



# Elucidating the performance of data assimilation neural networks for chaotic dynamics

Marc Bocquet<sup>1</sup>, Tobias Sebastian Finn<sup>1</sup>, Sibongwe Cheng<sup>1</sup>, and Alban Farchi<sup>2</sup>

<sup>1</sup>CEREA, ENPC, EDF R&D, Institut Polytechnique de Paris, Île-de-France, France

<sup>2</sup>ECMWF, Shinfield Park, Reading, UK

**Correspondence:** Marc Bocquet (marc.bocquet@enpc.fr)

## Abstract.

Recent work has shown that the analysis operator in sequential data assimilation designed to track chaotic dynamics, can be learned with deep learning from the sole knowledge of a true state trajectory and observations thereof. This approach to learning the analysis is computationally more challenging, yet conceptually more fundamental than approaches that learn a direct mapping from forecasts and observations to the corresponding analysis increments. Such learned scheme has been demonstrated to achieve accuracy comparable to that of the ensemble Kalman filter when applied to low-order dynamics. Strikingly, the same accuracy can be reached with a single state forecast instead of an ensemble, hence bypassing the need to explicitly represent forecast uncertainty.

In this study, we extend the investigation of such learned analysis operators beyond the preliminary experiments reported so far. First, we analyse the emergence of local patterns encoded in the operator, which accounts for the remarkable scalability of the approach to high-dimensional state spaces. Second, we assess the performance of the learned operators in stronger nonlinear regimes of the chaotic dynamics. We show that they can match the efficiency of the iterative ensemble Kalman filter, the baseline in this context, while avoiding the need for nonlinear iterative optimisation. Throughout the paper, we seek underlying reasons for the efficiency of the approach, drawing on insights from both machine learning and nonlinear data assimilation.

## 1 Introduction

Accurate prediction of geophysical flows relies on the continual correction of model trajectories using observations. This sequential data assimilation process is essential in high-dimensional chaotic systems, where errors amplify rapidly and model imperfections accumulate over time (Kalnay, 2003; Asch et al., 2016; Carrassi et al., 2018). In operational meteorology and oceanography, ensemble-based Kalman filters and ensemble variational methods provide reliable and well-understood frame-



works for these updates, but they remain computationally demanding and rely on explicit representations of flow-dependent forecast error covariances.

Recent developments have suggested that part of this complexity can be replaced by learned operators. It has been shown that the analysis step of a data assimilation cycle can be learned using deep neural networks, either by emulating existing schemes (e.g., Härter and de Campos Velho, 2012; Cintra and de Campos Velho, 2018; Maddy et al., 2024), or more fundamentally by learning an end-to-end update rules (McCabe and Brown, 2021). Within this latter family of approaches, termed data assimilation networks (DANs, Boudier et al., 2023), a surprising result has emerged: as demonstrated by Bocquet et al. (2024), a learned analysis operator can match the accuracy of a well-tuned ensemble Kalman filter (EnKF) even when it uses only single forecast trajectories, without any ensemble. This result challenges the long-standing assumption that explicit ensemble representations are indispensable to estimate flow-dependent uncertainties in chaotic systems.

Explaining this phenomenon is a central question for DA methodology. Initial investigations suggest that the learned operator implicitly reconstructs aspects of the analysis error covariances normally diagnosed from an ensemble, effectively uncovering key uncertainty directions from the forecast state alone (Bocquet et al., 2024). This behaviour is consistent with viewing the full DA cycle as a random dynamical system, for which generalised forms of the multiplicative ergodic theorem (Oseledec, 1968) offer a state-dependent structure linking model trajectories to dominant error-growth directions. At the same time, numerical experiments show that the neural network is likely to identify local patterns rather than memorising global states, which explains why its performance scales to larger systems and remains robust across different model dimensions.

This paper deepens the investigation of these mechanisms. In Sect. 2, we dive into the methodology of DANs, and recall the main results and questions raised in Bocquet et al. (2024). In Sect. 3, building on a more thorough analysis of the performance dependence on the batch size, dataset length, and the number of assimilation cycles used during backpropagation, we study how the operator behaves when interpreted as a diagnostic of uncertainty, and we propose methods to expose and interpret the local structures it extracts. In Sect. 4, we test the limits of the approach in more strongly nonlinear regimes, where the iterative ensemble Kalman filter (Sakov et al., 2012) often serves as the most accurate baseline. We show that the learned operator can reach comparable performance without requiring an ensemble or a nonlinear optimisation, and we offer a data assimilation-based interpretation for this ability. Section 5 presents our conclusions. Supporting numerical and mathematical results are collated in the appendices of this paper.

Throughout the paper, we will illustrate our results with the Lorenz-96 model (L96, Lorenz and Emanuel, 1998). More broadly, our goal is not only to document the performance of the learned analysis operators but also to clarify the mechanisms that underlie them, thereby contributing to the growing theoretical understanding of how deep learning and sequential data assimilation interact in chaotic geophysical systems (Cheng et al., 2023).

## 2 Theory and methods

In this section, we provide a deeper description of the problem, its context, and its mathematical formulation.



## 2.1 Sequential data assimilation for chaotic dynamics

Mathematically, data assimilation (DA), and in particular filtering algorithms, are meant to accurately estimate the state vector  $\mathbf{x}_k^t \in \mathbf{E}_x \triangleq \mathbb{R}^{N_x}$  of a physical system, where “t” stands for truth, at times  $\tau_k$  for  $k = 1, \dots, K$  along a trajectory of the dynamical system. These states are evolved according to

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

where  $\mathcal{M}$  is the integrated model over  $\tau_{k+1} - \tau_k = \Delta\tau$ . The dynamical system  $\mathcal{M}$  is assumed to be chaotic, such as for most geofluids, which is a prime incentive for frequently updating our knowledge of the system. It is furthermore assumed ergodic and autonomous, i.e. does not explicitly depend on time. Furthermore, the physical system is observed through

$$\mathbf{y}_k = \mathcal{H}_k(\mathbf{x}_k^t) + \varepsilon_k, \quad \varepsilon_k \sim \mathcal{N}(\mathbf{0}, \mathbf{R}_k), \quad (1b)$$

where  $\mathbf{y}_k \in \mathbf{E}_y^k \triangleq \mathbb{R}^{N_{y,k}}$  is the observation vector at  $\tau_k$  obtained from the hidden state  $\mathbf{x}_k^t$  via an observation operator  $\mathcal{H}_k$ , and perturbed by a white-in-time Gaussian noise  $\varepsilon_k$  of mean  $\mathbf{0}$  and covariance matrix  $\mathbf{R}_k$ . This very common but simplified formulation of the filtering problem with additive Gaussian noise is sufficient for the goals of this paper.

A *sequential* filtering DA scheme estimates the state at  $\tau_k$  from observation  $\mathbf{y}_k$  and from background information about the state at  $\tau_{k-1}$ . The scheme is recursive and unfolds as time flows and observations are collected. Its estimator,  $\mathbf{x}_k^a$  at  $\tau_k$ , called the *analysis*, is meant as an estimator of the conditional probability density function  $p(\mathbf{x}_k^t | \mathbf{y}_k, \mathbf{y}_{k-1}, \dots, \mathbf{y}_1)$ .

## 2.2 Learning the data assimilation analysis

### 2.2.1 Learning an incremental analysis

The approach developed in Bocquet et al. (2024), subsequently referred to as Boc24, in the wake of McCabe and Brown (2021); Boudier et al. (2023) is summarised in the following since it is the foundation of the present paper. The analysis step of the DA scheme is assumed to be given by the (incremental) analysis operator  $\mathbf{a}_\theta$ , typically a (deep) neural network, which depends on a set of weights and biases, stacked in the  $\theta$  vector, and which is defined, at time  $\tau_k$ , through

$$\mathbf{x}_k^a = \mathbf{x}_k^f + \mathbf{a}_\theta(\mathbf{x}_k^f, \mathbf{H}_k^T \mathbf{R}_k^{-1} \delta_k), \quad (2a)$$

$$\delta_k \triangleq \mathbf{y}_k - \mathcal{H}_k(\mathbf{x}_k^f), \quad (2b)$$

where  $\delta_k$  is the innovation,  $\mathbf{x}_k^a$  is the analysis state mentioned previously, and  $\mathbf{x}_k^f$  is the forecast state, to be defined shortly.  $\mathbf{H}_k$  is the tangent linear operator of  $\mathcal{H}_k$ . If  $\mathbf{a}_\theta$  was a function of  $\delta_k$  rather than  $\mathbf{H}_k^T \mathbf{R}_k^{-1} \delta_k$ , it could only handle static observation configurations, and would require to be retrained whenever that configuration changes. Although not the focus of this paper, this important issue is circumvented here by using the mapping from observation to state space  $\mathbf{H}_k^T \mathbf{R}_k^{-1}$ , such that both inputs of  $\mathbf{a}_\theta$  are in the same static state space  $\mathbf{E}_x$ . The DA forecast step propagates the analysis state to the subsequent date:

$$\mathbf{x}_{k+1}^f = \mathcal{M}(\mathbf{x}_k^a). \quad (3)$$



The operator  $\mathbf{a}_\theta$  is trained using supervised learning by comparing the analysis states to the true states, through the loss

$$\mathcal{L}(\theta) = \sum_{r=1}^{N_r} \sum_{k=1}^{N_c} \|\mathbf{x}_k^{t,r} - \mathbf{x}_k^{a,r}(\theta)\|^2, \quad (4)$$

85 where  $\mathbf{x}_k^{t,r}$  are the state vectors of the true trajectory, indexed by the time index  $k$  and a trajectory index  $r$  as  $N_r$  of them are processed concurrently in the training of  $\mathbf{a}_\theta$ .

The  $N_c$  parameter counts the number of cycles of each DA run. It can potentially be infinite since trajectories can be generated online as the training progresses. The implementation of similar losses is detailed in McCabe and Brown (2021); Boudier et al. (2023), and Sect. II of Boc24. This notably requires to truncate backpropagation in time (Tang and Glass, 2018)

90 which restricts the dependence on  $\theta$  over  $N_{\text{iter}} \ll N_c$  cycles only so as to reduce the computational cost and the excessive GPU memory (VRAM) requirement. Alternative losses based on probability density functions (pdf) diagnostics are possible and discussed in, e.g., Boudier et al. (2023). A semi-supervised loss where the training dataset reduces to the sparse and noisy observations is also possible by generalising a proposal of McCabe and Brown (2021), but is out of the scope of this paper.

### 2.2.2 Accuracy of the discovered assimilation scheme

95 Note that an ensemble variant of the update Eq. (2) was first considered in McCabe and Brown (2021); Boudier et al. (2023), and later by Boc24. Experimenting with the same low-order chaotic model, they found an accuracy of the learned DA scheme close to that of a well-tuned EnKF, which is per se very promising. In those experiments, DAN is built on an ensemble of analyses and forecasts. However, Boc24 demonstrated that this accuracy remains unchanged when these ensembles are reduced to a single state. This is a very surprising result since it is expected (and verified in the L96 context) that the EnKF-like methods  
100 critically relying on an ensemble representing the *errors of the day*, have a significant edge over other sequential DA methods that leverage a single forecast state such as 3D-Var and 4D-Var (see Bocquet and Sakov, 2013, for a quantitative comparison with the same model). This explains why numerical weather prediction (NWP) centres operating 4D-Var actually rely on an ensemble of such 4D-Var, a technique called EDA (see, e.g., Chapter 7 in Asch et al., 2016; Bannister, 2017, and references within), or use information from a concurrent EnKF (Buehner et al., 2015). That is why achieving the accuracy of a well-tuned  
105 EnKF with a single forecast state should have far-reaching implications on the mechanisms and designs of DA algorithms for chaotic dynamics, and warrants investigating the reason for such success.

### 2.3 Investigating the reason for this success

To unveil the mechanisms leveraged by the learned analysis operator to achieve this accuracy, Boc24 performed a numerical expansion of the operator  $\mathbf{a}_\theta$  in terms of associated classical DA objects, such as the analysis error covariance matrix  $\mathbf{P}^a$ :

$$110 \quad \mathbf{a}_\theta(\mathbf{x}, \zeta) \approx \mathbf{P}^a(\mathbf{x})\zeta, \quad (5)$$

where  $\mathbf{P}^a(\mathbf{x})$  is the analysis error covariance matrix as diagnosed by the learned analysis operator  $\mathbf{a}_\theta$ . This covariance matrix was numerically obtained through a linear regression in between a large ensemble of  $\zeta$  samples (corresponding to the projected



innovations  $\mathbf{H}^T \mathbf{R}^{-1} \delta$ ) and  $\mathbf{a}_\theta(\mathbf{x}, \zeta)$  outputs. It was found that this error covariance matrix is remarkably close to that of a well-tuned EnKF, especially for its dominant eigenvectors (which carry most of the uncertainty in the analysis).

115 Hence, the performance of DAN must be to a large extent due to its ability to estimate the analysis error covariance matrix from a single forecast state, hence the sole dependence of  $\mathbf{P}^a$  on  $\mathbf{x}^f$ . Such estimation is pivotal in maintaining the accuracy of the sequential DA over time. Other alternatives to ensemble forecasting, such as deriving the dynamics of the statistical moments of the errors (Pannekoucke et al., 2016, 2018), or estimating these dynamics through machine learning (Pannekoucke and Fablet, 2020; Sacco et al., 2024; Lu, 2025) are either numerically very costly or in their infancy. Hence, theoretical support  
 120 was needed to ascertain that such a feat is not chimeric.

Boc24 conjectured that this acquired ability can be fundamentally explained by the existence of an ad-hoc *multiplicative ergodic theorem* (MET). From the seminal MET result by Oseledec (1968), we know that, for an autonomous ergodic dynamical system such as  $\mathcal{M}$ , there exists a mapping from each of the system's states to the corresponding Lyapunov covariant vectors. Generalising, one can consider the whole sequential DA process as a dynamical system on its own (Carrassi et al., 2008). Such  
 125 DA process is not autonomous because it indirectly depends on the truth trajectory and the time-dependent observation operators. Moreover, it is a random process, since stochasticity is injected via the noisy observations. It turns out that generalised variants of the MET for non-autonomous random dynamics are possible (Arnold, 1998; Chekroun et al., 2011; Flandoli and Tonello, 2021; Ghil and Sciamarella, 2023) and are potentially applicable to such sequential DA process. Hence, Boc24 stated that  $\mathbf{a}_\theta$  must learn such mapping from the forecast state to the analysis error covariance matrix as seen by  $\mathbf{a}_\theta$ , together with  
 130 how to process this information and combine it with the innovation.

To explain the success of  $\mathbf{a}_\theta$ , one may suggest that the neural network memorises global configurations of the forecast state, with very limited ability to generalise. On the contrary, Boc24 showed via indirect scaling experiments that  $\mathbf{a}_\theta$  learns to identify local patterns (i.e. with a limited range in space), which was made easier by the architecture of  $\mathbf{a}_\theta$  being a residual convolutional neural network. Indeed, when the dimension  $N_x$  of the L96 state space is increased, and new  $\mathbf{a}_\theta$  operators are  
 135 learned but with a fixed number of weights and biases of the backbone architecture, the accuracy remains that of a well tuned EnKF of matching dimension. Yet, in the large  $N_x$  limit,  $\mathbf{a}_\theta$  with the same number of degrees of freedom should not be able to memorise increasingly numerous global patterns. Hence,  $\mathbf{a}_\theta$  must extract local patterns, and a limited number of them. This is further supported by learning  $\mathbf{a}_\theta$  with the original L96 dimension,  $N_x = 40$ , but applying it to L96 with significantly different  $N_x$  in successful DA runs (this is allowed by the convolutional architecture which does not explicitly depend on  $N_x$ ),  
 140 performing on par with a well tuned EnKF. Hence, any local pattern learned in the case  $N_x = 40$ , must still be spotted by  $\mathbf{a}_\theta$  where  $N_x \neq 40$ . This outcome is consistent with the existence of such local patterns, since L96 is an extensive model when  $N_x$  is increased, with the number of nonlinear interacting waves in the model being proportional to  $N_x$ .

### 3 Exploration of data assimilation networks

With the previously established context in mind, we now explore what data assimilation networks ( $\mathbf{a}_\theta$ ) learn in mild nonlinear  
 145 regimes. The neural network for  $\mathbf{a}_\theta$  implements Eq. (2). The architecture for  $\mathbf{a}_\theta$  we choose in the present paper is the same as



the one reported in Boc24, i.e. a simple residual convolutional neural network that accounts well for the spatial homogeneity of L96 (and hence its statistical stationarity). The deep learning architecture is described in Boc24. All the training tasks are carried out over a training dataset with a minimisation of the loss controlled by computing the loss over a validation dataset.

The numerical experiments of the present paper are performed on L96, a chaotic model abundantly used for benchmarking new sequential data assimilation algorithms. As a reminder, L96 represents a mid-latitude zonal circle of the global atmosphere. It is governed by  $N_x = 40$  ordinary differential equations:

$$\frac{dx_n}{dt} = (x_{n+1} - x_{n-2})x_{n-1} - x_n + F, \quad (6)$$

where  $F = 8$ , and with cyclic boundary conditions. The resulting model is chaotic with 13 positive and 1 neutral Lyapunov exponents. Its Lyapunov time, defined as the inverse of the first Lyapunov exponent, is about 0.60, which corresponds to 3 days of a typical weather forecasting model (Lorenz and Emanuel, 1998).

Although the results are generalisable to sparse observations, the observation operator will mostly be chosen to be the identity  $\mathbf{H}_k \triangleq \mathbf{I}_x$ . The observations are read off the true state and perturbed with an unbiased white-in-time Gaussian additive noise of covariance matrix  $\mathbf{R}_k \triangleq \mathbf{I}_x$  following Eq. (1b). This will be our reference DA configuration. Examples of alternative sparseness and noise levels are given in Boc24, but here as well when relevant.

All the  $\mathbf{a}_\theta$  operators learned in this experimental configuration are subsequently evaluated with an analysis RMSE whose root mean square is averaged over the  $N_x$  variables, averaged over time, and assessed over a dataset of independent test trajectories, similarly to traditional DA twin experiments. For brevity, this score computed for each trained DAN scheme will simply be called *test RMSE* of the DAN scheme in the rest of this paper.

In this configuration, running a well-tuned EnKF in a twin experiment and comparing its analysis to the truth yield a test RMSE between 0.18 and 0.20 depending on the ensemble size  $N_e$  and whether localisation is used or not. By contrast, a basic but well-tuned 3D-Var or a reasonably short window basic but well tuned 4D-Var yields a test RMSE of about 0.40. They largely underperform the EnKF because they fail to capture the errors of day (Bocquet and Sakov, 2013; Fillion et al., 2018).

In order to be able to perform a large number of training experiments on a limited number of GPUs and limited VRAM, we carried out a sensitivity study on the batch size, the size of the datasets, and the backpropagation truncation, beyond the restricted set of experiments reported in Boc24. Since the results are technical and mostly of practical interest, they are reported in Appendix A. They were nonetheless instrumental in obtaining the main numerical results of this paper.

### 3.1 Linear expansion in the innovations

In this section, we discuss the relevance of expanding  $\mathbf{a}_\theta$  linearly in the innovations, as in Eq. (5). We recall that the focus is on the analysis step of the DAN process, where  $N_e = 1$ , i.e. a single forecast state is propagated in between updates. With such premises, we can assume that the analysis state  $\mathbf{x}^a$  of a close to optimal DA method is given by the maximum a posteriori of the conditional pdf, hence by the minimum of the cost function associated to the analysis. In the following, the observation operator is assumed linear (or linearised) for simplicity, i.e.  $\mathcal{H} \triangleq \mathbf{H}$ .



### 3.1.1 A Gaussian standpoint on the analysis

Here, we further assume that the background errors are Gaussian. Nonetheless, as opposed to a basic 3D-Var, the background error covariance matrix depends on the forecast state. Hence, the typical analysis cost function associated to the analysis at any given time step has the form:

$$\mathcal{J}(\mathbf{x}) = \frac{1}{2} \|\mathbf{y} - \mathbf{H}\mathbf{x}\|_{\mathbf{R}^{-1}}^2 + \frac{1}{2} \|\mathbf{x} - \mathbf{x}^f\|_{(\mathbf{P}^f(\mathbf{x}^f))^{-1}}^2. \quad (7)$$

As a consequence,  $\mathcal{J}$  is quadratic in  $\mathbf{x}$ , strictly convex, and its minimum argument is (Daley, 1991)

$$\mathbf{x}^a = \mathbf{x}^f + \left( \mathbf{H}^T \mathbf{R}^{-1} \mathbf{H} + (\mathbf{P}^f(\mathbf{x}^f))^{-1} \right)^{-1} \mathbf{H}^T \mathbf{R}^{-1} (\mathbf{y} - \mathbf{H}\mathbf{x}^f) \quad (8a)$$

$$= \mathbf{x}^f + \mathbf{P}^a(\mathbf{x}^f) \mathbf{H}^T \mathbf{R}^{-1} (\mathbf{y} - \mathbf{H}\mathbf{x}^f), \quad (8b)$$

which matches the expansion Eq. (5). Hence, the analysis equation of DANs with a single ensemble member, Eq. (2), and supplemented Gaussian hypotheses, are sufficient assumptions for yielding Eq. (5).

### 3.1.2 Numerical evidence

To numerically test whether such an expansion is a good approximation for  $\mathbf{a}_\theta$ , we created a modified DAN which concurrently and explicitly learns the mapping  $\mathbf{x}^f \mapsto \mathbf{P}^f(\mathbf{x}^f)$  and linearly combines it with the projected innovations, strictly following the update equation Eq. (8b). Details on our scalable implementation of Eq. (8b), for a DAN linear in the projected innovations, can be found in Appendix B. One must keep in mind that such a *deconstruction* of  $\mathbf{a}_\theta$  is numerically inefficient since it requires to build a representation of the covariance matrix  $\mathbf{P}^f$  before being applied to the projected innovations, an operation which is likely to be achieved with our standard DAN without ever computing  $\mathbf{P}^f$  explicitly. Yet, this now directly connects with the heuristic developed by Sacco et al. (2024); Sakov (2025), where the mapping  $\mathbf{x}^f \mapsto \mathbf{P}^f(\mathbf{x}^f)$  is learned and then successfully used within a classical EnKF.

This linear-in- $\zeta$  DAN scheme turns out to be as accurate as a well-tuned EnKF with  $N_e = 40$  for  $\Delta_t = 0.05$ , with a test RMSE of 0.19, which is in line with the results by Sacco et al. (2024); Sakov (2025). Hence, we can claim that the good performance of DAN in this mild nonlinear regime importantly relies on the (implicit) estimation of the mapping  $\mathbf{x}^f \mapsto \mathbf{P}^f(\mathbf{x}^f)$ .

### 3.1.3 On the importance of the $\mathbf{x}^f \mapsto \mathbf{P}^f(\mathbf{x}^f)$ map

Bocquet et al. (2017); Bocquet and Carrassi (2017) showed that the ensemble of the EnKF applied to chaotic dynamics asymptotically collapses onto the unstable subspace of the dynamics. Their results are exact when the model dynamics are linear between two observation times, even in the non-autonomous case. Their theoretical framework provides rigorous mathematical arguments that support, if not fully justify, the heuristic ideas put forward by Sacco et al. (2024); Sakov (2025).

Let us focus on Sakov (2025), whose algorithmic constructions targeted at estimating  $\mathbf{P}^f(\mathbf{x}^f)$  are especially transparent. Their Algorithm A1 proceeds as follows: (i) the state  $\mathbf{x}^f$  is backtracked by  $T$  time steps; (ii) the tangent linear model, evaluated along the resulting state trajectory from  $\tau_{-T}$  to  $\tau_0$ , is applied to a matrix of state perturbations  $\varepsilon \mathbf{I}_x$ ; and (iii) the resulting





perturbations at time  $\tau_0$  are used to estimate  $\mathbf{P}^f(\mathbf{x}^f)$ . This procedure is closely related to the Assimilation in the Unstable Space (AUS) methods (Palatella et al., 2013; Carrassi et al., 2022), since for sufficiently large  $T$  the output of the tangent linear  
 210 model at  $\tau_0$  provides a square-root factor of the backward Lyapunov vectors.

Algorithm A2 of Sakov (2025) also backtracks the state  $\mathbf{x}^f$  by  $T$  time steps but differs in the subsequent step: instead of propagating perturbations, it applies a (square-root) Kalman filter to an initial covariance matrix  $\varepsilon^2 \mathbf{I}_x$  from  $\tau_{-T}$  to  $\tau_0$ , yielding a more refined estimate of  $\mathbf{P}^f(\mathbf{x}^f)$  at  $\tau_0$ . Within the degenerate Kalman filter framework of Bocquet et al. (2017), Algorithm A2 can be directly interpreted as the covariance propagation step of the degenerate Kalman filter itself.

215 In this setting, the covariance evolution is formally decoupled from the state update, with the important caveat that the forecast error covariance  $\mathbf{P}^f$  depends on the state  $\mathbf{x}^f$  at time  $\tau_0$ . Indeed, Bocquet et al. (2017) showed that, asymptotically (i.e. for large  $T$ ),  $\mathbf{P}^f$  depends on the system dynamics only through the Lyapunov vectors. By the multiplicative ergodic theorem (MET), these vectors depend solely on the state  $\mathbf{x}^f$  at  $\tau_0$ .

220 Taken together, these theoretical results provide a clear rationale for the existence of a map  $\mathbf{x}^f \mapsto \mathbf{P}^f(\mathbf{x}^f)$  and its critical importance in estimating the errors-of-the-day in DA schemes based on a single forecast state  $N_e = 1$ .

### 3.2 Patterns

In this section, we visualise a footprint of the local patterns leveraged by DAN to make its inferences. To that end, we perform a linear sensitivity analysis and study the dependence of  $\nabla_{\zeta} \mathbf{a}_{\theta}(\mathbf{x}, \zeta)|_{\zeta=0} \approx \mathbf{P}^a(\mathbf{x})$  on  $\mathbf{x}$ , i.e. the forecast state.

#### 3.2.1 Mean marginal gain

225 The sensitivity analysis will focus on Jacobians of the analysis operator with respect to  $\mathbf{x}$ :

$$\mathbf{\Gamma}(\mathbf{x}, \zeta) \triangleq \nabla_{\mathbf{x}} \nabla_{\zeta} \mathbf{a}_{\theta}(\mathbf{x}, \zeta), \quad (9a)$$

$$\mathbf{\Gamma}(\mathbf{x}) \triangleq \mathbf{\Gamma}(\mathbf{x}, 0) = \nabla_{\mathbf{x}} \nabla_{\zeta} \mathbf{a}_{\theta}(\mathbf{x}, \zeta)|_{\zeta=0}. \quad (9b)$$

We call  $\mathbf{\Gamma}(\mathbf{x})$ , the *marginal error covariance matrix* since it is a proxy to  $\nabla_{\mathbf{x}} \mathbf{P}^a(\mathbf{x})$  as per Eq. (5). This is a *tensor field*, i.e. a map from the state space  $\mathbf{E}_x$  to  $\mathbf{E}_x^3$ . Component-wise, using a conventional placement of indices, its definition reads

$$230 [\mathbf{\Gamma}(\mathbf{x})]_{ijl} \triangleq \left( \partial_{x_i} \partial_{\zeta_j} a_{\theta, i}(\mathbf{x}, \zeta) \right)_{|\zeta=0}. \quad (10)$$

It can be interpreted as the marginal variation of  $\mathbf{P}^a(\mathbf{x})$  when the forecast state  $\mathbf{x}$  is perturbed. If the DA system is such that  $\mathbf{H}^T \mathbf{R}^{-1} = \mathbf{I}_x$ ,  $\mathbf{P}^a(\mathbf{x})$  coincides with the Kalman gain and  $\mathbf{\Gamma}(\mathbf{x})$  is hence the *marginal Kalman gain*. For the sake of brevity,  $\mathbf{\Gamma}(\mathbf{x})$  will hence be called the *marginal gain* in the rest of the paper, irrespective of the observation configuration.

To mitigate the complexity in studying the map  $\mathbf{x} \mapsto \mathbf{\Gamma}(\mathbf{x})$ , we introduce the *mean marginal gain*,  $\bar{\mathbf{\Gamma}}$ . This 3-tensor is defined  
 235 as the average of  $\mathbf{\Gamma}(\mathbf{x})$  over the  $K$  states of a long trajectory  $\mathcal{T}_x = \{\mathbf{x}_k^t\}_{k=1, \dots, K} \subset \mathbf{E}_x^K$  of the  $\mathcal{M}$ -based ergodic chaotic dynamics:

$$\bar{\mathbf{\Gamma}} \triangleq \langle \mathbf{\Gamma}(\mathbf{x}) \rangle_{\mathbf{x} \in \mathcal{T}_x} \underset{K \rightarrow \infty}{=} \int d\mathbf{x} \pi(\mathbf{x}) \mathbf{\Gamma}(\mathbf{x}) = \mathbb{E}_{\mathbf{x} \sim \pi} [\mathbf{\Gamma}(\mathbf{x})], \quad (11)$$





where  $\pi$  is the invariant distribution of the ergodic chaotic dynamics.

As recalled in Sect. 2.3,  $\nabla_{\zeta} \mathbf{a}_{\theta}(\mathbf{x}, \zeta)|_{\zeta=0} \approx \mathbf{P}^a(\mathbf{x})$  was heuristically estimated in Boc24 using samples of  $\zeta$  followed by  
 240 a linear regression. This estimation through a regression can be formally justified using the following argument. We wish to  
 average  $\nabla_{\zeta} \mathbf{a}_{\theta}(\mathbf{x}, \zeta)$  over the pdf  $\rho(\zeta)$  of the projected innovations  $\zeta$ , which is assumed to be a Gaussian with a positive-definite  
 covariance matrix  $\Sigma_{\rho}$ . Then it can be shown that

$$\mathbb{E}_{\zeta \sim \rho} [\nabla_{\zeta} \mathbf{a}_{\theta}(\mathbf{x}, \zeta)] = \int d\zeta \rho(\zeta) \nabla_{\zeta} \mathbf{a}_{\theta}(\mathbf{x}, \zeta) = \Sigma_{\rho}^{-1} \text{Cov}_{\zeta \sim \rho} [\zeta, \mathbf{a}_{\theta}(\mathbf{x}, \zeta)], \quad (12)$$

where  $\text{Cov}[\cdot, \cdot]$  is the covariance operator, and  $\Sigma_{\rho}^{-1}$  operates on the first tensor factor. A proof is given in Appendix C,  
 245 along with a generalisation to the case where  $\Sigma_{\rho}$  is only semi positive-definite, which is made necessary because the set  
 $\{\zeta_k = \mathbf{H}_k^T \mathbf{R}_k^{-1} \delta_k\}_k$  may only span a subspace of  $\mathbf{E}_{\mathbf{x}}$  if  $\mathbf{H}_k$  is not injective. This result is none other than Stein's lemma  
 (Liu, 1994) applied to the Gaussian approximation of  $\rho$  and  $\mathbf{a}_{\theta}(\mathbf{x}, \cdot)$ . It aligns with the connection established by Lemma 2 of  
 Agarwal et al. (2021) between  $\Sigma_{\rho}^{-1} \text{Cov}_{\zeta \sim \rho} [\zeta, \mathbf{a}_{\theta}(\mathbf{x}, \zeta)]$ , which is a perturbations-based estimator (Ribeiro et al., 2016), and  
 a SmoothGrad estimator (Smilkov et al., 2017). In the limit where  $\zeta$ , as a random vector, is concentrated around  $\mathbf{0}$  and  $\rho$  can  
 250 be approximated by its second-order moment truncation, which corresponds to the most common *weak assimilation* regime  
 where the information content of the innovation is small compared to that of the background, we have

$$\Sigma_{\rho}^{-1} \text{Cov}_{\zeta \sim \rho} [\zeta, \mathbf{a}_{\theta}(\mathbf{x}, \zeta)] = \int d\zeta \rho(\zeta) \nabla_{\zeta} \mathbf{a}_{\theta}(\mathbf{x}, \zeta) \approx \nabla_{\zeta} \mathbf{a}_{\theta}(\mathbf{x}, \zeta)|_{\zeta=0}, \quad (13)$$

which, with Eq. (12), relates the integral form  $\mathbb{E}_{\zeta \sim \rho} [\nabla_{\zeta} \mathbf{a}_{\theta}(\mathbf{x}, \zeta)]$  to the gradient form  $\nabla_{\zeta} \mathbf{a}_{\theta}(\mathbf{x}, \zeta)|_{\zeta=0}$  of the sensitivity of  
 $\mathbf{a}_{\theta}(\mathbf{x}, \zeta)$  with respect to  $\zeta$ .

255 Building on this identification and Eq. (9a), we can connect the *mean marginal error covariance matrix*

$$\tilde{\Gamma} \triangleq \mathbb{E}_{\mathbf{x} \sim \pi, \zeta \sim \rho} [\Gamma(\mathbf{x}, \zeta)] = \int d\mathbf{x} d\zeta \pi(\mathbf{x}) \rho(\zeta) \Gamma(\mathbf{x}, \zeta), \quad (14)$$

to the mean marginal gain  $\bar{\Gamma}$ :

$$\tilde{\Gamma} = \mathbb{E}_{\mathbf{x} \sim \pi, \zeta \sim \rho} [\Gamma(\mathbf{x}, \zeta)] = \mathbb{E}_{\mathbf{x} \sim \pi} [\nabla_{\mathbf{x}} \mathbb{E}_{\zeta \sim \rho} [\nabla_{\zeta} \mathbf{a}_{\theta}(\mathbf{x}, \zeta)]] \approx \mathbb{E}_{\mathbf{x} \sim \pi} [\nabla_{\mathbf{x}} \nabla_{\zeta} \mathbf{a}_{\theta}(\mathbf{x}, \zeta)|_{\zeta=0}] = \mathbb{E}_{\mathbf{x} \sim \pi} [\Gamma(\mathbf{x})] = \bar{\Gamma}. \quad (15)$$

However, the *empirical mean* of  $\Gamma(\mathbf{x}, \zeta)$  from which to evaluate  $\tilde{\Gamma}$ , and denoted  $\hat{\Gamma}$ , and which is the numerical estimation of the  
 260 sensitivity out of a long enough DA run, may differ from both theoretical means  $\tilde{\Gamma}$  and  $\bar{\Gamma}$ . That is why we show in Appendix D  
 how  $\hat{\Gamma}$  approximates  $\tilde{\Gamma}$ .

### 3.2.2 Invariance and equivariance

An important means to mitigate the computational cost and complexity in interpreting the mean marginal gain  $\tilde{\Gamma}$  is through the  
 symmetries of the DA process, if applicable. For instance, the one-dimensional periodic L96 model is translationally invariant  
 265 on the discretised circle. If, furthermore, the model's variables are homogeneously and homecedastically observed, then the  
 associated DA processes are also invariant with respect to these symmetries.



Let us generically denote  $\mathcal{G}$  such group of symmetries for the DA process, which are assumed to be isometries. It is chosen to be the maximal group for which both the dynamics and the observation process are equivariant (see Appendix G for details). We assume that for each symmetry  $g \in \mathcal{G}$  which acts in state space, there is a single induced symmetry  $g_y$  which acts in  
 270 observation space; and that the group  $\mathcal{G}_y$  of such induced isometries is isomorphic to  $\mathcal{G}$ .

The action of  $g \in \mathcal{G}$  on  $\Gamma(\mathbf{x}, \zeta)$  is denoted  $g \odot \Gamma(\mathbf{x}, \zeta)$ . It is a tensor field for  $\mathcal{G}$ , in the sense that it is *equivariant* under the action of this symmetry group following the transformation rule:

$$\forall g \in \mathcal{G}, \forall \mathbf{x} \in \mathbf{E}_x, \forall \zeta \in \mathbf{E}_x : \quad g \odot \Gamma(\mathbf{x}, \zeta) \triangleq \Gamma(g \circ \mathbf{x}, g \circ \zeta) = g \otimes g^T \otimes g^T \circ \Gamma(\mathbf{x}, \zeta), \quad (16)$$

where the ordering of the three tensorial factors follows the convention of Eq. (10). A proof of the equivariance Eq. (16) is  
 275 proposed in Appendix E.

Invoking the ergodicity of the dynamics, any symmetry of  $\mathcal{G}$  leaves the invariant distribution of the dynamics  $\pi$ , and the projected innovation distribution  $\rho$ , unchanged:

$$\forall g \in \mathcal{G}, \forall \mathbf{x} \in \mathbf{E}_x, \forall \zeta \in \mathbf{E}_x : \quad \pi(g \circ \mathbf{x}) = \pi(\mathbf{x}), \quad \rho(g \circ \zeta) = \rho(\zeta). \quad (17)$$

where  $g \circ (\cdot)$  denotes the action of  $g$  on a field.

280 Leveraging the equivariance of the marginal gain and the symmetries of the invariant distribution, we now demonstrate the invariance of the mean marginal gain  $\tilde{\Gamma}$ ; for  $g \in \mathcal{G}$ , we have

$$g \odot \tilde{\Gamma} = \int d\mathbf{x} d\zeta \pi(\mathbf{x}) \rho(\zeta) g \odot \Gamma(\mathbf{x}, \zeta) = \int d\mathbf{x} d\zeta \pi(\mathbf{x}) \rho(\zeta) \Gamma(g \circ \mathbf{x}, g \circ \zeta) \quad (18a)$$

$$= \int d(g^{-1} \circ \mathbf{x}) d(g^{-1} \circ \zeta) \pi(g^{-1} \circ \mathbf{x}) \rho(g^{-1} \circ \zeta) \Gamma(\mathbf{x}, \zeta) \quad (18b)$$

$$= \int d\mathbf{x} d\zeta \pi(\mathbf{x}) \rho(\zeta) \Gamma(\mathbf{x}, \zeta) = \tilde{\Gamma}, \quad (18c)$$

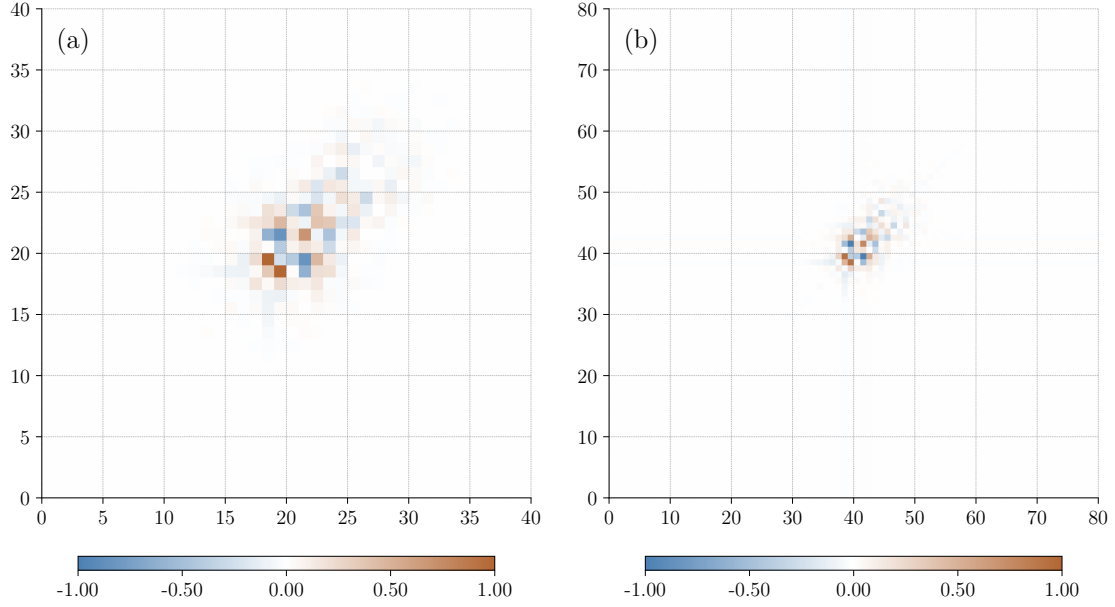
285 where a change of variables was carried out from Eq. (18a) to Eq. (18b). The invariance of  $\pi$  and  $\rho$  under  $\mathcal{G}$ , and the fact that the determinant of the Jacobian of  $g^{-1}$  is 1 since  $g^{-1}$  is an isometry were utilised from Eq. (18b) to Eq. (18c). We finally conclude:

$$\forall g \in \mathcal{G} : \quad g \odot \tilde{\Gamma} = \tilde{\Gamma}. \quad (19)$$

The same result can be obtained for  $\bar{\Gamma}$ , with a simpler derivation only involving the invariant distribution  $\pi$  of the underlying  
 290 dynamics. However, the symmetry group  $\bar{\Gamma}$  must be the same as that of  $\tilde{\Gamma}$ , and not the potentially larger group associated to  $\pi$ :

$$\forall g \in \mathcal{G} : \quad g \odot \bar{\Gamma} = \bar{\Gamma}. \quad (20)$$

As a consequence, in the rest of this section, we assume that the results indifferently applies to either  $\tilde{\Gamma}$  or  $\bar{\Gamma}$ .



**Figure 1.** Plot of the normalised mean marginal gain  $[\bar{\Omega}_r]_{ij} / \max |[\bar{\Omega}_r]_{ij}|$  in the fully observed DA configuration of the L96 model, obtained with a model size of  $N_x = 40$  (panel a) and  $N_x = 80$  (panel b).

### 3.2.3 Numerical illustrations

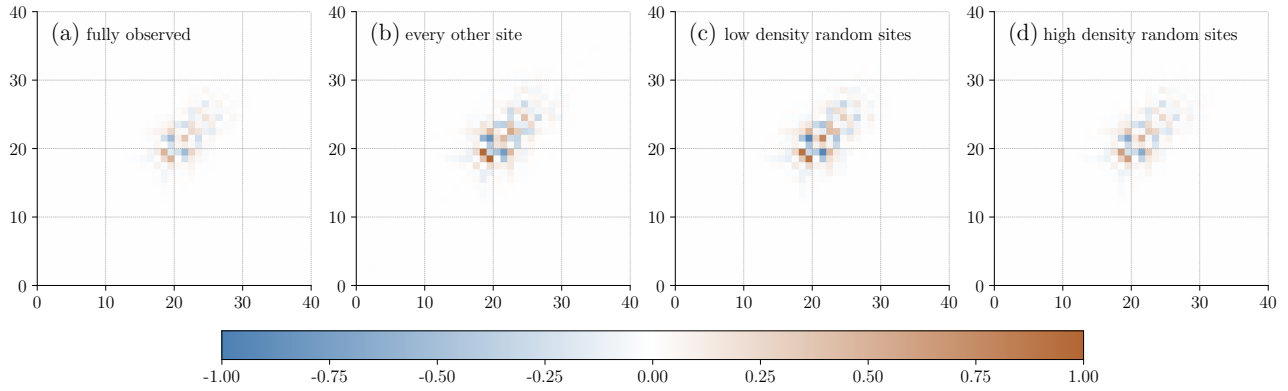
Leveraging the equivariance induced by the translational invariance of the L96 model,  $\Gamma(\mathbf{x})$  can be reduced to a 2-tensor  $\bar{\Omega}_r$  with  $r$  an arbitrary model grid point (or site) index; it is defined for all  $i, j$  by

$$[\bar{\Omega}_r]_{ij} \triangleq [\Gamma]_{ijr}. \quad (21)$$

For clarity,  $r$  is chosen to be in the middle of the domain  $r = \lfloor N_x/2 \rfloor$  in the L96 case. The simple but tedious details justifying Eq. (21) are given in Appendix F. How to numerically compute this  $\bar{\Omega}_r$  matrix is then discussed in Appendix G.

To start with, we choose the fully observed DA configuration  $\mathbf{H} = \mathbf{R} = \mathbf{I}_x$ , which applies to both the training of  $\mathbf{a}_\theta$  and the subsequent DA tests. The computation of  $\bar{\Omega}_r$  is carried out through the composite approach (see Appendix G). The results are shown in Fig. 1.

The dominant values of  $|\bar{\Omega}_r|$  form a pattern. They are concentrated in the vicinity of the perturbed variable of the forecast state, which makes them local. To further qualify the local patterns, we also trained  $\mathbf{a}_\theta$  on a L96 model but with  $N_x = 80$ , i.e. twice the size of the standard L96 model, while all other parameters either related to the dynamics or the DA experiments, remain the same. Then,  $\bar{\Omega}_r$  is similarly computed and plotted in Fig. 1. As expected, the same pattern emerges with the same spatial extension. This further supports the EnKF-like performance of  $\mathbf{a}_\theta$  trained with  $N_x = 40$  when tested with  $N_x = 80$  (as recalled in Sect. 2.3).



**Figure 2.** Plot of the mean marginal gain  $[\bar{\Omega}_r]_{ij}$  normalised by the maximum of its absolute value over all entries and over the four panels, in several sparse observation DA configuration of the L96 model.

### 3.2.4 Patterns with non-trivial observation operators

We carried out the exact same experiments described above, but now with three different observation configurations instead of the full observation setup. The corresponding  $\bar{\Omega}_r$  are plotted in Fig. 2. Panel (a) corresponds to the fully observed configuration for reference. The first configuration (panel b) corresponds to the observation of every other site:  $[\mathbf{H}]_{ji} = \delta_{i,2j}$  for  $1 \leq i \leq N_x$  and  $1 \leq j \leq \lfloor N_x/2 \rfloor$ . The second (panel c) corresponds, at each time step  $\tau_k$ , to an observation at  $N_y^k$  random but distinct sites, where  $N_y^k$  is uniformly drawn at each time step in between 0 and  $N_x$  (bounds included). The third configuration (panel d) is the same as the second but with  $N_y^k$  uniformly drawn at each time step in between  $\lfloor N_x/2 \rfloor$  and  $N_x$ . It is remarkable that the same local pattern emerges in all configurations, even when one every other site is never observed. However, the magnitude of the sensitivities (values of the patterns) changes depending on the information balance in the analysis and hence in the gain magnitude. Finally, note that  $\bar{\Omega}_r$  only represents an average pattern. It is possible to exhibit a family of patterns but it would go beyond the aim of the present paper.

## 4 Data assimilation networks in stronger nonlinear regimes

In Sect. 3, we made several contributions to the understanding of DAN. The numerical experiments were set in a mild nonlinear regime of L96 where  $\Delta_t = 0.05$  in between analyses, known to correspond to 6 hours of a synoptic meteorological model (Lorenz and Emanuel, 1998), and a forcing  $F$  set to 8. We now examine the ability of the method to learn efficient analysis schemes under stronger nonlinear conditions. While we exclude cases in which the nonlinearity is so severe that it induces implicit or explicit recurrent multi-modal priors, we do consider regimes that exhibit substantial departures from mild nonlinearity.

Probing such mild to stronger nonlinearity regimes can be achieved by less frequent observation, typically increasing the update time-step  $\Delta_t$ . The dimensional analysis of Appendix 1 in Bocquet and Carrassi (2017) shows that varying  $\Delta_t$  is indeed



relevant to achieve such objective for the L96 model versus, e.g., increasing the observation error amplitude. The forcing  $F$  can also be varied to that end. However, it rather stands as a signature of the magnitude of the instability of the dynamics (e.g., a covariant function of the Kaplan-Yorke dimension, a measure of the fractal dimension of the dynamics attractor) rather than the signature of the deviation from Gaussianity.

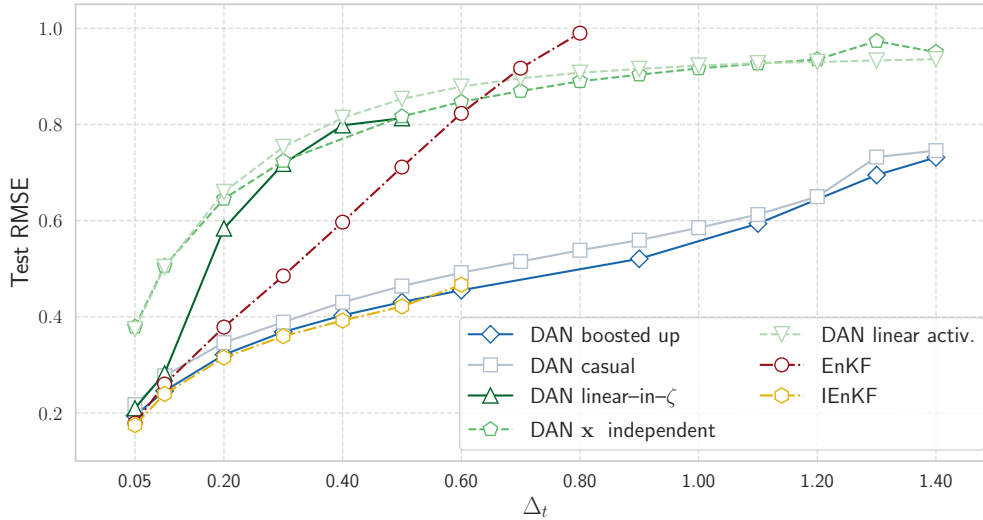
#### 4.1 Performance as a function of the update time-step

Increasing  $\Delta_t$  to multiples of 0.05 is the experimental design chosen in Sakov et al. (2012); Bocquet and Sakov (2012) to evaluate the performance of the iterative Ensemble Kalman Filter (IEnKF). As the ensemble variant of the iterative Kalman filter (Wishner et al., 1969; Jazwinski, 1970), the IEnKF still stands, to our best knowledge, as the most accurate scalable DA method in mild to stronger nonlinear conditions (Bocquet and Sakov, 2013). It is hence a hard-to-beat baseline for learning an advanced DA analysis scheme in such conditions. The remarkable performance of the IEnKF stems from its ensemble-variational formulation obtained from Bayesian first principles: an ensemble is used to construct a time-dependent prior, and the analysis is performed through a nonlinear iterative optimisation (Bocquet and Sakov, 2014).

The IEnKF can be made even more accurate (both for smoothing and filtering) by choosing longer DA windows, yielding the iterative ensemble Kalman smoother (IEnKS, Bocquet and Sakov, 2014). Since the IEnKS naturally derives from the IEnKF and can be more costly, we focus here on the IEnKF as a sufficiently strong baseline.

We choose a simple common observational configuration for the many DA methods we intend to compare. All sites are observed through  $\mathbf{H}_k \triangleq \mathbf{I}_x$  and  $\mathbf{R}_k \triangleq \mathbf{I}_x$ , a well documented setup in the literature. In this configuration, any useful, but not necessarily accurate, DA method must exhibit a test RMSE slightly below 1. We consider the following DA methods:

- A well tuned EnKF with an ensemble of size  $N_e = N_x = 40$ . Optimal multiplicative inflation is addressed through the finite-size EnKF (Bocquet, 2011; Bocquet et al., 2015). This first contender is meant to illustrate the progressive failure of the EnKF in stronger nonlinear conditions.
- A well tuned IEnKF with an ensemble of size  $N_e = N_x = 40$ , whose optimal multiplicative inflation is addressed through the finite-size IEnKF (Bocquet and Sakov, 2012). This is our hard baseline.
- A *casual* learned DAN scheme using the neural network and the training parameters values set in Appendix A.
- A *boosted* learned DAN scheme using  $N_f = 80$ ,  $N_r = 2^{19}$ ,  $N_{\text{iter}} = 32$ ,  $S_b = 512$ , which is parameter and data intensive, and hence much more time-consuming to train on a single GPU, while its inference remains cheap.
- A learned DAN scheme where the explicit dependence on the forecast state is discarded, while the dependence on the projected innovations is maintained. We expect the resulting DA method to perform similarly to a well tuned 3D-Var, see Boc24.
- A learned DAN scheme where the activation functions are all linear, such that  $\mathbf{a}_\theta$  is linear in its inputs. Like the previous approach, we expect the resulting DA method to perform similarly to a well tuned 3D-Var, see Boc24.



**Figure 3.** Test RMSEs of DA methods as a function of the update time-step  $\Delta_t$  in between analyses. See text for details.

To avoid early divergences in the training when  $\Delta_t \gg 0.05$ , all the DAN schemes benefited from a modified Eq. (2):

$$\mathbf{x}_k^a = \mathbf{x}_k^f + \alpha \mathbf{H}_k^T \mathbf{R}_k^{-1} (\mathbf{y}_k - \mathbf{H}_k \mathbf{x}_k^f) + \mathbf{a}_\theta (\mathbf{x}_k^f, \mathbf{H}_k^T \mathbf{R}_k^{-1} \delta_k), \quad (22)$$

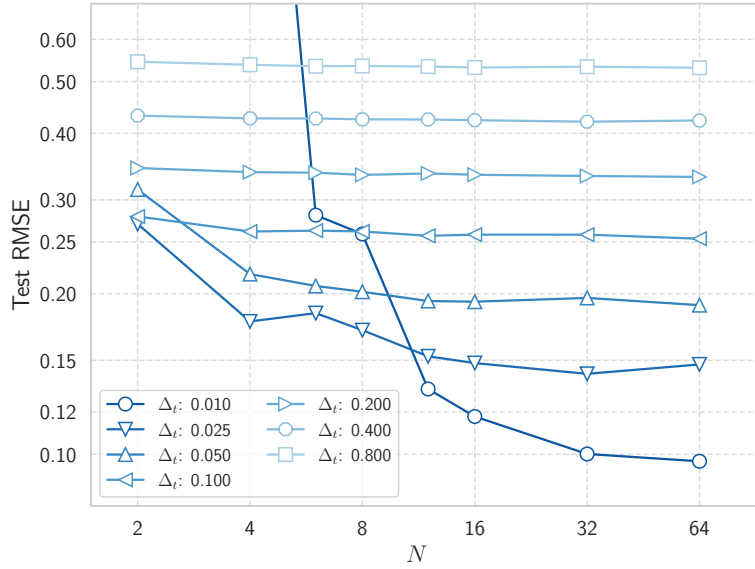
where  $\alpha$  is a trainable scalar. Note that this formulation is mathematically equivalent to the original Eq. (2). This is only meant to explicitly enforce the solution  $\mathbf{x}_k^a = \mathbf{y}_k$  when  $\mathbf{H}_k = \mathbf{I}_x$ , which guides the training in its first few epochs.

The test RMSEs of those DA methods, over a long DA run, are displayed in Fig. 3 as a function of  $\Delta_t$ . Note that  $\Delta_t = 0.60$  is already very significantly nonlinear as it corresponds to the Lyapunov time of L96, i.e. the time horizon beyond which the DA system becomes significantly non-Gaussian. Note that, beyond  $\Delta_t = 0.60$ , the IEnKF is trickier to stabilise and hence its performance is not reported. Moreover, passed  $\Delta_t = 0.80$ , the EnKF becomes uninformative and its test RMSE is not reported.

As expected, the two degraded DAN operators severely underperform the other DA schemes, that leverage non-static priors, with a test RMSE that ranges from 0.38 for  $\Delta_t = 0.05$  to an asymptotic RMSE below 1 for much larger  $\Delta_t$ . This is consistent with the findings of Boc24 and mirrors the performance of a 3D-Var with static background covariance matrix.

The EnKF offers a very good performance in the mild nonlinear regime  $\Delta_t \approx 0.05$  but gradually degrades as  $\Delta_t$  is increased. The method becomes uninformative beyond  $\Delta_t = 0.80$ . As already reported by Sakov et al. (2012), the IEnKF offers significantly better performance, from a marginal improvement over the EnKF at  $\Delta_t = 0.05$  that gets more and more significant as  $\Delta_t$  is increased.

Remarkably, the boosted DAN operator achieves performance very similar to the IEnKF, but can still be learned for much larger  $\Delta_t$ , and remains informative with still a very significant edge over the static prior methods. Again, this is achieved without the use of an ensemble but of a single forecast state. Moreover, the learned DAN does not explicitly resort to a nonlinear (Gauss-Newton) iterative minimisation, in contrast to the IEnKF. This aspect of such DAN is reminiscent of approaches meant



**Figure 4.** Test RMSEs of DANs as a function of both  $\Delta_t$  and  $N_{\text{iter}}$ .

to learn a solver for a variational DA problem (Fablet et al., 2021; Frerix et al., 2021; Lafon et al., 2023; Filoche et al., 2023; Keller and Potthast, 2024). The casual DAN operator is slightly less performing but follows the same trend. We could have  
 380 evaluated the methods for even larger  $\Delta_t > 1.40$ ; however,  $\Delta_t = 1.20$  already stands for twice the Lyapunov time, which is equivalent to 6 days in the L96 correspondence.

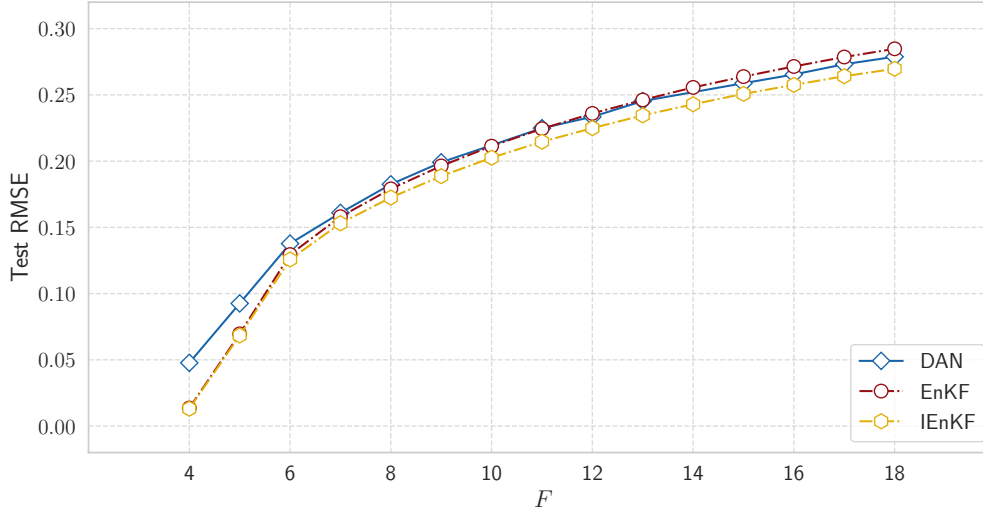
We suspect that the test RMSEs should primarily be a function of  $N_{\text{iter}}\Delta_t$  rather than just  $\Delta_t$ , accounting for the forgetful effect associated to the chaotic dynamics. Hence, to reach the same test RMSEs,  $N_{\text{iter}}$  could be roughly chosen inversely proportional to  $\Delta_t$ . To test this hypothesis, we compare the test RMSE on a large number of trained DANs, varying  $\Delta_t$   
 385 (including for  $\Delta_t \leq 0.05$ ) and  $N_{\text{iter}}$  with results shown in Fig. 4. Leveraging the findings of Appendix A and the use of smaller batches to mitigate the numerical cost, we have chosen for their training,  $N_r = 2^{16} = 65,536$  and  $S_b = 2^7 = 128$ . The test RMSEs are shown in Fig. 4. The results corroborate the intuition with the weaker and weaker dependence on  $N_{\text{iter}}$  of the performance when  $\Delta_t$  is increased. Conversely, a much larger  $N_{\text{iter}}$  is required when  $\Delta_t$  gets very small ( $\Delta_t \approx 0.01$ ).

Hence, DAN is numerically more difficult to train when  $\Delta_t < 0.05$  requiring larger  $N_{\text{iter}}$ , which corresponds to quasi-linear  
 390 regimes where DAN is nonetheless of limited interest.

## 4.2 Performance as a function of the energy forcing

The energy forcing  $F$  pumps in and out energy off the L96 dynamics and feeds their instability. Besides  $\Delta_t$ , this is another tunable parameter of the nonlinearity of the dynamics. There is a wealth of dynamical phenomenology of the dynamics for





**Figure 5.** Test RMSEs of DA methods as a function of the forcing  $F$ . See text for details.

a wide range of  $F$  (see, e.g., Barone et al., 2025, and references therein). For the range of  $F$  we focus on here, instabilities progressively develop in between  $1 \lesssim F \lesssim 4$ , while chaos fully sets in beyond  $F \gtrsim 4$  (van Kekem and Sterk, 2018). The number of Lyapunov exponents, and similarly the Kaplan-Yorke dimension, increases monotonically from 4 to about 30 where it saturates (Karimi and Paul, 2010).

We choose the same experimental setup as in the previous experiments and consider the following DA methods:

- A well tuned EnKF with an ensemble of size  $N_e = N_x = 40$ . Optimal multiplicative inflation is addressed through the finite-size EnKF (Bocquet, 2011; Bocquet et al., 2015), in its Dirac-Jeffreys variant (Bocquet et al., 2015) required to handle the weakly nonlinear regime ( $4 \lesssim F \lesssim 8$ ).
- A well tuned IEnKF with an ensemble of size  $N_e = N_x = 40$ , whose optimal multiplicative inflation is addressed through the finite-size IEnKF (Bocquet and Sakov, 2012) in its Dirac-Jeffreys variant. This is our hard baseline.
- A (casual) learned DAN scheme using the neural network and training parameters as defined in Appendix. A.

The test RMSEs are displayed in Fig. 5. The IEnKF has a slight edge over the EnKF with larger  $F$  with the increasing Kaplan-Yorke dimension. Hence, the DA system does not deviate much from non-Gaussianity as  $F$  increases, only the magnitude of the instabilities are. The DAN achieves a performance in between that of the EnKF and that of the IEnKF, which is patent for large  $F$ . Note that it turns out trickier to train DANs for  $F$  getting close to  $F = 4$ . There, the dynamics become more and more laminar and exhibit almost periodic waves, with patterns that, although simpler, are very different from those learned in the regimes explored so far.



### 4.3 Reasons for the success of data assimilation networks

In the light of the previous numerical results, we discuss the reasons why DAN can be as accurate as the IEnKF, without an ensemble, without any experimental tuning, and even in mild to strong nonlinear conditions.

We have already shown in Sect. 3.1 that in the mild nonlinear regime, i.e.  $\Delta_t \approx 0.05$ , the success of DAN mainly resides  
 415 in its implicit estimation of the mapping  $\mathbf{x}^f \mapsto \mathbf{P}^f(\mathbf{x}^f)$ , in line with the results by Sacco et al. (2024); Sakov (2025). As  
 showed in Sect. 3.1.1, this amounts to assume a non-static but Gaussian prior in the analysis: this is equivalent to having  
 $\|\mathbf{x} - \mathbf{x}^f\|_{(\mathbf{P}^f(\mathbf{x}^f))^{-1}}^2$  as the background term in the analysis cost function. Whether this mechanism is sufficient to ensure the  
 performance of DAN when  $\Delta_t \geq 0.05$  is doubtful.

#### 4.3.1 The linear-in- $\zeta$ data assimilation network beyond mild nonlinearity

420 To address this question, let us assess the linear-in- $\zeta$  DAN, see Sect. 3.1, in stronger nonlinear regimes  $\Delta \geq 0.05$ . We could  
 anticipate that it accounts well for the errors of the day, but that it may not be able to handle stronger nonlinearity/non-  
 Gaussianity, e.g., if  $\Delta_t \gg 0.05$ , because of the implied Gaussian prior. It should hence match the EnKF, rather than the IEnKF.

Let us check that hypothesis numerically. The test RMSEs of this specific DAN scheme, which mirrors Eq. (5), are reported  
 in Fig. 3, as the linear-in- $\zeta$  DAN. Let us first remark that, in accordance with the claims of Sect. 3.1, it performs as well as  
 425 the EnKF for  $\Delta_t = 0.05$ , again corroborating the results by Sacco et al. (2024); Sakov (2025). However, as  $\Delta_t$  increases, its  
 performance significantly degrades compared to the EnKF, not to mention the IEnKF. Even though it relies on an estimate of  
 the  $\mathbf{x}^f \mapsto \mathbf{P}^f(\mathbf{x}^f)$  map, its underlying Gaussian assumption penalises it beyond the mild nonlinear regime, as expected.

#### 4.3.2 Deviation of the data assimilation network prior from Gaussianity

We conclude that the full DAN scheme not only implicitly learns the  $\mathbf{x}^f \mapsto \mathbf{P}^f(\mathbf{x}^f)$  map but also a non-Gaussian prior, which  
 430 is more informative than the Gaussian prior associated to the analysis cost function background term  $\|\mathbf{x} - \mathbf{x}^f\|_{(\mathbf{P}^f(\mathbf{x}^f))^{-1}}^2$ . As  
 shown, both abilities are required for DAN to perform so well in the extended range of mild to strong nonlinear regimes.

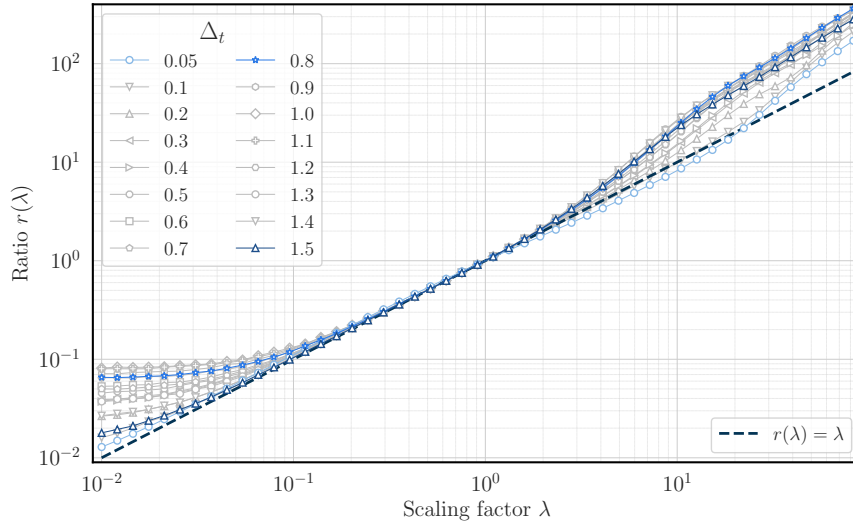
We investigate the deviation of  $\mathbf{a}_\theta$  from the presumed quasi-linearity in  $\zeta$  by defining the scalar function

$$r(\lambda) = \mathbb{E}_{\mathbf{x}, \zeta} \left[ \frac{\|\mathbf{P}^a(\mathbf{x})^{-1} (\mathbf{a}_\theta(\mathbf{x}, \lambda\zeta) - \mathbf{a}_\theta(\mathbf{x}, \mathbf{0}))\|}{\|\mathbf{P}^a(\mathbf{x})^{-1} (\mathbf{a}_\theta(\mathbf{x}, \zeta) - \mathbf{a}_\theta(\mathbf{x}, \mathbf{0}))\|} \right], \quad (23)$$

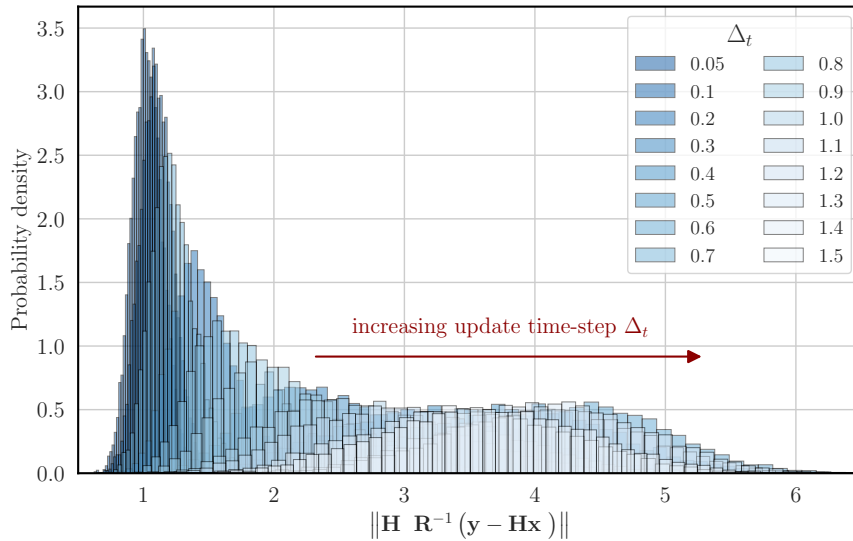
where  $\|\cdot\|$  is the Euclidean norm. It depends on a dimensionless scale parameter  $\lambda > 0$  and measures deviations from the  
 435 linearity in  $\zeta$ . In Eq. (23),  $\mathbf{P}^a(\mathbf{x})$  is meant to standardise the deviations from linearity. Indeed, for the quasi-linear regime, we  
 have

$$\mathbf{P}^a(\mathbf{x})^{-1} (\mathbf{a}_\theta(\mathbf{x}, \lambda\zeta) - \mathbf{a}_\theta(\mathbf{x}, \mathbf{0})) \approx \lambda\zeta, \quad (24)$$

such that  $r(\lambda) \simeq \lambda$  should hold. The ratio  $r(\lambda)$  is estimated using perturbations as in Boc24:  $\mathbf{x}$  and  $\zeta$  are sampled from the  
 forecast state  $\mathbf{x}^f$  and the projected innovations  $\zeta = \mathbf{H}^T \mathbf{R}^{-1} (\mathbf{y} - \mathbf{H}\mathbf{x}^f)$  that are obtained from a long trajectory.



**Figure 6.** Scaling of the analysis increments of  $\mathbf{a}_\theta$  in the projected innovations, for a large range of  $\Delta_t$  values.



**Figure 7.** Histograms of the projected innovations, for a large range of  $\Delta_t$  values.

440 The nonlinearity is demonstrated in Fig. 6 with the function  $\lambda \mapsto r(\lambda)$  from Eq. (23) shown for several values of  $\Delta_t$  in the range  $[0.05, 1.5]$ . These curves should be appreciated knowing the range of values taken by the projected innovations in a long DA run for the selected values of  $\Delta_t$  which, for each  $\Delta_t$ , points to the relevant range of  $\lambda$  values to consider and which mainly contributes to the computation of  $r(\lambda)$ . This is shown in Fig. 7 in the form of histograms of those values. Given the range of innovation values and the logarithmic scale of Fig. 6, the most relevant part of the scaling behaviour sits in the range  $\lambda \in [1, 10]$ .



445 In this range of the scaling, as seen in Fig. 6, the more nonlinear the DA run, the steeper  $\lambda \mapsto r(\lambda)$ : the innovation impact is stronger than the one expected in the linear regime, especially for large magnitude innovations.

## 5 Discussion and conclusions

In this paper, we have continued to explore the potential of learning sequential data assimilation operators with neural networks for tracking chaotic dynamical systems, following on the initial results and conclusions of Boc24 who built on the concepts of  
450 data assimilation networks (DAN) introduced by McCabe and Brown (2021); Boudier et al. (2023). Assuming the dynamics to be known, the focus is on learning the analysis. Compared to learning the analysis from a dataset of inputs and outputs of DA runs, the training is numerically challenging since the analysis operator is learned through several DA cycles and from trajectories of the dynamics and its observations only.

The resulting DANs are robust, in the sense that they do not require inflation and any other correction and regularisation.  
455 They can operate without an ensemble and only use the forecast state as prior information. And yet, they are as accurate as a well tuned EnKF in mildly nonlinear regime, and as accurate as a well tuned IEnKF from mild to stronger nonlinear conditions.

### 5.1 Abilities of data assimilation networks

We have previously shown that, to achieve this level of performance, a data assimilation network (DAN) must implicitly learn a map  $\mathbf{x}^f \mapsto \mathbf{P}^f(\mathbf{x}^f)$ , whose existence is supported by a multiplicative ergodic theorem applied to the entire DA process  
460 viewed as a dynamical system. The network must implicitly identify spatial local patterns in  $\mathbf{x}^f$  in order to internally represent components of  $\mathbf{P}^f(\mathbf{x}^f)$ , a property that makes the learned DA method scalable.

In this paper, we have further examined the reliance of DANs on the map  $\mathbf{x}^f \mapsto \mathbf{P}^f(\mathbf{x}^f)$  by constructing an ad hoc learnable DAN that explicitly incorporates this mapping. For the L96 dynamics, we identified an average characteristic local pattern by leveraging both the invariant distribution of the dynamics and the translational symmetry of the model. We also showed  
465 that DAN can learn non-Gaussian priors that depend solely on  $\mathbf{x}^f$ , and that these priors are necessary for DAN to match the performance of the IEnKF under more stringent nonlinear conditions. In this regime, we demonstrated that implicit or explicit knowledge of the mapping alone is insufficient. Thus, both mechanisms must operate within DAN, paralleling IEnKF's critical reliance on its ensemble for flow-dependent error estimation and on its Gauss-Newton iterative solver to probe departures of the dynamics from linearity.

### 470 5.2 Implication for the end-to-end processors in numerical weather prediction

Taken in the context of past DA literature, these results are somewhat surprising. For a long time, it was believed that an ensemble was essential for estimating flow-dependent errors. Moreover, the proper construction of non-Gaussian priors has remained a long-standing challenge, whereas even a simple DAN can rapidly learn an effective one. These findings should inform the (re-)design of future DA algorithms.



475 These considerations have implications for what DAN-like operators are capable of achieving. For example, recently de-  
 veloped end-to-end atmospheric processors (McNally et al., 2024; Alexe et al., 2024; Allen et al., 2025; Lean et al., 2025;  
 Laloyaux et al., 2025) that ingest observations and predict future observations, implicitly construct their own internal repre-  
 sentation of the system and repeatedly compare this latent state to newly assimilated observations. In doing so, they implicitly  
 learn an analysis operator in their latent space. Lean et al. (2025) concluded that their processor, GraphDOP, must be able to  
 480 learn not only a climatological background but also dynamical priors, a result that may seem surprising given that GraphDOP  
 does not rely on any explicit background information. However, in light of our results, GraphDOP must internally learn a map-  
 ping from the 12-hour observational snapshot used as input of the processor to an estimate of the underlying error covariances.  
 Consequently, it must be able to construct its own background, incorporating advanced error statistics that go beyond a mere  
 climatology.

485 This challenges the view that the accuracy of such end-to-end processors is fundamentally limited by the absence of a back-  
 ground (such as a forecast ensemble or climatological information) in their inputs. While an explicit background representation  
 would certainly provide additional information, the extent of the achievable performance gains remains a subtler question.

*Code availability.* The key algorithmic pieces of code will be made publicly available upon acceptance of the manuscript.

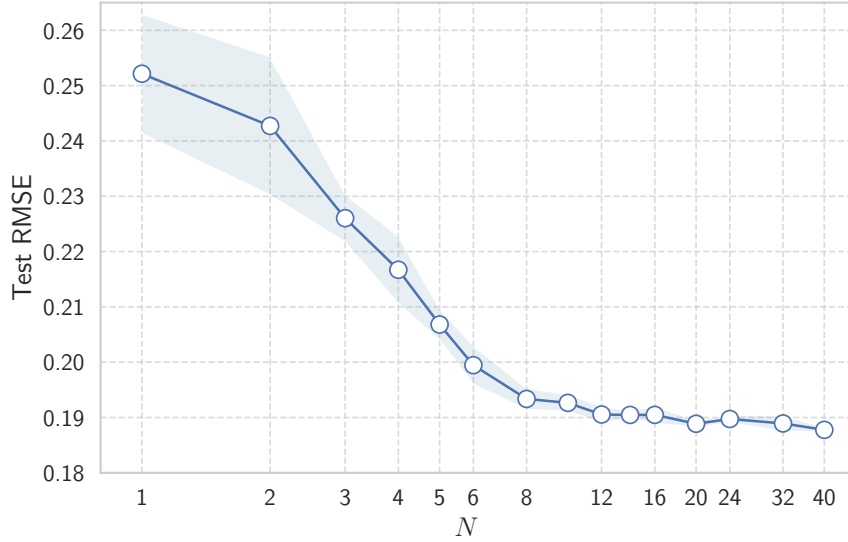
## Appendix A: Sensitivity to the batch size, the size of the datasets, and the backpropagation truncation

490 The key parameters in the design of  $\mathbf{a}_\theta$  and in its training are:

- the number of channels/filters  $N_f$  processed by the convolutional neural network. This essentially gives away the com-  
 plexity of the neural network, how many features it can identify and process.
- the number of trajectory  $N_r$  processed in parallel. The larger  $N_r$ , the more DA runs the neural network can learn from  
 and be made robust against. 10% of these are reserved for validation.
- 495 – the number of cycles  $N_{\text{iter}}$  through which the gradient is computed. This corresponds to the limit imposed by the  
 truncated backpropagation through time. The larger  $N_{\text{iter}} \ll N_c$ , the better the information transmission from one cycle  
 to the next should be learned, but the more costly and less accurate the gradients.
- the size of the batch  $S_b$ . Hence, the number of steps in each epoch is about  $N_r/S_b$  (training and validation).

The dependence of the performance of  $\mathbf{a}_\theta$  on those parameters was studied in Boc24. The values  $N_{\text{iter}} = 16$ ,  $N_r = 2^{18}$ ,  $S_b =$   
 500 2048 were chosen for the  $\mathbf{a}_\theta$  hyperparameters of the reference configuration, as a compromise between training speed and  
 accuracy of the resulting  $\mathbf{a}_\theta$ . However, the dependence on  $N_{\text{iter}}$  and  $S_b$  was barely reported and discussed, so that we focus on  
 them in what follows.

Using this reference configuration, we compute the test RMSE of the DAN schemes as a function of the truncation number  
 $N_{\text{iter}}$ . The results are shown in Fig. A1. As expected,  $N_{\text{iter}} \leq 5$  prevents DAN from learning an efficient prior that relies on the



**Figure A1.** Test RMSE of DAN as a function of the truncation cycles number  $N_{\text{iter}} = 1, \dots, 40$  in the truncated backpropagation. For each value of  $N_c$ , an ensemble of 5  $\mathbf{a}_\theta$  operators is learned. Plus and minus one standard deviation of the RMSE are displayed as shades around the RMSE curve.

errors of the day. It is however remarkable that the improvement in the test RMSE as  $N_{\text{iter}}$  increases is noticeable up to about  $N_{\text{iter}} \simeq 40$ , which corresponds to about 6 times the doubling time of the L96 model in the reference configuration.

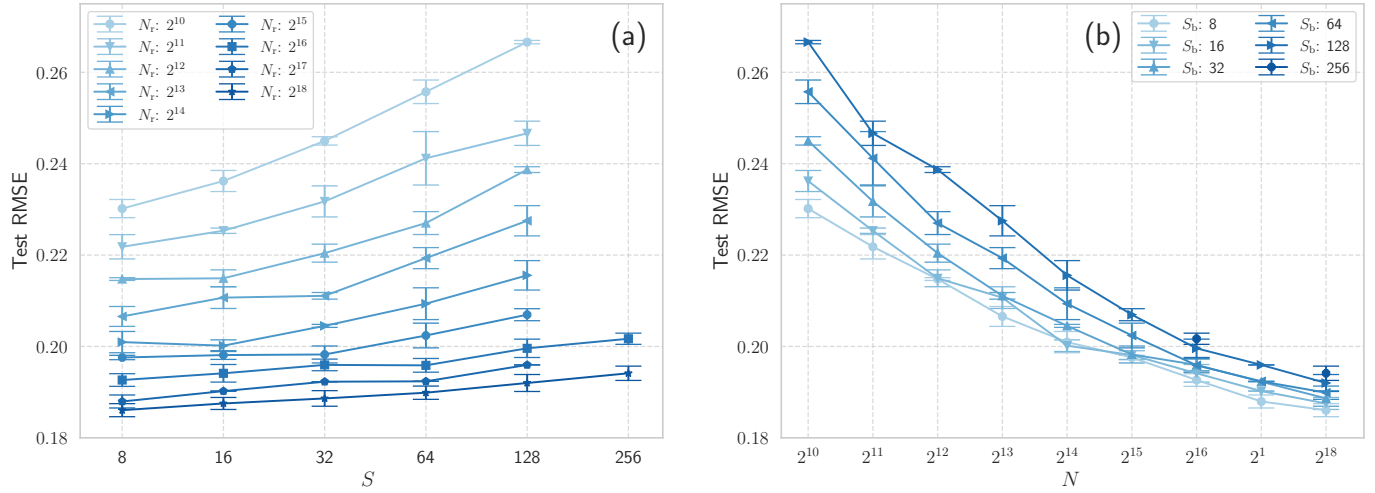
We have also experimented with the batch size  $S_b$  much more thoroughly than in Boc24. We found that using smaller batches is beneficial to the performance of the trained  $\mathbf{a}_\theta$ . This may also enable reducing the number of trajectories  $N_r$  in the dataset. Nonetheless, we empirically found that the number of steps, i.e.  $\sim N_r/S_b$  in each epoch still needs to remain large to achieve high accuracy. Experimenting, we learned a large set of  $\mathbf{a}_\theta$  operators with  $N_r$  in the range  $2^{10} - 2^{18}$ , while  $S_b$  is chosen in the range  $2^3 - 2^8$ . The corresponding test RMSEs are plotted as a function of either  $S_b$ , as shown in Fig. A2a, or  $N_r$  as shown in Fig. A2b.

For this specific configuration a scaling law can be numerically estimated. The following Ansatz is assumed:

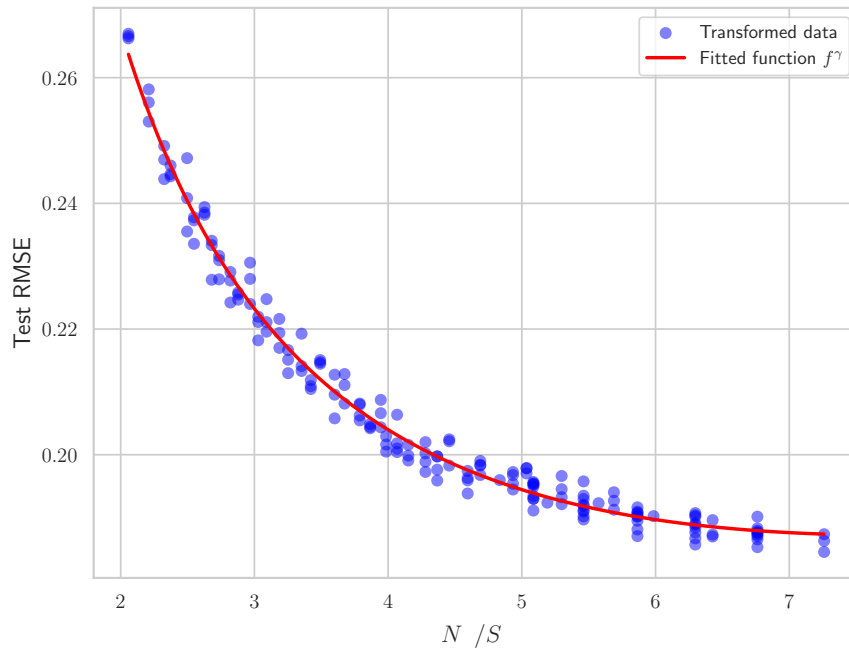
$$\text{RMSE}^{1/\gamma} = f\left(\frac{N_r^\alpha}{S_b^\beta}\right), \quad (\text{A1})$$

where  $f$  is a polynomial of order 2. This Ansatz (3 polynomial coefficients and 3 exponents) is fitted to the test RMSE results. The fit relevance can be visualised through the plot of  $f^\gamma$  in Fig. A3. The fitted exponents are  $\gamma = 1.435$ ,  $\alpha = 0.103$ , and  $\beta = 0.060$ . What really matters is the ratio  $\rho = \beta/\alpha = 0.58$ , since the test RMSE turns out to show a strong dependence on  $N_r/S_b^\rho$ .

In practice, the interest of using smaller batches  $S_b$  and hence smaller  $N_r$  must be weighted against the overheads created by the dataset pipeline but also by the transfer of the batches from RAM to VRAM. It is however possible to run several training experiments in parallel over the same GPU with smaller batches. Hence, the efficiency of using mini-batches is very dependent



**Figure A2.** Test RMSEs of DAN as a function of the batch size  $S_b$  for a selection of  $N_r$  values (panel a) and as a function of the number of dataset trajectories  $N_r$  for a selection of  $S_b$  values (panel b). For each pair  $(S_b, N_r)$ , an ensemble of 3  $\mathbf{a}_\theta$  operators is learned, from which error bars (plus or minus the standard deviation) are estimated and added to the curve plots.



**Figure A3.** Checking the relevance of the functional regression by plotting the fit function  $f^\gamma$  with respect for  $N_r^\alpha / S_b^\beta$  (curve), and of the test RMSE results (dots).





on the accelerator device(s) and the type of experiment to conduct. Note that, when training  $\mathbf{a}_\theta$  on much higher dimensional systems, relying on mini-batches may turn mandatory so as to fit into VRAM.

Finally, the scaling law Eq. (A1) can be leveraged to reduce both  $N_r$  and  $S_b$  and still be able to extrapolate to values yielding better test RMSEs, using the dependence on  $N_r/S_b^\rho$ . This scaling could nonetheless change with a different DA configuration.

## Appendix B: Neural network implementation of the $\mathbf{x}^f \mapsto \mathbf{P}^f(\mathbf{x}^f)$ map

Here, we detail how to build the linear-in- $\zeta$  DAN meant to enforce Eq. (8b) in Sect. 3.1. Implementing the map  $(\mathbf{x}^f, \zeta) \mapsto \mathbf{P}^a(\mathbf{x}^f) \cdot \zeta$  as a neural network is non-trivial if we wish to make it scalable. First,  $\mathbf{P}^a(\mathbf{x}^f)$  is symmetric positive definite; to enforce such constraint, the usual most efficient approach is to write:  $\mathbf{P}^a(\mathbf{x}^f) = \mathbf{X}^a(\mathbf{x}^f) \mathbf{X}^a(\mathbf{x}^f)^\top$ , where  $\mathbf{X}^a(\mathbf{x}^f)$  is a matrix of anomalies which may be easier to interpret than  $\mathbf{P}^a(\mathbf{x}^f)$ . The first difficulty towards scalability is the fact that  $\mathbf{P}^a(\mathbf{x}^f)$ , or  $\mathbf{X}^a(\mathbf{x}^f)$ , is of size  $N_x \times N_x$ , which is considered non-scalable. However, if  $\mathbf{P}^a(\mathbf{x}^f)$  can be approximated as low-rank, then  $\mathbf{P}^a(\mathbf{x}^f) \approx \mathbf{X}^a(\mathbf{x}^f) \mathbf{X}^a(\mathbf{x}^f)^\top$  with  $\mathbf{X}^a(\mathbf{x}^f)$  of size  $N_x \times N_r$ , with  $N_r$  remaining small enough when  $N_x$  is increased. Unfortunately,  $\mathbf{P}^a(\mathbf{x}^f)$  should not realistically be considered low-rank. However, one can exploit the locality of the covariances and write  $\mathbf{P}^a(\mathbf{x}^f) \approx \boldsymbol{\rho} \circ (\mathbf{X}^a(\mathbf{x}^f) \mathbf{X}^a(\mathbf{x}^f)^\top)$ , where  $\boldsymbol{\rho}$  is the localisation correlation matrix and  $\circ$  is the Schur/Hadamard product. In addition to making  $\mathbf{P}^a(\mathbf{x}^f)$  full rank in spite of a manageable number of parameters in  $\mathbf{X}^a(\mathbf{x}^f)$ , it also tapers spurious correlations that could be learned in the course of the training. Indeed, the implementation of the localisation significantly accelerates the training. This is reminiscent of the proposal by Bocquet and Farchi (2019) to estimate  $\mathbf{X}^a$  through a loss involving a Schur product with  $\boldsymbol{\rho}$ .

Furthermore, given  $\mathbf{P}^a = \boldsymbol{\rho} \circ (\mathbf{X}^a (\mathbf{X}^a)^\top)$  and the projected innovation  $\zeta$ , the implementation of such mapping can be efficiently coded using (see, e.g., Desroziers et al., 2014)

$$\mathbf{P}^a \cdot \zeta = \boldsymbol{\rho} \circ (\mathbf{X}^a (\mathbf{X}^a)^\top) \cdot \zeta = \sum_{l=1}^{N_r} \mathbf{X}_l^a \circ (\boldsymbol{\rho} \cdot (\mathbf{X}_l^a \circ \zeta)), \quad (\text{B1})$$

where the  $\mathbf{X}_l^a$  are the columns of  $\mathbf{X}^a$ , and  $\cdot$  denotes the usual matrix/vector multiplication. The matrix multiplication by the localisation matrix  $\boldsymbol{\rho}$  is scalable since  $\boldsymbol{\rho}$  is assumed to be a banded matrix. With L96 in mind, i.e. in a one-dimensional context, with a localisation matrix support of (band-)width  $N_l$ , the numerical complexity of  $\mathbf{P}^a \cdot \zeta$  is  $N_r N_x (2 + N_l)$ .

If localisation is not useful and  $\mathbf{P}^a(\mathbf{x}^f)$  is low-rank, then it is easy to implement  $\mathbf{P}^a \cdot \zeta = \mathbf{X}^a \cdot ((\mathbf{X}^a)^\top \cdot \zeta)$ , which is reminiscent of the very popular machine learning *attention* mechanism. The numerical complexity is then  $2N_r N_x$ .

An alternative linear-in- $\zeta$  DAN is to implement an *hypernetwork* which would map  $\mathbf{x}^f$  to a set of weights and biases of a linear neural network that would then be applied to  $\zeta$ . Naively, the number of weights and biases could scale like  $N_x^2$ . However, we can instead map  $\mathbf{x}^f$  to weights and biases of a sequence of convolutional neural networks that we later apply to  $\zeta$ . Yet, we did not test this more sophisticated approach since the former approach is scalable and successful.



### Appendix C: Stein lemma for $\nabla_{\zeta} \mathbf{a}_{\theta}(\mathbf{x}, \zeta)$

Here, we offer a proof of the Stein lemma in the degenerate case where the Gaussian density is singular within the embedding space. This is useful with sparse observations resulting in  $\zeta$  confined within a subspace of  $\mathbf{E}_{\mathbf{x}}$ . This subsumes the full-rank case. Let us assume that  $\rho$  is the pdf defined over  $\mathbf{E}_{\mathbf{x}}$  of the Gaussian random vector  $\zeta$ , whose mean is  $\bar{\zeta}$  and whose covariance matrix is  $\Sigma_{\rho}$  which is assumed positive semi-definite. Hence, the support of  $\zeta$  may be singular in  $\mathbf{E}_{\mathbf{x}}$ . That is why we resort to the singular value decomposition  $\Sigma_{\rho} = \mathbf{U} \Sigma_{\lambda} \mathbf{U}^{\top}$ , where  $\mathbf{U}$  is an orthonormal (though not necessarily orthogonal) matrix such that  $\mathbf{U}^{\top} \mathbf{U} = \mathbf{I}_{\mathbf{x}}$ , and  $\Sigma_{\lambda}$  is a positive definite diagonal matrix of rank lower or equal to  $N_{\mathbf{x}}$ . We can parametrise the random vector  $\zeta$  by

$$\zeta(\omega) = \bar{\zeta} + \mathbf{U} \omega, \quad (\text{C1})$$

where the random vector  $\omega$  has  $\lambda(\omega) = n(\mathbf{0}, \Sigma_{\lambda})$  for Gaussian pdf. Then, denoting  $f(\zeta) \triangleq \mathbf{a}_{\theta}(\mathbf{x}, \zeta)$  for brevity, we have

$$\text{Cov}_{\zeta \sim \rho} [\zeta, f(\zeta)] = \text{Cov}_{\omega \sim \lambda} [\zeta(\omega), f(\zeta(\omega))] \quad (\text{C2a})$$

$$= \int d\omega \lambda(\omega) (\zeta(\omega) - \bar{\zeta}) \otimes (f(\zeta(\omega)) - f(\bar{\zeta})) \quad (\text{C2b})$$

$$= \int d\omega \lambda(\omega) (\mathbf{U} \omega) \otimes f(\zeta(\omega)) \quad (\text{C2c})$$

$$= \mathbf{U} \Sigma_{\lambda} \int d\omega \lambda(\omega) (\Sigma_{\lambda}^{-1} \omega) \otimes f(\zeta(\omega)) \quad (\text{C2d})$$

$$= -\mathbf{U} \Sigma_{\lambda} \int d\omega \lambda(\omega) \nabla_{\omega} \ln \lambda(\omega) \otimes f(\zeta(\omega)) \quad (\text{C2e})$$

$$= -\mathbf{U} \Sigma_{\lambda} \int d\omega \nabla_{\omega} \lambda(\omega) \otimes f(\zeta(\omega)) \quad (\text{C2f})$$

$$= \mathbf{U} \Sigma_{\lambda} \int d\omega \lambda(\omega) \nabla_{\omega} (f(\zeta(\omega))) \quad (\text{C2g})$$

$$= \mathbf{U} \Sigma_{\lambda} \mathbf{U}^{\top} \int d\omega \lambda(\omega) (\nabla_{\zeta} f)(\zeta(\omega)) \quad (\text{C2h})$$

$$= \Sigma_{\rho} \mathbb{E}_{\omega \sim \lambda} [(\nabla_{\zeta} f)(\zeta(\omega))] \quad (\text{C2i})$$

$$= \Sigma_{\rho} \mathbb{E}_{\zeta \sim \rho} [\nabla_{\zeta} f]. \quad (\text{C2j})$$

Hence, we conclude:

$$\mathbb{E}_{\zeta \sim \rho} [\nabla_{\zeta} f] = \Sigma_{\rho}^{\dagger} \text{Cov}_{\zeta \sim \rho} [\zeta, f(\zeta)], \quad (\text{C3})$$

where  $(\cdot)^{\dagger}$  is the Moore-Penrose inverse operator, which comes with the regularisation choice to taper  $\mathbb{E}_{\zeta \sim \rho} [\nabla_{\zeta} f]$  outside of the range of  $\zeta$ . Applied to  $f(\zeta) = \mathbf{a}_{\theta}(\mathbf{x}, \zeta)$ , this yields:

$$\mathbb{E}_{\zeta \sim \rho} [\nabla_{\zeta} \mathbf{a}_{\theta}(\mathbf{x}, \zeta)] = \Sigma_{\rho}^{\dagger} \text{Cov}_{\zeta \sim \rho} [\zeta, \mathbf{a}_{\theta}(\mathbf{x}, \zeta)], \quad (\text{C4a})$$

which, if  $\Sigma_{\rho}$  is full rank, can be written (usual Stein lemma)

$$\mathbb{E}_{\zeta \sim \rho} [\nabla_{\zeta} \mathbf{a}_{\theta}(\mathbf{x}, \zeta)] = \Sigma_{\rho}^{-1} \text{Cov}_{\zeta \sim \rho} [\zeta, \mathbf{a}_{\theta}(\mathbf{x}, \zeta)]. \quad (\text{C4b})$$



## Appendix D: Comparing the empirical and theoretical mean marginal gains

In Sect. 3.2, we showed that the theoretical mean marginal gain  $\bar{\Gamma}$  is expected to be a good approximation of the less simple  
 580 theoretical mean

$$\tilde{\Gamma} = \int d\mathbf{x} d\zeta \pi(\mathbf{x}) \rho(\zeta) \nabla_{\mathbf{x}} \nabla_{\zeta} \mathbf{a}_{\theta}(\mathbf{x}, \zeta), \quad (\text{D1})$$

The latter can now be related to the empirical mean of  $\Gamma$  over a long DA run, denoted  $\hat{\Gamma}$ , and which is obtained from numerical experiments. This is meant to ensure that a pattern emerging from our approximation of  $\bar{\Gamma}$ , is nonetheless consistent with those learned through  $\mathbf{a}_{\theta}$  over the training DA dataset. Hence, given a long trajectory of true states and projected innovations  
 585  $\mathcal{T}_{\mathbf{x}, \zeta} = \{(\mathbf{x}_k^t, \zeta_k)\}_{k=1, \dots, K} \subset (\mathbf{E}_{\mathbf{x}} \otimes \mathbf{E}_{\zeta})^K$ , the empirical sensitivity associated to  $\tilde{\Gamma}$ , and hence  $\bar{\Gamma}$ , should be

$$\hat{\Gamma} \triangleq \langle \Gamma(\mathbf{x}^t, \zeta) \rangle_{(\mathbf{x}^t, \zeta) \in \mathcal{T}_{\mathbf{x}, \zeta}} \xrightarrow{K \rightarrow \infty} \int d\mathbf{x}^t d\zeta p(\mathbf{x}^t, \zeta) \Gamma(\mathbf{x}^t, \zeta), \quad (\text{D2})$$

where  $p(\mathbf{x}^t, \zeta)$  is the joint distribution of  $\mathbf{x}^t$  and  $\zeta$ . However,  $\zeta_k$  not only depends on  $\mathbf{x}_k^t$  but also on the forecast  $\mathbf{x}_k^f$ , so that:

$$\hat{\Gamma} = \int d\mathbf{x}^t d\mathbf{x}^f d\zeta p(\mathbf{x}^t, \mathbf{x}^f, \zeta) \Gamma(\mathbf{x}^t, \zeta), \quad (\text{D3})$$

or, introducing the forecast error  $\mathbf{e}^f \triangleq \mathbf{x}^f - \mathbf{x}^t$ ,

$$590 \quad \hat{\Gamma} = \int d\mathbf{x}^t d\mathbf{e}^f d\zeta p(\mathbf{x}^t, \mathbf{e}^f, \zeta) \Gamma(\mathbf{x}^t + \mathbf{e}^f, \zeta). \quad (\text{D4})$$

By marginalising over  $\mathbf{x}^t$ , we have

$$\hat{\Gamma} = \int d\mathbf{x}^t \pi(\mathbf{x}^t) \int d\mathbf{e}^f d\zeta p(\mathbf{e}^f, \zeta | \mathbf{x}^t) \Gamma(\mathbf{x}^t + \mathbf{e}^f, \zeta) \quad (\text{D5a})$$

$$= \int d\mathbf{x}^t \pi(\mathbf{x}^t) \int d\mathbf{e}^f d\zeta p(\zeta | \mathbf{e}^f) p(\mathbf{e}^f | \mathbf{x}^t) \Gamma(\mathbf{x}^t + \mathbf{e}^f, \zeta) \quad (\text{D5b})$$

$$\approx \int d\mathbf{x}^t \pi(\mathbf{x}^t) \int d\zeta [p(\zeta | \mathbf{e}^f) \Gamma(\mathbf{x}^t + \mathbf{e}^f, \zeta)]_{\mathbf{e}^f=0} \quad (\text{D5c})$$

$$595 \quad = \int d\mathbf{x}^t \pi(\mathbf{x}^t) \int d\zeta \rho(\zeta) \Gamma(\mathbf{x}^t, \zeta) = \tilde{\Gamma}. \quad (\text{D5d})$$

From Eq. (D5a) to Eq. (D5b), we used  $p(\mathbf{e}^f, \zeta | \mathbf{x}^t) = p(\zeta | \mathbf{e}^f) p(\mathbf{e}^f | \mathbf{x}^t)$  since the full dependence of  $\zeta$  on  $\mathbf{x}^t$  is in  $\mathbf{e}^f$ . From Eq. (D5b) to Eq. (D5c), we assumed that the forecast error norm is small compared to the norm of  $\mathbf{x}^t$ , which should indeed be the case in the weak assimilation regime. From Eq. (D5c) to Eq. (D5d), we identified  $p(\zeta | \mathbf{e}^f = 0)$  to  $\rho(\zeta)$ . This points to the (reasonable) approximations made when identifying  $\tilde{\Gamma}$ , Eq. (D1), with the empirical marginal gain  $\hat{\Gamma}$ ,  $\langle \Gamma(\mathbf{x}, \zeta) \rangle_{(\mathbf{x}, \zeta) \in \mathcal{T}}$ ,  
 600 emerging from a long DA run.

Moreover, a formal expression for  $\rho(\zeta | \mathbf{e}^f)$  which appeared in the previous derivation is

$$p(\zeta | \mathbf{e}^f) = \int d\boldsymbol{\varepsilon} d\mathbf{H} d\mathbf{R} p(\boldsymbol{\varepsilon}, \mathbf{H}, \mathbf{R} | \mathbf{e}^f) \delta(\zeta - \mathbf{H}^T \mathbf{R}^{-1} (\boldsymbol{\varepsilon} - \mathbf{H} \mathbf{e}^f)), \quad (\text{D6})$$



where  $p(\varepsilon, \mathbf{H}, \mathbf{R} | \mathbf{e}^f)$  is the pdf of the joint distribution for the observation error  $\varepsilon$ , the observation operator  $\mathbf{H}$ , and the observation error covariance matrix  $\mathbf{R}$ , given the forecast error  $\mathbf{e}^f$ . Marginalising over  $\varepsilon^f$ , we obtain

$$\rho(\zeta) = \int d\mathbf{e}^f p(\zeta | \mathbf{e}^f) p(\mathbf{e}^f) = \int d\mathbf{e}^f d\varepsilon d\mathbf{H} d\mathbf{R} p(\mathbf{e}^f, \varepsilon, \mathbf{H}, \mathbf{R}) \delta(\zeta - \mathbf{H}^T \mathbf{R}^{-1} (\varepsilon - \mathbf{H} \mathbf{e}^f)). \quad (\text{D7})$$

This expression is helpful to formally investigate the symmetries of the distribution of  $\zeta$ , for which Eq. (17) applies.

## Appendix E: Proof of the equivariance of the marginal gain

We wish to prove the equivariance of the marginal gain, i.e. that for all  $\mathbf{x}$  and  $\zeta$ ,  $\Gamma(g \circ \mathbf{x}, g \circ \zeta) = g \otimes g^T \otimes g^T \circ \Gamma(\mathbf{x}, \zeta)$ , under the action of an isometry  $g \in \mathcal{G}$ . The action of  $g$  on either state vector  $\mathbf{x}$  or  $\zeta$  is represented here by the orthogonal matrix  $\mathbf{G}$ :

$$g \circ \mathbf{x} = \mathbf{G}\mathbf{x}, \quad g \circ \zeta = \mathbf{G}\zeta. \quad \text{The action of its associated induced isometry } g_y \in \mathcal{G}_y \text{ is represented by the orthogonal matrix } \mathbf{G}_y: \\ g_y \circ \mathbf{y} = \mathbf{G}_y \mathbf{y}. \text{ Because } \mathbf{G} \text{ and } \mathbf{G}_y \text{ are orthogonal, one has } \mathbf{G}^{-1} = \mathbf{G}^T \text{ and } \mathbf{G}_y^{-1} = \mathbf{G}_y^T.$$

We further assume that (i) the autonomous dynamics  $\mathcal{M}$  commute with the elements of  $\mathcal{G}$ :  $\mathbf{G}\mathcal{M}(\mathbf{G}^T(\cdot)) = \mathcal{M}(\cdot)$ , which stands for, e.g., the L96 model and the group of discrete translations, (ii) the observation operator, and the observation errors are subject to  $\mathbf{G}_y \mathcal{H}_k(\mathbf{G}^T(\cdot)) = \mathcal{H}_k(\cdot)$ ,  $\mathbf{G}_y \mathbf{H}_k \mathbf{G}^T = \mathbf{H}_k$ , and  $\mathbf{G}_y \mathbf{R} \mathbf{G}_y^T = \mathbf{R}$ . To prove the equivariance, we first consider the

$$\text{optimisation problem that defines } (\mathbf{x}, \zeta) \mapsto \mathbf{a}_\theta(g \circ \mathbf{x}, g \circ \zeta):$$

$$\mathcal{L}(\theta | \{\mathbf{x}_k^t, \mathbf{y}_k\}) = \sum_{k=1}^K \|\mathbf{x}_k^t - \mathbf{x}_k^a(\theta)\|^2, \quad (\text{E1a})$$

$$\mathbf{x}_k^a = \mathbf{x}_k^f + \mathbf{a}_\theta(\mathbf{G}\mathbf{x}_k^f, \mathbf{G}\zeta_k), \quad (\text{E1b})$$

$$\zeta_k = \mathbf{H}_k^T \mathbf{R}_k^{-1} (\mathbf{y}_k - \mathcal{H}_k(\mathbf{x}_k^f)), \quad (\text{E1c})$$

$$\mathbf{x}_{k+1}^f = \mathcal{M}(\mathbf{x}_k^a), \quad (\text{E1d})$$

which is equivalent to

$$\mathcal{L}(\theta | \{\mathbf{x}_k^t, \mathbf{y}_k\}) = \sum_{k=1}^K \|\mathbf{G}\mathbf{x}_k^t - \mathbf{G}\mathbf{x}_k^a(\theta)\|^2, \quad (\text{E2a})$$

$$\mathbf{G}\mathbf{x}_k^a = \mathbf{G}\mathbf{x}_k^f + \mathbf{G}\mathbf{a}_\theta(\mathbf{G}\mathbf{x}_k^f, \mathbf{G}\zeta_k), \quad (\text{E2b})$$

$$\mathbf{G}\zeta_k = \mathbf{G}\mathbf{H}_k^T \mathbf{G}_y^T \mathbf{G}_y \mathbf{R}_k^{-1} \mathbf{G}_y^T \mathbf{G}_y (\mathbf{y}_k - \mathcal{H}_k(\mathbf{G}^T \mathbf{G}\mathbf{x}_k^f)), \quad (\text{E2c})$$

$$\mathbf{G}\mathbf{x}_{k+1}^f = \mathbf{G}\mathcal{M}(\mathbf{G}^T \mathbf{G}\mathbf{x}_k^a), \quad (\text{E2d})$$

where Eq. (E2a) is obtained from Eq. (E1a) because  $\mathbf{G}$  is orthogonal, and Eqs. (E2b, E2c, E2d) are obtained from a multiplication on the left by  $\mathbf{G}$  and insertion of  $\mathbf{G}^T \mathbf{G} = \mathbf{I}_x$  and  $\mathbf{G}_y^T \mathbf{G}_y = \mathbf{I}_y$  in Eqs. (E1b, E1c, E1d). Hence, denoting  $\tilde{\mathbf{x}}_k^a = \mathbf{G}\mathbf{x}_k^a$ ,



$\tilde{\mathbf{x}}_k^f = \mathbf{G}\mathbf{x}_k^f$ ,  $\tilde{\boldsymbol{\zeta}}_k = \mathbf{G}\boldsymbol{\zeta}_k$ , the problem is reformulated as

$$\mathcal{L}(\boldsymbol{\theta} | \{\mathbf{x}_k^t, \mathbf{y}_k\}) = \sum_{k=1}^K \|\mathbf{G}\mathbf{x}_k^t - \tilde{\mathbf{x}}_k^a(\boldsymbol{\theta})\|^2, \quad (\text{E3a})$$

$$\tilde{\mathbf{x}}_k^a = \tilde{\mathbf{x}}_k^f + \mathbf{G}\mathbf{a}_{\boldsymbol{\theta}}(\tilde{\mathbf{x}}_k^f, \tilde{\boldsymbol{\zeta}}_k), \quad (\text{E3b})$$

$$\tilde{\boldsymbol{\zeta}}_k = (\mathbf{G}_y \mathbf{H}_k \mathbf{G}^T)^T (\mathbf{G}_y \mathbf{R}_k^{-1} \mathbf{G}_y^T) (\mathbf{G}_y \mathbf{y}_k - \mathbf{G}_y \mathcal{H}_k(\mathbf{G}^T \tilde{\mathbf{x}}_k^f)), \quad (\text{E3c})$$

$$\tilde{\mathbf{x}}_{k+1}^f = \mathbf{G}\mathcal{M}(\mathbf{G}^{-1} \tilde{\mathbf{x}}_k^a). \quad (\text{E3d})$$

Leveraging the assumptions on  $\mathcal{G}$ , we finally obtain

$$\mathcal{L}(\boldsymbol{\theta} | \{\mathbf{x}_k^t, \mathbf{y}_k\}) = \sum_{k=1}^K \|\mathbf{G}\mathbf{x}_k^t - \tilde{\mathbf{x}}_k^a(\boldsymbol{\theta})\|^2, \quad (\text{E4a})$$

$$\tilde{\mathbf{x}}_k^a = \tilde{\mathbf{x}}_k^f + \mathbf{G}\mathbf{a}_{\boldsymbol{\theta}}(\tilde{\mathbf{x}}_k^f, \tilde{\boldsymbol{\zeta}}_k), \quad (\text{E4b})$$

$$\tilde{\boldsymbol{\zeta}}_k = \mathbf{H}_k^T \mathbf{R}_k^{-1} (\mathbf{G}_y \mathbf{y}_k - \mathcal{H}_k(\tilde{\mathbf{x}}_k^f)), \quad (\text{E4c})$$

$$\tilde{\mathbf{x}}_{k+1}^f = \mathcal{M}(\tilde{\mathbf{x}}_k^a). \quad (\text{E4d})$$

This shows that  $\mathbf{a}_{\boldsymbol{\theta}}(\mathbf{G}\cdot, \mathbf{G}\cdot)$  is the solution of Eq. (E1) whose input is the dataset  $\{\mathbf{x}_k^t, \mathbf{y}_k\}$ , while  $\mathbf{G}\mathbf{a}_{\boldsymbol{\theta}}(\cdot, \cdot)$  is the solution of Eq. (E4) whose input is the dataset  $\{\mathbf{G}\mathbf{x}_k^t, \mathbf{G}_y \mathbf{y}_k = \mathbf{G}_y \mathcal{H}_k(\mathbf{x}_k^t) + \mathbf{G}_y \varepsilon_k = \mathcal{H}_k(\mathbf{G}\mathbf{x}_k^t) + \mathbf{G}_y \varepsilon_k\}$ . Since both the invariant distribution of the dynamics and the distribution of the observations errors are invariant under  $\mathcal{G}$ , the datasets  $\{\mathbf{x}_k^t, \mathbf{y}_k\}$  and  $\{\mathbf{G}\mathbf{x}_k^t, \mathbf{G}_y \mathbf{y}_k\}$  must asymptotically yield the same solution for a large enough number of samples  $K$ . This proves the equiv-  
ariance

$$\forall g \in \mathcal{G}, \forall \mathbf{x} \in \mathbf{E}_x, \forall \boldsymbol{\zeta} \in \mathbf{E}_x : \quad g \odot \mathbf{a}_{\boldsymbol{\theta}}(g \circ \mathbf{x}, g \circ \boldsymbol{\zeta}) = g \circ \boldsymbol{\Gamma}(\mathbf{x}, \boldsymbol{\zeta}). \quad (\text{E5})$$

Then, taking the gradient with respect to  $\mathbf{x}$  and  $\boldsymbol{\zeta}$  of  $\mathbf{a}_{\boldsymbol{\theta}}(g \circ \mathbf{x}, g \circ \boldsymbol{\zeta})$  yields a contravariant action  $g^T \otimes g^T$  onto the tensor factors for  $\mathbf{x}$  and  $\boldsymbol{\zeta}$ , which proves the equivariance of  $\boldsymbol{\Gamma}$ , Eq. (16). The equivariance assumption on the observation operator is rather  
stringent. A weaker assumption is to assume that  $\{\mathbf{G}_y \mathcal{H}_k(\mathbf{G}^T \mathbf{x}_k^f)\}_{k=1, \dots, K}$  almost spans the same set as  $\{\mathcal{H}_k(\mathbf{x}_k^f)\}_{k=1, \dots, K}$  for  $K \rightarrow \infty$ . Then the optimisation problems Eq. (E1) and Eq. (E4) should almost coincide. This is for instance useful when one considers random observation operators for which operator instances have no specific symmetry, while their distribution do exhibit the symmetry, a case occurring in Sect. 3.2.4.

## Appendix F: Sleek representation of the mean marginal gain tensor

Assume that the states and projected innovations are defined as fields over a physical manifold  $\mathcal{D}$ , and further discretised at  $N_x$  collocation space points of  $\mathcal{D}$  indexed by  $i \in \llbracket 1, \dots, N_x \rrbracket$ . Hence,  $\mathcal{G}$  is a discrete group of isometries. As a consequence,  $g \in \mathcal{G}$  can be seen as a bijection of  $\llbracket 1, \dots, N_x \rrbracket$  and the action of  $g \in \mathcal{G}$  on the fields  $\mathbf{x}$  and  $\boldsymbol{\zeta}$  reads

$$[g \circ \mathbf{x}]_i = [\mathbf{x}]_{g(i)}, \quad [g \circ \boldsymbol{\zeta}]_j = [\boldsymbol{\zeta}]_{g(j)}, \quad (\text{F1})$$



respectively. Likewise, the action of  $g \in \mathcal{G}$  on  $\bar{\Gamma}$  is, for all  $i, j, l$ :

$$655 \quad [g \odot \bar{\Gamma}]_{ijl} = [g \otimes g^T \otimes g^T \circ \bar{\Gamma}]_{ijl} = [\bar{\Gamma}]_{g(i)g^T(j)g^T(l)}, \quad (\text{F2})$$

so that Eq. (20) reads, for all  $i, j, l$ :

$$[\bar{\Gamma}]_{ijl} = [\bar{\Gamma}]_{g(i)g^T(j)g^T(l)}. \quad (\text{F3})$$

Let us choose one of the collocation point in  $\mathcal{D}$  with index  $r \in \llbracket 1, \dots, N_x \rrbracket$ . With the above assumptions, the orbit of the site indexed by  $r$  under the action of  $\mathcal{G}$  is  $\llbracket 1, \dots, N_x \rrbracket$ . We can then define a 2-tensor  $\bar{\Omega}_r$  by, for all  $i, j$ :

$$660 \quad [\bar{\Omega}_r]_{ij} = [\bar{\Gamma}]_{ijr}. \quad (\text{F4})$$

For all  $l \in \llbracket 1, \dots, N_x \rrbracket$ , we can pick at least one  $g_r^l \in \mathcal{G}$  such that  $g_r^l(r) = l$  and let us denote its inverse by  $g_l^r$  which coincides with it adjoint  $g^T$  and satisfies  $g_l^r(l) = r$  in particular. Hence, we have from Eq. (20) and for all  $i, j, l$ :

$$[\bar{\Gamma}]_{ijl} = [\bar{\Gamma}]_{g_r^l(i)g_l^r(j)g_l^r(l)} = [\bar{\Gamma}]_{g_r^l(i)g_l^r(j)r} = [\bar{\Omega}_r]_{g_r^l(i)g_l^r(j)}. \quad (\text{F5})$$

As a consequence of the symmetry and Eq. (F5), the 3-tensor  $\bar{\Gamma}$  can be entirely specified by the 2-tensor  $\bar{\Omega}_r$ , where  $r$ , a reference site index, is arbitrarily chosen. In the case where  $\mathcal{D}$  is one-dimensional (as for L96),  $\bar{\Omega}_r$  is a matrix, hence depictable and more easily interpretable.

## Appendix G: Numerical computation of the mean marginal gain

The mean marginal gain can be computed from states of a trajectory  $\mathcal{T}_x$  of the ergodic dynamics, and the ability to evaluate  $x \mapsto \Gamma(x)$  as defined by Eq. (9b). The trajectory should be long enough so that its states adequately sample the invariant distribution  $\pi$ . From Eq. (11), we hence have the empirical estimator:

$$670 \quad \bar{\Gamma} = \langle \Gamma(x) \rangle_{x \in \mathcal{T}_x} = \frac{1}{K} \sum_{k=1}^K \Gamma(x_k). \quad (\text{G1})$$

As a result, the computational complexity of the mean marginal gain is proportional to  $K$ , but may be significantly alleviated by the presence of symmetries as discussed before. Such symmetries must make both  $\pi$  and  $\rho$  invariant even though the definition of  $\bar{\Gamma}$  only implicitly depends on  $\rho$ .

675 We now turn to the estimation of the marginal gain  $\Gamma(x)$ . Its computation can be achieved through several routes with distinct numerical complexities which, as approximations, may not be equivalent and may lead to mildly differing results.

The first way to compute  $\Gamma(x)$  is through automatic differentiation. As a second-order sensitivity of  $\mathbf{a}_\theta$  with respect to  $x$  and  $\zeta$ ,  $\Gamma(x) = \nabla_x \nabla_\zeta \mathbf{a}_\theta(x, \zeta)|_{\zeta=0}$  requires taking the Jacobian of  $\mathbf{a}_\theta$  twice. Hence, such computation through either JAX, Pytorch or Tensorflow, can be prohibitive, with a substantial need for GPU memory. On an NVIDIA RTX5000 Ada GPU with 32 Go

680 of memory, using the JAX-inspired Pytorch `torch.func` module,<sup>1</sup> we found it to be achievable with the L96 model, difficult

<sup>1</sup><https://pytorch.org/docs/stable/func.html>



with a Kuramoto-Sivashinsky model (Kuramoto and Tsuzuki, 1976; Sivashinsky, 1977), but prohibitive with a single-layer QG model on the sphere. Hence, it is likely to be impractical with high-dimensional models.

Note that the mean Eq. (G1) can be computed through updates whenever a new  $\Gamma(\mathbf{x}_k)$  is computed, preventing the need to store them. Moreover, when exploiting symmetries of  $\mathcal{G}$ , the intermediate tensor

$$685 \quad [\Omega_r]_{ij}(\mathbf{x}_k) = \frac{1}{N_x} \sum_{l=1}^{N_x} [\Gamma]_{g_i^r(i) g_r^l(j) l}(\mathbf{x}_k), \quad (\text{G2a})$$

can be computed, and will contribute to the computation of  $\bar{\Omega}_r$  through the update of

$$\bar{\Omega}_r = \langle \Omega_r(\mathbf{x}) \rangle_{\mathbf{x} \in \mathcal{T}_x} = \frac{1}{K} \sum_{k=1}^K \Omega_r(\mathbf{x}_k). \quad (\text{G2b})$$

In Boc24, either the gain  $\mathbf{K}(\mathbf{x})$  or  $\mathbf{P}^a(\mathbf{x})$  were obtained by generating an ensemble of perturbations to feed a regression. The same idea can be used for  $\Gamma(\mathbf{x})$ , once again assuming a quasi-linear behaviour of  $\mathbf{a}_\theta(\mathbf{x}, \zeta)$  and  $\nabla_{\mathbf{x}} \mathbf{a}_\theta(\mathbf{x}, \zeta)$  as functions  
690 of  $\zeta$ . From the results in Sect. 3.2.1 and Appendix D, we infer that

$$\Gamma(\mathbf{x}) = \nabla_{\mathbf{x}} \nabla_{\zeta} \mathbf{a}_\theta(\mathbf{x}, \zeta)|_{\zeta=0} \approx \mathbb{E}_{\zeta \sim \rho} [\nabla_{\mathbf{x}} \nabla_{\zeta} \mathbf{a}_\theta(\mathbf{x}, \zeta)] \approx \Sigma_{\rho}^{\dagger} \text{Cov}_{\zeta \sim \rho} [\zeta, \nabla_{\mathbf{x}} \mathbf{a}_\theta(\mathbf{x}, \zeta)], \quad (\text{G3})$$

which tells that a sampling approach can be applied to  $\nabla_{\mathbf{x}} \mathbf{a}_\theta(\mathbf{x}, \cdot)$ . Assuming  $\mathbf{a}_\theta$  is regular enough, the gradients with respect to  $\mathbf{x}$  and  $\zeta$  commute and we also have  $\Gamma(\mathbf{x}) \approx \nabla_{\mathbf{x}} \mathbb{E}_{\zeta \sim \rho} [\nabla_{\zeta} \mathbf{a}_\theta(\mathbf{x}, \zeta)]$ , although this could considerably complexify backpropagation if automatic differentiation is used to handle  $\nabla_{\mathbf{x}}$ .

695 Automatic differentiation is hence used only once for the computation of the Jacobian  $\nabla_{\mathbf{x}} \mathbf{a}_\theta(\mathbf{x}, \cdot)$ , as opposed to the full automatic differentiation approach. Hence,  $\Gamma(\mathbf{x})$  can be computed using a *composite* Monte Carlo/differentiation approach. The details of the subsequent regression are reported in Appendix H.

## Appendix H: Regression for the composite mean marginal gain

An ensemble of  $N_p$  perturbations  $\partial \mathbf{a}_p$  generated from  $N_p$  samples  $\zeta_p \sim N(\mathbf{0}, \Xi)$  for  $p = 1, \dots, N_p$  should first be computed,

$$700 \quad [\partial \mathbf{a}_p]_{il} = [\partial_{x_l} \mathbf{a}_\theta(\mathbf{x}, \zeta_p)]_i, \quad (\text{H1})$$

via an ensemble of first-order Jacobians. In the best linear unbiased estimator framework, the covariance matrix  $\Xi$  should roughly match  $\mathbf{H}^T \mathbf{R}^{-1} (\mathbf{R} + \mathbf{H} \mathbf{P}^f \mathbf{H}^T) \mathbf{R}^{-1} \mathbf{H}$  where  $\mathbf{P}^f$  and  $\mathbf{R} + \mathbf{H} \mathbf{P}^f \mathbf{H}^T$  are the forecast error and innovation covariance matrices, respectively. Hence, in the weak assimilation regime, we can use the approximation  $\Xi \approx \mathbf{H}^T \mathbf{R}^{-1} \mathbf{H}$ , which should be regarded as a scale for the perturbations anyway, and generate samples with  $\zeta = \mathbf{H}^T \mathbf{R}^{-\frac{1}{2}} \xi$ , where  $\xi \sim \mathcal{N}(\mathbf{0}, \mathbf{I}_y)$ . Introducing

705 for  $p = 1, \dots, N_p$ , the recentred samples

$$\partial \mathbf{a}'_p = \mathbf{a}_p - N_p^{-1} \sum_{p=1}^{N_p} \partial \mathbf{a}_p, \quad \zeta'_p = \zeta_p - N_p^{-1} \sum_{p=1}^{N_p} \zeta_p, \quad (\text{H2})$$





we have from the definitions  $\mathbf{C} \triangleq \text{Cov}_{\zeta \sim \rho} [\zeta, \nabla_{\mathbf{x}} \mathbf{a}_{\theta}(\mathbf{x}, \zeta)]$  and  $\mathbf{D} \triangleq \Sigma_{\rho}$ :

$$[\mathbf{C}]_{ijl} \approx \sum_{p=1}^{N_p} [\zeta'_p]_j [\partial \mathbf{a}'_p]_{il}, \quad [\mathbf{D}]_{jl} \approx \sum_{p=1}^{N_p} [\zeta'_p]_j [\zeta'_p]_l, \quad (\text{H3})$$

which, from Eq. (G3), yields

$$[\mathbf{\Gamma}]_{ijl} = \sum_{m=1}^{N_x} [\mathbf{C}]_{iml} [\mathbf{D}^{-1}]_{mj}. \quad (\text{H4})$$

Obviously, in a high-dimensional context, the approach would necessitate reduction methods such as Lanczos vectors or (randomised) singular value decompositions, and the generation of the ensemble would require a massive vectorisation on GPUs.  $\mathbf{\Omega}_r$  can then be computed by averaging  $\mathbf{\Gamma}$  using, e.g., Eqs. (G2). Moreover, it is not difficult to show that  $\mathbf{\Omega}_r$ , as defined by Eq. (G2a), can alternatively be obtained by first averaging over  $\mathbf{C}$  and  $\mathbf{D}$  before performing the inversion of the regression, that is:

$$\mathbf{\Omega}_r = \overline{\mathbf{C}} \overline{\mathbf{D}}^{-1}, \quad [\mathbf{C}]_{ijl} = \frac{1}{N_x} \sum_{s=1}^{N_x} [\mathbf{C}]_{g_s^l(i) g_s^s(j) s}, \quad [\mathbf{D}]_{ij} = \frac{1}{N_x} \sum_{s=1}^{N_x} [\mathbf{D}]_{g_s^j(i) s}. \quad (\text{H5})$$

Either way, the composite approach may turn out numerically cheaper than the full differentiation approach.

*Author contributions.* MB performed the mathematical, algorithmic, and numerical analysis. MB and TF discussed of the implications of the results. All the authors worked on the structure of the manuscript. All the authors reviewed and edited the manuscript.

*Competing interests.* The contact author has declared that none of the authors has any competing interests

*Acknowledgements.* This project—Learning efficient Data Assimilation from Artificial Intelligence (DAbyAI)—has been supported by NVIDIA and their Academic Grant Program through the grant of two RTX 6000 Ada GPUs that were intensively used in the numerical experiments of this work. This paper is also a contribution to the DRUIDS project, supported by France 2030 PEPR Maths-Vives, grant ANR-24-EXMA-0002. Tobias S. Finn acknowledges the support of the project SASIP (grant no. G-24-66154) funded by Schmidt Sciences—a philanthropic initiative that seeks to improve societal outcomes through the development of emerging science and technologies, and the support of France 2030 PEPR Maths-Vives, grant ANR-24-EXMA-0001, project Climaths/GenClimEx. Sibio Cheng acknowledges the support of the French Agence Nationale de la Recherche (ANR) under reference ANR-22-CPJ2-0143-01. CEREIA is a member of Institut Pierre-Simon Laplace (IPSL). AI was used to polish the English of a limited number of paragraphs of the manuscript.



## References

- 730 Agarwal, S., Jabbari, S., Agarwal, C., Upadhyay, S., Wu, S., and Lakkaraju, H.: Towards the Unification and Robustness of Perturbation and Gradient Based Explanations, in: *Proceedings of the 38th International Conference on Machine Learning*, edited by Meila, M. and Zhang, T., vol. 139 of *Proceedings of Machine Learning Research*, pp. 110–119, PMLR, <https://proceedings.mlr.press/v139/agarwal21c.html>, 2021.
- Alexe, M., Boucher, E., Lean, P., Pinnington, E., Laloyaux, P., McNally, A., Lang, S., Chantry, M., Burrows, C., Chrust, M., Pinault, F., Villeneuve, E., Bormann, N., and Healy, S.: GraphDOP: Towards skilful data-driven medium-range weather forecasts learnt and initialised directly from observations, <https://doi.org/10.48550/arXiv.2412.15687>, 2024.
- 735 Allen, A., Markou, S., Tebbutt, W., Requeima, J., Bruinsma, W. P., Andersson, T. R., Herzog, M., Lane, N. D., Chantry, M., Hosking, J. S., and Turner, R. E.: End-to-end data-driven weather prediction, *Nature*, <https://doi.org/10.1038/s41586-025-08897-0>, 2025.
- Arnold, L.: *Random Dynamical Systems*, Springer Berlin, Heidelberg, ISBN 978-3-540-63758-5, [https://doi.org/10.1007/978-3-662-12878-](https://doi.org/10.1007/978-3-662-12878-7)
- 740 7, 1998.
- Asch, M., Bocquet, M., and Nodet, M.: *Data Assimilation: Methods, Algorithms, and Applications*, Fundamentals of Algorithms, SIAM, Philadelphia, ISBN 978-1-611974-53-9, <https://doi.org/10.1137/1.9781611974546>, 2016.
- Bannister, R. N.: A review of operational methods of variational and ensemble-variational data assimilation, *Q. J. R. Meteorol. Soc.*, 143, 607–633, <https://doi.org/10.1002/qj.2982>, 2017.
- 745 Barone, A., Carrassi, A., Savary, T., Demayer, J., and Vannitsem, S.: Structural origins and real-time predictors of intermittency, *Chaos*, 35, 103 119, <https://doi.org/10.1063/5.0287572>, 2025.
- Bocquet, M.: Ensemble Kalman filtering without the intrinsic need for inflation, *Nonlin. Processes Geophys.*, 18, 735–750, <https://doi.org/10.5194/npg-18-735-2011>, 2011.
- Bocquet, M. and Carrassi, A.: Four-dimensional ensemble variational data assimilation and the unstable subspace, *Tellus A*, 69, 1304 504, <https://doi.org/10.1080/16000870.2017.1304504>, 2017.
- 750 Bocquet, M. and Farchi, A.: On the consistency of the perturbation update of local ensemble square root Kalman filters, *Tellus A*, 71, 1–21, <https://doi.org/10.1080/16000870.2019.1613142>, 2019.
- Bocquet, M. and Sakov, P.: Combining inflation-free and iterative ensemble Kalman filters for strongly nonlinear systems, *Nonlin. Processes Geophys.*, 19, 383–399, <https://doi.org/10.5194/npg-19-383-2012>, 2012.
- 755 Bocquet, M. and Sakov, P.: Joint state and parameter estimation with an iterative ensemble Kalman smoother, *Nonlin. Processes Geophys.*, 20, 803–818, <https://doi.org/10.5194/npg-20-803-2013>, 2013.
- Bocquet, M. and Sakov, P.: An iterative ensemble Kalman smoother, *Q. J. R. Meteorol. Soc.*, 140, 1521–1535, <https://doi.org/10.1002/qj.2236>, 2014.
- Bocquet, M., Raanes, P. N., and Hannart, A.: Expanding the validity of the ensemble Kalman filter without the intrinsic need for inflation, *Nonlin. Processes Geophys.*, 22, 645–662, <https://doi.org/10.5194/npg-22-645-2015>, 2015.
- 760 Bocquet, M., Gurumoorthy, K. S., Apte, A., Carrassi, A., Grudzien, C., and Jones, C. K. R. T.: Degenerate Kalman filter error covariances and their convergence onto the unstable subspace, *SIAM/ASA J. Uncertain. Quantif.*, 5, 304–333, <https://doi.org/10.1137/16M1068712>, 2017.
- Bocquet, M., Farchi, A., Finn, T. S., Durand, C., Cheng, S., Chen, Y., Pasmans, I., and Carrassi, A.: Accurate deep learning-based filtering for chaotic dynamics by identifying instabilities without an ensemble, *Chaos*, 29, 091 104, <https://doi.org/10.1063/5.0230837>, 2024.
- 765



- Boudier, P., Fillion, A., Gratton, S., Gürol, S., and Zhang, S.: Data Assimilation Networks, *J. Adv. Model. Earth Syst.*, 15, e2022MS003 353, <https://doi.org/10.1029/2022MS003353>, 2023.
- Buehner, M., McTaggart-Cowan, R., Beaulne, A., Charette, C., Garand, L., Heillette, S., Lapalme, E., Laroche, S., Macpherson, S. R., Morneau, J., and Zadra, A.: Implementation of Deterministic Weather Forecasting Systems based on Ensemble-Variational Data Assimilation at Environment Canada. Part I: The Global System, *Mon. Wea. Rev.*, 143, 2532–2559, <https://doi.org/10.1175/MWR-D-14-00354.1>, 2015.
- Carrassi, A., Ghil, M., Trevisan, A., and Ubaldi, F.: Data assimilation as a nonlinear dynamical systems problem: Stability and convergence of the prediction-assimilation system, *Chaos*, 18, 023 112, <https://doi.org/10.1063/1.2909862>, 2008.
- Carrassi, A., Bocquet, M., Bertino, L., and Evensen, G.: Data Assimilation in the Geosciences: An overview on methods, issues, and perspectives, *WIREs Climate Change*, 9, e535, <https://doi.org/10.1002/wcc.535>, 2018.
- Carrassi, A., Bocquet, M., Demaeyer, J., Gruzien, C., Raanes, P. N., and Vannitsem, S.: Data assimilation for chaotic dynamics, in: *Data Assimilation for Atmospheric, Oceanic and Hydrologic Applications (Vol. IV)*, edited by P., S. K. and X., L., pp. 1–42, Springer International Publishing, Cham, ISBN 978-3-030-77722-7, [https://doi.org/10.1007/978-3-030-77722-7\\_1](https://doi.org/10.1007/978-3-030-77722-7_1), 2022.
- Chekroun, M. D., Simonnet, E., and Ghil, M.: Stochastic climate dynamics: Random attractors and time-dependent invariant measures, *Physica D*, 240, 1685–1700, <https://doi.org/10.1016/j.physd.2011.06.005>, 2011.
- Cheng, S., Quilodran-Casas, C., Ouala, S., Farchi, A., Liu, C., Tandeo, P., Fablet, R., Lucor, D., Iooss, B., Brajard, J., Xiao, D., Janjic, T., Ding, W., Guo, Y., Carrassi, A., Bocquet, M., and Arcucci, R.: Machine learning with data assimilation and uncertainty quantification for dynamical systems: a review, *IEEE/CAA J. Autom. Sin.*, 10, 1361–1387, <https://doi.org/10.1109/JAS.2023.123537>, 2023.
- Cintra, R. S. and de Campos Velho, H. F.: Data assimilation by artificial neural networks for an atmospheric general circulation model, chap. 17, pp. 265–286, IntechOpen, <https://doi.org/10.5772/intechopen.70791>, 2018.
- Daley, R.: *Atmospheric Data Analysis*, Cambridge University Press, New-York, ISBN 9780521458252, 1991.
- Desroziers, G., Camino, J.-T., and Berre, L.: 4DEnVar: link with 4D state formulation of variational assimilation and different possible implementations, *Q. J. R. Meteorol. Soc.*, 140, 2097–2110, <https://doi.org/10.1002/qj.2325>, 2014.
- Fablet, R., Chapron, B., Drumetz, L., Mémin, E., Pannekoucke, O., and Rousseau, F.: Learning Variational Data Assimilation Models and Solvers, *J. Adv. Model. Earth Syst.*, 13, e2021MS002 572, <https://doi.org/10.1029/2021MS002572>, 2021.
- Fillion, A., Bocquet, M., and Gratton, S.: Quasi static ensemble variational data assimilation: a theoretical and numerical study with the iterative ensemble Kalman smoother, *Nonlin. Processes Geophys.*, 25, 315–334, <https://doi.org/10.5194/npg-25-315-2018>, 2018.
- Filoche, A., Brajard, J., Charantonis, A., and Béréziat, D.: Learning 4DVAR Inversion Directly from Observations, in: *Computational Science – ICCS 2023*, edited by Mikyška, J., de Mulatier, C., Paszynski, M., Krzhizhanovskaya, V. V., Dongarra, J. J., and Sloot, P. M., pp. 414–421, Springer Nature Switzerland, Cham, ISBN 978-3-031-36027-5, [https://doi.org/10.1007/978-3-031-36027-5\\_32](https://doi.org/10.1007/978-3-031-36027-5_32), 2023.
- Flandoli, F. and Tonello, E.: An introduction to random dynamical systems for climate, <https://pagine.dm.unipi.it/flandoli/Part1bis.pdf>, 2021.
- Frerix, T., Kochkov, D., Smith, J., Cremers, D., Brenner, M., and Hoyer, S.: Variational Data Assimilation with a Learned Inverse Observation Operator, in: *Proceedings of the 38th International Conference on Machine Learning*, edited by Meila, M. and Zhang, T., vol. 139 of *Proceedings of Machine Learning Research*, pp. 3449–3458, PMLR, <https://proceedings.mlr.press/v139/frerix21a.html>, 2021.
- Ghil, M. and Sciamarella, D.: Review article: Dynamical systems, algebraic topology and the climate sciences, *Nonlin. Processes Geophys.*, 30, 399–434, <https://doi.org/10.5194/npg-30-399-2023>, 2023.
- Härter, T. P. and de Campos Velho, H. F.: Data Assimilation Procedure by Recurrent Neural Network, *Eng. Appl. Comput. Fluid Mech.*, 6, 224–233, <https://doi.org/10.1080/19942060.2012.11015417>, 2012.



- Jazwinski, A. H.: Stochastic Processes and Filtering Theory, Academic Press, New-York, 1970.
- 805 Kalnay, E.: Atmospheric Modeling, Data Assimilation and Predictability, Cambridge University Press, Cambridge, ISBN 9780521796293, 2003.
- Karimi, A. and Paul, M. R.: Extensive chaos in the Lorenz-96 model, *Chaos*, 20, 043 105, <https://doi.org/10.1063/1.3496397>, 2010.
- Keller, J. D. and Potthast, R.: AI-based data assimilation: Learning the functional of analysis estimation, <https://doi.org/10.48550/arXiv.2406.00390>, 2024.
- 810 Kuramoto, Y. and Tsuzuki, T.: Persistent propagation of concentration waves in dissipative media far from thermal equilibrium, *Progr. Theoret. Phys.*, 55, 356–369, <https://doi.org/10.1143/PTP.55.356>, 1976.
- Lafon, N., Fablet, R., and Naveau, P.: Uncertainty Quantification When Learning Dynamical Models and Solvers With Variational Methods, *J. Adv. Model. Earth Syst.*, 15, e2022MS003 446, <https://doi.org/10.1029/2022MS003446>, 2023.
- Laloyaux, P., Alexe, M., Boucher, E., Lean, P., Pinnington, E., Lang, S., Necker, T., and McNally, A.: Using data assimilation tools to dissect
- 815 GraphDOP, <https://doi.org/10.48550/arXiv.2510.27388>, 2025.
- Lean, P., Alexe, M., Boucher, E., Pinnington, E., Lang, S., Laloyaux, P., Bormann, N., and McNally, A.: Learning from nature: insights into GraphDOP's representations of the Earth System, <https://doi.org/10.48550/arXiv.2508.18018>, 2025.
- Liu, J. S.: Siegel's formula via Stein's identities, *Statistics & Probability Letters*, 21, 247–251, [https://doi.org/https://doi.org/10.1016/0167-7152\(94\)90121-X](https://doi.org/https://doi.org/10.1016/0167-7152(94)90121-X), 1994.
- 820 Lorenz, E. N. and Emanuel, K. A.: Optimal sites for supplementary weather observations: simulation with a small model, *J. Atmos. Sci.*, 55, 399–414, [https://doi.org/10.1175/1520-0469\(1998\)055<0399:OSFSWO>2.0.CO;2](https://doi.org/10.1175/1520-0469(1998)055<0399:OSFSWO>2.0.CO;2), 1998.
- Lu, F.: U-Net Kalman Filter (UNetKF): An Example of Machine Learning-Assisted Data Assimilation, *J. Adv. Model. Earth Syst.*, 17, e2023MS003 979, <https://doi.org/10.1029/2023MS003979>, 2025.
- Maddy, E. S., Boukabara, S. A., and Iturbide-Sanchez, F.: Assessing the Feasibility of an NWP Satellite Data Assimilation System Entirely Based on AI Techniques, *IEEE Journal of Selected Topics in Applied Earth Observations and Remote Sensing*, 17, 9828–9845, <https://doi.org/10.1109/JSTARS.2024.3397078>, 2024.
- 825 McCabe, M. and Brown, J.: Learning to Assimilate in Chaotic Dynamical Systems, in: *Advances in Neural Information Processing Systems*, edited by Ranzato, M., Beygelzimer, A., Dauphin, Y., Liang, P., and Vaughan, J. W., vol. 34, pp. 12 237–12 250, Curran Associates, Inc., [https://proceedings.neurips.cc/paper\\_files/paper/2021/file/65cc2c8205a05d7379fa3a6386f710e1-Paper.pdf](https://proceedings.neurips.cc/paper_files/paper/2021/file/65cc2c8205a05d7379fa3a6386f710e1-Paper.pdf), 2021.
- 830 McNally, T., Lessig, C., Lean, P., Chantry, M., Alexe, M., and Lang, S.: Red sky at night... producing weather forecasts directly from observations, <https://doi.org/10.21957/tmc81jo4c7>, 2024.
- Oseledec, V. I.: A multiplicative ergodic theorem. Ljapunov characteristic numbers for dynamical systems., *Trans. Moscow Math Soc.*, 19, 197–231, 1968.
- Palatella, L., Carrassi, A., and Trevisan, A.: Lyapunov vectors and assimilation in the unstable subspace: theory and applications, *J. Phys. A: Math. Theor.*, 46, 254 020, <https://doi.org/10.1088/1751-8113/46/25/254020>, 2013.
- 835 Pannekoucke, O. and Fablet, R.: PDE-NetGen 1.0: from symbolic partial differential equation (PDE) representations of physical processes to trainable neural network representations, *Geosci. Model Dev.*, 13, 3373–3382, <https://doi.org/10.5194/gmd-13-3373-2020>, 2020.
- Pannekoucke, O., Ricci, S., Barthelemy, S., Ménard, R., and Thual, O.: Parametric Kalman filter for chemical transport model, *Tellus A*, 68, 31 457, <https://doi.org/10.3402/tellusa.v68.31547>, 2016.
- 840 Pannekoucke, O., Bocquet, M., and Ménard, R.: Parametric covariance dynamics for the nonlinear diffusive Burgers equation, *Nonlin. Processes Geophys.*, 25, 481–495, <https://doi.org/10.5194/npg-25-481-2018>, 2018.



- Ribeiro, M. T., Singh, S., and Guestrin, C.: Why Should I Trust You?: Explaining the Predictions of Any Classifier, in: Proceedings of the 22nd ACM SIGKDD International Conference on Knowledge Discovery and Data Mining, KDD '16, pp. 1135—1144, Association for Computing Machinery, New York, NY, USA, ISBN 9781450342322, <https://doi.org/10.1145/2939672.2939778>, 2016.
- 845 Sacco, M. A., Pulido, M., Ruiz, J. J., and Tandeo, P.: On-line machine-learning forecast uncertainty estimation for sequential data assimilation, *Q. J. R. Meteorol. Soc.*, 150, 2937–2954, <https://doi.org/10.1002/qj.4743>, 2024.
- Sakov, P.: On building the state error covariance from a state estimate, <https://doi.org/10.48550/arXiv.2411.14809>, 2025.
- Sakov, P., Oliver, D. S., and Bertino, L.: An iterative EnKF for strongly nonlinear systems, *Mon. Wea. Rev.*, 140, 1988–2004, <https://doi.org/10.1175/MWR-D-11-00176.1>, 2012.
- 850 Sivashinsky, G. I.: Nonlinear analysis of hydrodynamic instability in laminar flames-I. Derivation of basic equations, *Acta Astronaut.*, 4, 1177–1206, [https://doi.org/10.1016/0094-5765\(77\)90096-0](https://doi.org/10.1016/0094-5765(77)90096-0), 1977.
- Smilkov, D., Thorat, N., Kim, B., Viégas, F. B., and Wattenberg, M.: SmoothGrad: removing noise by adding noise, *CoRR*, abs/1706.03825, <https://doi.org/10.48550/arXiv.1706.03825>, 2017.
- Tang, H. and Glass, J.: On Training Recurrent Networks with Truncated Backpropagation Through time in Speech Recognition, in: 2018 IEEE Spoken Language Technology Workshop (SLT), pp. 48–55, <https://doi.org/10.1109/SLT.2018.8639517>, 2018.
- 855 van Kekem, D. L. and Sterk, A. E.: Travelling waves and their bifurcations in the Lorenz-96 model, *Physica D*, 367, 38–60, <https://doi.org/10.1016/j.physd.2017.11.008>, 2018.
- Wishner, R. P., Tabaczynski, J. A., and Athans, M.: A Comparison of Three Non-Linear Filters, *Automatica*, 5, 487–496, [https://doi.org/10.1016/0005-1098\(69\)90110-1](https://doi.org/10.1016/0005-1098(69)90110-1), 1969.