans, and A. Carrassi. "Accurate deep learning-based filtering for chaotic dynamics by identifying instabilities without an ensemble". In: Chaos 29 (2024), p. 091104.
- [7] M. Bocquet, T. S. Finn, S. Cheng, W. Yu, and A. Farchi. "On the performance of data assimilation neural networks in nonlinear conditions". In: (2025). in preparation.
- [8] M. Bocquet, P. N. Raanes, and A. Hannart. "Expanding the validity of the ensemble Kalman filter without the intrinsic need for inflation". In: Nonlin. Processes Geophys. 22 (2015), pp. 645–662.
- Nonlin. Processes Geophys. 22 (2015), pp. 645–662.

 M. Bocquet and P. Sakov. "Joint state and parameter estimation with an iterative ensemble Kalman smoother". In: Nonlin. Processes Geophys. 20
- [10] P. Boudier, A. Fillion, S. Gratton, S. Gürol, and S. Zhang, "Data Assimilation Networks". In: J. Adv. Model. Earth Syst. 15 (2023), e2022MS003353.
- P. Bouder, A. Fillott, S. Gracton, S. Guror, and S. Zhang. Data Assimilation Networks 1 II. S. Adv. Woder. Latti Syst. 15 (2023), e2022Wi300333
- [11] A. Carrassi, M. Ghil, A. Trevisan, and F. Uboldi. "Data assimilation as a nonlinear dynamical systems problem: Stability and convergence of the prediction-assimilation system". In: Chaos 18 (2008), p. 023112.
- [12] Y. Chen, D. Sanz-Alonso, and R. Willett. "Autodifferentiable Ensemble Kalman Filters". In: SIAM J. Math. Data Sci. 4 (2022), pp. 801–833.
- [13] R. S. Cintra and H. F. de Campos Velho. "Data assimilation by artificial neural networks for an atmospheric general circulation model". In: Advanced applications for artificial neural networks. Ed. by A. ElShahat. IntechOpen, 2018. Chap. 17, pp. 265–286.
- [14] F. Flandoli and E. Tonello. An introduction to random dynamical systems for climate. 2021.
- [15] T. Haarnoja, A. Ajay, S. Levine, and P. Abbeel. "Backprop KF: Learning Discriminative Deterministic State Estimators". In: Advances in Neural Information Processing Systems. Ed. by D. Lee, M. Sugiyama, U. Luxburg, I. Guyon, and R. Garnett. Vol. 29. Curran Associates, Inc., 2016.

References II

- [16] T. P. Härter and H. F. de Campos Velho. "Data Assimilation Procedure by Recurrent Neural Network". In: Eng. Appl. Comput. Fluid Mech. 6 (2012), pp. 224–233.
- [17] H.S. Hoang, P. De Mey, and O. Talagrand. "A simple adaptive algorithm of stochastic approximation type for system parameter and state estimation". In: Proceedings of 1994 33rd IEEE Conference on Decision and Control. Vol. 1. 1994, 747–752 vol.1.
- [18] S. Hoang, R. Baraille, O. Talagrand, X. Carton, and P. De Mey. "Adaptive filtering: application to satellite data assimilation in oceanography". In: Dynam. Atmos. Ocean 27 (1998), pp. 257–281.
- [19] P. Lean, M. Alexe, E. Boucher, E. Pinnington, S. Lang, P. Laloyaux, N. Bormann, and A. McNally. Learning from nature: insights into GraphDOP's representations of the Earth System. 2025. arXiv: 2508.18018 [physics.ao-ph].
- [20] E. Luk, E. Bach, R. Baptista, and A. Stuart, Learning Optimal Filters Using Variational Inference, 2024, arXiv: 2406.18066 [cs.LG].
- M. McCabe and J. Brown. "Learning to Assimilate in Chaotic Dynamical Systems". In: Advances in Neural Information Processing Systems. Ed. by M. Ranzato, A. Beygelzimer, Y. Dauphin, P.S. Liang, and J. Wortman Vaughan. Vol. 34. Curran Associates, Inc., 2021, pp. 12237–12250.
- [22] M. A. Sacco, M. Pulido, J. J. Ruiz, and P. Tandeo. "On-line machine-learning forecast uncertainty estimation for sequential data assimilation". In: Q. J. R. Meteorol. Soc. 150 (2024), pp. 2937–2954.
- [23] P. Sakov. On building the state error covariance from a state estimate, 2025. arXiv: 2411.14809 [nlin.CD].
- [24] H. Tang and J. Glass. "On Training Recurrent Networks with Truncated Backpropagation Through time in Speech Recognition". In: 2018 IEEE Spoken Language Technology Workshop (SLT). 2018, pp. 48–55.