# DEEP BAYESIAN FILTER FOR BAYES-FAITHFUL DATA ASSIMILATION

Anonymous authors

Paper under double-blind review

#### Abstract

State estimation for nonlinear state space models (SSMs) is a challenging task. Existing assimilation methodologies predominantly assume Gaussian posteriors on physical space, where true posteriors become inevitably non-Gaussian. We propose Deep Bayesian Filtering (DBF) for data assimilation on nonlinear SSMs. DBF constructs new latent variables  $h_t$  in addition to the original physical variables  $z_t$  and assimilates observations  $o_t$ . By (i) constraining the state transition on the new latent space to be linear and (ii) learning a Gaussian inverse observation operator  $r(h_t|o_t)$ , posteriors remain Gaussian. Notably, the structured design of test distributions enables an analytical formula for the recursive computation, eliminating the accumulation of Monte Carlo sampling errors across time steps. DBF trains the Gaussian inverse observation operators  $r(h_t|o_t)$  and other latent SSM parameters (e.g., dynamics matrix) by maximizing the evidence lower bound. Experiments demonstrate that DBF outperforms model-based approaches and latent assimilation methods in tasks where the true posterior distribution on physical space is significantly non-Gaussian.

024 025 026

004

010 011

012

013

014

015

016

017

018

019

021

### 027 028

#### 1 INTRODUCTION

029

Data assimilation (DA) is a crucial technique across various scientific domains. Its primary objective is to estimate the trajectory and current state of a system by integrating an imperfect model with partially informative observations. Specifically, given a series of observations T time steps  $o_{1:T}$ , the goal is to infer the posterior distribution of the system's physical variables  $z_t$ :  $p(z_t|o_{1:t})$ . DA has been widely applied in fields such as weather forecasting (Hunt et al., 2007; Lorenc, 2003; Andrychowicz et al., 2023), ocean research analysis (Ohishi et al., 2024), sea surface temperature prediction (Larsen et al., 2007), seismic wave analysis (Alfonzo & Oliver, 2020), multi-sensor fusion localization (Bach & Ghil, 2023), and visual object tracking (Awal et al., 2023).

A key challenge in DA arises from the non-Gaussian nature of the posterior distributions  $p(z_t|o_{1:t})$ , which results from the inherent nonlinearity in both the system dynamics and observation models. Despite this, many operational DA systems, such as those used in weather forecasting, rely on meth-040 ods like the ensemble Kalman Filter (EnKF) (Evensen, 1994; Bishop et al., 2001) for sequential 041 state filtering (i.e.,  $p(z_t|o_{1:t})$ ) and the four-dimensional variational method (4D-Var) for retrospec-042 tive state analysis (i.e.,  $p(z_t|o_{1:T}), t < T$ ). These approaches assume Gaussianity in their test 043 distributions  $q(z_t|o_{1:t})$  or  $q(z_t|o_{1:T})$ , a simplification driven by computational constraints. While 044 exact methods such as bootstrap Particle Filters (PF) or sequential Monte Carlo (SMC) (Chopin & Papaspiliopoulos, 2020; Daum & Huang, 2007; Hu & van Leeuwen, 2021) could compute the true posterior, their performance degrades significantly when the number of particles is insuffi-046 cient (Beskos et al., 2014). This issue is exacerbated in high-dimensional systems, making SMC 047 approaches impractical for many physical problems. 048

To address these limitations, we propose a novel variational inference approach called Deep
Bayesian Filtering (DBF) for posterior estimation. Our strategy consists of two main components:
(i) constraining the test distribution to remain Gaussian to ensure computational tractability, and,
in cases where the original dynamics are nonlinear, (ii) leveraging a nonlinear mapping to enhance
the expressive capability of the test distribution. The DBF methodology diverges into two paths
depending on the nature of the system dynamics, whether linear Gaussian or nonlinear:

**Linear dynamics** When the system's dynamics  $p(z_{t+1}|z_t)$  are linear, DBF assumes Gaussianity in the original space, similar to traditional methods. However, DBF introduces the concept of the inverse observation operator (IOO; see also Frerix et al. 2021) to construct Gaussian test distributions  $q(z_t|o_{1:t})$ . The IOO, along with any unknown system parameters, are trained to minimize the Kullback-Leibler divergence between the test distribution  $q(z_t|o_{1:t})$  and the true posterior  $p(z_t|o_{1:t})$ . The IOO and the system parameters are trained without teacher signals  $z_t$ .

**Nonlinear dynamics** In the more common case of nonlinear dynamics, DBF operates in a latent space, assuming Gaussianity in the latent variables  $h_t$ . The original physical variables are recovered through a nonlinear mapping function  $\phi$ , implemented via neural networks (NNs). This nonlinear mapping allows for a more flexible representation of the test distribution  $q(z_t|o_{1:t})$ . The IOO and other parameters are trained in a supervised manner (i.e.,  $z_t$  is used during training).

For state space models (SSMs) with nonlinear dynamics, DBF functions as a variational autoencoder (VAE) that adheres to the Markov property. Posterior distributions of the latent variables  $h_t$  are expressed in a Bayesian framework. This approach is closely related to dynamical VAEs (DVAEs, Girin et al. 2021 for a review), which use VAEs to model time-series data. However, DBF distinguishes itself by its posterior design. Unlike DVAEs, where Monte Carlo sampling is required for inference (see Sec. 2.6.1), DBF allows for the analytical computation of the prediction step, recursively computing posteriors through closed-form expressions.

When applied to problems with nonlinear or unknown dynamics, DBF can be interpreted as learning the Koopman operator (Koopman, 1931) using NNs. The discovery of such latent spaces and
operators through machine learning has been extensively studied (Takeishi et al., 2017; Lusch et al.,
2018; Azencot et al., 2020) and will be experimentally validated through the handling of nonlinear
filtering tasks involving chaotic dynamics.

- 078 Key contributions of the proposed DBF methodology include:
  - DBF is the first VAE-based model for time-series data that maintains a posterior structure faithful to the Markov property in SSMs.
  - For linear dynamics, DBF extends the Kalman Filter (KF) with a more flexible observation update. A simple object-tracking experiment is presented in Sec. 3.1. Additional experiments in Sec. B demonstrates that DBF can infer unknown system parameters via training.
  - For nonlinear dynamics, DBF constructs a new latent space for data assimilation, allowing for the analytical integration of time steps and preventing the accumulation of Monte Carlo sampling errors. This is accomplished through the application of Koopman operator theory, which ensures that the model's representational power is maintained, as long as the latent space is sufficiently high-dimensional (see Sec. 3.2 and 3.3).
  - As a generative model, DBF estimates the uncertainty of the physical variables  $z_t$ , in contrast to 3D- and 4D-Var, which yield only point estimates (see Sec. 3.2 and Fig. 3).
    - The linear constraint on dynamics stabilizes the training process, which is known to be unstable in standard recurrent NNs (see Sec. 3.3 and Fig. 6).

DBF has demonstrated superior performance over classical DA algorithms and latent assimilation methods in scenarios with highly non-Gaussian posteriors, particularly in the presence of strongly nonlinear observation operators or large observation noise.

### 2 Method

060

079

081

082

084

085

090

091

092

094

096

097 098

099

101

2.1 INFERENCE OF PHYSICAL VARIABLES IN A STATE-SPACE MODEL

102 A physical system is defined by variables  $z_t$ , with its evolution described by the dynamics model 103  $p(z_{t+1}|z_t) = \mathcal{N}(z_{t+1}; f(z_t), Q)$ , where  $\mathcal{N}(x|\mu, \Sigma)$  denotes a Gaussian whose mean and covari-104 ance are  $\mu$  and  $\Sigma$ . The nonlinear function f is the dynamics operator and Q is the system 105 covariance. The Markov property holds, as  $z_{t+1}$  depends only on  $z_t$ . An observation model 106  $p(o_t|z_t) = \mathcal{N}(o_t; h(z_t), R)$  relates observations to physical variables via the observation opera-107 tor h and covariance R. The panel (a) of Fig. 1 shows the system's graphical model. The objective of sequential DA is to compute the posterior of  $z_t$  given  $o_{1:t}$ .



Figure 1: Panel (a) shows the graphical model for the simplest SSM. If the dynamics of the original SSM is linear, DBF assimilates on that space. Panel (b) shows the graphical model for the SSM assumed for SSM with nonlinear dynamics. Panel (c) shows the inference structure of our methodology for SSM with nonlinear dynamics.

#### 2.2 KF FOR LINEAR DYNAMICS, LINEAR OBSERVATIONS

In the KF, the dynamics and observation models are both linear Gaussian. Given that the dynamics and observation operators f, h are linear, we can represent them using matrices A and C, respectively. All matrices (A, C, Q, and R) are constant. The filter distribution  $p(z_t|o_{1:t})$  remains Gaussian, provided that the initial distribution  $p(z_1)$  is Gaussian. We can recursively compute the posterior parameters (means  $\mu_t$  and covariance matrices  $\Sigma_t$ ) using the following equations:

$$\mu_t = \Sigma_t (A \Sigma_{t-1} A^T + Q)^{-1} A \mu_{t-1} + K_t (o_t - H A \mu_{t-1}), \tag{1}$$

$$\Sigma_t^{-1} = (A\Sigma_{t-1}A^T + Q)^{-1} + HR^{-1}H^T,$$
(2)

where  $K_t = (A\Sigma_{t-1}A^T + Q)H^T(H(A\Sigma_{t-1}A^T + Q)H^T + R^{-1})^{-1}$  is the Kalman Gain.

#### 2.3 DBF FOR LINEAR DYNAMICS, NONLINEAR OBSERVATIONS

In this scenario, Gaussianity of the test distribution is lost during the KF update step. We introduce an inverse observation operator (IOO)  $r(z_t|o_t)$  (see also Frerix et al. 2021):

$$p(z_t|o_{1:t}) = \frac{p(o_t|z_t)p(z_t|o_{1:t-1})}{p(o_t|o_{1:t-1})} \propto \frac{r(z_t|o_t)}{\rho(z_t)}p(z_t|o_{1:t-1}),$$
(3)

where  $r(z_t|o_t) = \frac{p(o_t|z_t)\rho(z_t)}{\int p(o_t|z_t)\rho(z_t)dz_t}$  and  $\rho(z_t)$  is a prior virtually introduced for the IOO. By approximating both the IOO and the virtual prior as Gaussians,  $r(z_t|o_t) = \mathcal{N}(f_{\theta}(o_t), G_{\theta}(o_t))$  and  $\rho(z_t) = \mathcal{N}(m, V)$ , respectively, the posterior  $q(z_t|o_{1:t})$  can be analytically computed as a Gaussian, where the mean  $\mu_t$  and covariance  $\Sigma_t$  are given as:

151

120

121

122

123 124 125

126

132 133

134 135

136 137

138 139

140

141 142

143

$$\mu_t = \Sigma_t (A \Sigma_{t-1} A^T + Q)^{-1} A \mu_{t-1} + G_\theta(o_t)^{-1} f_\theta(o_t) - V^{-1} m, \tag{4}$$

$$\Sigma_t^{-1} = (A\Sigma_{t-1}A^T + Q)^{-1} + G_\theta(o_t)^{-1} - V^{-1},$$
(5)

where  $f_{\theta}(o_t)$  and  $G_{\theta}(o_t)$  are NNs with parameters  $\theta$ , and m and V are constants set to m = 0and  $V = 10^8 I$ . These values bias the NNs' outputs without affecting performance. The initial distribution  $q(z_1)$  is taken to be a Gaussian with  $\mu_1 = 0$  and  $\Sigma_1 = 100I$ .

The recursive formula for the exact posterior (Equation 3) requires no approximation. Thus, DBF computes the exact posterior when the true IOO  $r_{true}(h_t|o_t)$  is Gaussian, i.e., the SSM is a Linear-Gaussian State Space (LGSS). In that case, the posterior update formula agrees with the KF (see Equations 1, 2 and 4, 5). The key difference is that nonlinear functions are applied to both the mean,  $f_{\theta}(o_t)$ , and the covariance,  $G_{\theta}(o_t)$ . In the KF,  $f_{\theta}(o_t)$  is linear, and  $G_{\theta}(o_t)$  is a constant matrix (see Equations 1 and 2).  $G_{\theta}(o_t)$ 's dependence on observations allows flexible adjustment of the new observation's impact on state estimation. The importance of adjusting the internal state updates based on observations has also been discussed in recent SSM-based approaches (Gu & Dao, 2023).

# 162 2.4 DBF FOR NONLINEAR DYNAMICS, LINEAR/NONLINEAR OBSERVATIONS

164 In this scenario, the Gaussianity of the test distribution is lost during the predict step, making it 165 impossible to apply the original dynamics over the physical variables  $z_t$ . Therefore, we introduce a new set of latent variables  $h_t$  and assume a dynamics model over  $h_t$ :  $p(h_{t+1}|h_t) = \mathcal{N}(h_{t+1}|Ah_t, Q)$ 166 (see panel (b) in Fig. 1). The IOO maps observations into the latent variables  $h_t$ :  $r(h_t|o_t)$ . The 167 recursive formula follows Equations 4 and 5. To retrieve the distribution of the original physical 168 variables  $z_t$ , we introduce an emission model  $p(z_t|h_t) = \mathcal{N}(z_t; \phi(h_t), R)$ , where  $\phi$  is represented 169 by a NN. By marginalizing over  $h_t$  with this emission model, a trained DBF can generate samples 170 of  $z_t$  that follow the test distribution  $q(z_t|o_{1:t})$  given observations  $o_{1:t}$ . 171

Although the dynamics operator A for the latent variables  $h_t$  is linear, it can express any nonlinear 172 dynamics if the latent space is sufficiently high-dimensional. The Koopman operator (Koopman, 173 1931) provides a framework for representing nonlinear systems by mapping observables—functions 174 of the system's state-into a higher-dimensional space where the dynamics are linear. For a system 175  $z_{t+1} = f(z_t)$ , the Koopman operator  $\mathcal{K}$  is a linear operator acting on a set of observables g(z), such 176 that  $\mathcal{K}g(z_t) = g(f(z_t))$ . This reformulates the system as  $h_{t+1} = Ah_t$  in the latent space, where A 177 is the dynamics matrix learned by DBF. While the physical dynamics f(z) are nonlinear, the Koop-178 man operator ensures the existence of an embedding that linearizes the dynamics, enabling recursive 179 computation of test distributions. Discovering such embeddings in finite dimensions has been widely 180 studied (Takeishi et al., 2017; Lusch et al., 2018; Azencot et al., 2020). In high-dimensional simula-181 tions, the true degrees of freedom are often far fewer than the simulated variables, making surrogate 182 modeling with the Koopman operator a promising approach to reducing computational costs.

#### 2.5 TRAINING

When assimilating in the physical space (i.e., when the dynamics are linear), we train the IOO (i.e.,  $f_{\theta}$  and  $G_{\theta}$ ) by optimizing the evidence lower bound (ELBO) without using the teacher signal  $z_t$ :

$$\log p(o_{1:T}) = \sum_{t=1}^{T} \log p(o_t|o_{1:t-1}) \ge -\mathcal{L}_{\text{ELBO}},$$
  
$$\mathcal{L}_{\text{ELBO}} = -\sum_{t=1}^{T} \int q(h_t|o_{1:t}) \log p(o_t|h_t) dh_t + KL[q(h_t|o_{1:t})||q(h_t|o_{1:t-1})], \quad (6)$$

194

195

183 184

185

187 188 189

> $\mathcal{L}_{\text{ELBO}} = -\sum_{t=1} \int q(h_t|o_{1:t}) \log p(o_t|h_t) dh_t + KL[q(h_t|o_{1:t})||q(h_t|o_{1:t-1})], \quad (6)$ where KL[p||q] denotes the Kullback-Leibler divergence between distributions p and q (see Sec. A.1 in the appendix for the derivation). Here,  $q(h_1|o_{1:0}) = q(h_1)$  is the initial distribution. If the SSM

contains any unknown parameters, we can train these parameters as well.
 For SSMs with poplinger or unknown dynamics, we have two approaches;

For SSMs with nonlinear or unknown dynamics, we have two approaches:

**Strategy 1** Pretrain the Koopman operator, which consists of the nonlinear mapping from  $z_t$  to  $h_t$ , the linear dynamics between  $h_t$  and  $h_{t+1}$  represented by matrix A, and the reverse nonlinear mapping from  $h_t$  to  $z_t$  denoted by  $\phi$ . With these components (A and  $\phi$ ) of the Koopman operator, the method designed for linear dynamics can be applied. For pretraining, we require samples of  $z_t$  or the SSM for the physical variables to generate these samples. Pairs of  $z_t$  and  $o_t$  are not necessary, as the training for the linear dynamics (A and  $\phi$ ) and the IOO ( $r(h_t|o_t)$ ) can be performed separately.

**Strategy 2** Train all components (the matrix A, the stochastic mapping  $p(z_t|h_t) = \mathcal{N}(z_t; \phi(h_t), \operatorname{diag}[\sigma^2])$ , and the IOO) simultaneously. In this case, samples of  $(z_t, o_t)$  pairs or the SSM for both physical and observation variables to generate these sample pairs are required during training. Note that the physical variables  $z_t$  are not required for inference, ensuring that real-time applications are not hindered by the need for  $z_t$  during training. The parameters are optimized by maximizing a joint ELBO,  $\mathcal{L}_{\text{ELBO, joint}}$ , via supervised training:

211 212

212  
213 
$$\log p(o_{1:T}, z_{1:T}) = \sum_{t=1}^{T} \log p(o_t, z_t | o_{1:t-1}, z_{1:t-1}) \ge -\mathcal{L}_{\text{ELBO,joint}},$$
214  
215 
$$\int_{T} \sum_{t=1}^{T} \int_{T} g(h_t | o_{t-1}) \log p(z_t | h_t) dh_t + KL[g(h_t | o_{t-1})] | g(h_t | o_{t-1})]$$
(7)

$$\mathcal{L}_{\text{ELBO,joint}} = -\sum_{t} \int q(h_t|o_{1:t}) \log p(z_t|h_t) dh_t + KL[q(h_t|o_{1:t})||q(h_t|o_{1:t-1})].$$
(7)

(See Sec. A.2 in the appendix for the derivation). We have replaced  $q(h_t|o_{1:t}, z_{1:t})$  with its special case  $q(h_t|o_{1:t})$  as our objective is to give the best estimate of  $z_t$  given observations  $o_{1:t}$ .

- 219 2.6 RELATED WORKS
- 221 2.6.1 DYNAMICAL VARIATIONAL AUTOENCODERS

DVAEs (see Girin et al. 2021 for a review) are a broad class of models incorporating time-series architectures into VAEs, with DBF as a specialized subcategory. Key differences include (i) the posterior design and realization of the dynamics step, and (ii) the loss function.

226 **posterior design** Our strategy for the test distribution is to incorporate an appropriate architecture 227 that reflects the Markov property in the time dimension of the test distribution. The IOO,  $r(h_t|o_t)$ , 228 and the linear dynamics model serve as key instruments in constructing the test posterior distribu-229 tions. A distinguishing feature of our methodology is that each component's role is defined with 230 respect to the Markov property of the state-space model (SSM) and is clearly differentiated from 231 other components involved in posterior construction. For example, the IOO influences only the up-232 date step and does not affect the prediction step. We refer to this methodology as "Bayes-Faithful" due to its tailored design for SSMs that exhibit the Markov property. 233

In contrast, the test posterior distributions in DVAEs are constructed using RNNs. The complexity
of the transition model prevents the analytical computation of latent variables across time steps. As
a result, these values can only be estimated via Monte Carlo sampling. Consequently, during inference, successive Monte Carlo sampling ("cascade trick"; Girin et al. 2021) becomes unavoidable.

loss function DBF takes the ELBO from factorized density  $\log p(o_t|o_{1:t-1})$  in  $\log p(o_{1:T}) = \sum_t \log p(o_t|o_{1:t-1})$ :

240 241 242

243 244

245 246

247

238 239

222

223

224

225

$$\log p(o_{1:T}) \geq \sum_{t=1}^{T} (E_{q(h_t|o_{1:t})}[\log p(o_t|h_t)] - KL[q(h_t|o_{1:t})|q(h_t|o_{1:t-1})]).$$
(8)

On the other hand, DVAEs take the ELBO from probability density with all the observations at once.

$$\log p(o_{1:T}) \geq E_{q(h_{1:T}|o_{1:T})}[\log p(h_{1:T}, o_{1:T}) - \log q(h_{1:T}|o_{1:T})].$$
(9)

Therefore, DBF seeks for the filtered distributions  $q(h_t|o_{1:t})$  whereas DVAEs model the smoother distributions  $q(h_t|o_{1:T})$ . Again, for DVAEs, to evaluate the expected values in Equation 9, we need to undergo successive Monte-Carlo sampling over T variables  $(h_{1:T})$  (see also Sec. A.3).

Assuming linear Gaussian dynamics and a Gaussian IOO, DBF allows for the analytical integration of  $q(h_t|o_{1:t-1})$ , resulting in a structured encoder. This structured posterior enables the recursive computation of the filtered distribution  $q(h_t|o_{1:t})$  without relying on Monte Carlo sampling, setting it apart from other DVAEs. By constraining the dynamics to be linear, DBF ensures exact integration without the accumulation of Monte Carlo sampling errors across time steps.

256 Moreover, the linear assumption helps DBF mitigate the instability issues commonly faced when training standard RNNs. The linearity of the latent dynamics is also assumed in normalizing Kalman 257 Filter (de Bézenac et al., 2020) and Kalman variational auto-encoder (Fraccaro et al., 2017). SSMs 258 are increasingly favored for modeling long-range dependencies (Gu & Dao, 2023). S4 (Gu et al., 259 2022) learns linear dynamics in the latent space, proposing an efficient computation algorithm that 260 outperforms transformers on datasets with long-range dependencies. LS4 (Zhou et al., 2023) extends 261 S4 by introducing stochasticity through a VAE-like structure. Both LS4 and DBF employ linear 262 SSMs and Gaussian posterior approximations, but DBF updates the mean and covariance using a 263 recursive formula based on Bayes' rule, while the construction of posteriors in LS4 is not recursive. 264

- 265
- 265 2.6.2 KF-BASED METHODS

Various approaches have been explored to address LGSS limitations, including linearizing the model
via first-order approximations like the extended Kalman Filter (EKF), approximating populations
with a Gaussian distribution in the ensemble Kalman Filter (EnKF; Evensen 1994), and using NNs to approximate the Kalman gain (Revach et al., 2022). EKFNet (Xu & Niu, 2024) assumes EKF for

the construction of test distribution and train the SSM parameters. Auto-EnKF (Chen et al., 2022;
2023) leverages EnKF and train the model by optimizing ELBO. The EnKF and its variants (e.g., ETKF; Bishop et al. 2001) are commonly used in real-time data assimilation for weather forecasting.
However, these methods rely on the KF's posterior update equations, limiting the expressivity of
the distributions. Additionally, computations for covariance matrices become challenging in highdimensional spaces, requiring specialized techniques for computational efficiency.

276 277

2.6.3 SAMPLING-BASED METHODS

The Particle Filter is a popular method for assimilating any posterior. However, achieving adequate particle density in high-dimensional state spaces poses significant challenges. Insufficient density of particles leads to particle degeneracy, where few particles explain the observed data (Beskos et al., 2014). In contrast, DBF directly learns to position density through the IOO, offering advantages for high-dimensional tasks. The Particle Flow Filter (PFF; Daum & Huang 2007; Hu & van Leeuwen 2021) addresses particle degeneracy by moving particles according to gradient flow and effectively scales to nonlinear SSMs with hidden state dimensions up to 1000 (Hu & van Leeuwen, 2021).

285

#### 286 2.6.4 APPROXIMATE MAP ESTIMATION METHOD

MAP estimation is used to identify the high-density point of the posterior in high-dimensional space, such as in weather forecasting (Lorenc, 2003; Frerix et al., 2021). Even if the computation of the posterior  $p(h_t|o_{1:t})$  is intractable, we can optimize  $\log p(h_t|o_{1:t}) = \log p(o_t|h_t) + \log p(h_t|o_{1:t-1})$  if we can describe  $p(o_t|h_t)$  and  $p(h_t|o_{1:t-1}) = \int p(h_t|h_{t-1})p(h_{t-1}|o_{1:t-1})dh_{t-1}$  explicitly. In practice, we cannot access  $p(h_{t-1}|o_{1:t-1})$  and therefore the integral  $\int p(h_t|h_{t-1})p(h_{t-1}|o_{1:t-1})dh_{t-1}$ , so we only compute the mean. The downside is that sequential computation of the covariance matrix of  $p(h_t|o_{1:t-1})$  is impossible.

295 2.6.5 NN-BASED PDE SURROGATE

Recently, there have been attempts to approximate partial differential equations (PDEs) using NNs.
In this study, we experimented with one of the latest methods, PDE-refiner (Lippe et al., 2023), but
its performance was poor and was excluded from the experiments. We suspect this is because PDE-refiner, designed for constructing PDE surrogates, does not handle noisy observations well, making
it sensitive to noise. However, we confirmed that it performs well under noiseless conditions.

302 303

304

308

### 3 EXPERIMENTS

We evaluate the performance of DBF on three tasks: a linear dynamics problem (object tracking) and two nonlinear dynamics problems (double pendulum and Lorenz96). An additional experiment on linear dynamics (moving MNIST) is presented in Sec. B of the appendix.

Linear dynamics: object tracking Linear dynamics are applicable in real-world object tracking tasks involving objects with continuous motion, either stationary or nearly uniform linear. An object tracker can be easily constructed by replacing the IOO with an object detector that performs independent detection for each frame. DBF adaptively weights the detected object positions based on the confidence in the current detector's estimates, significantly enhancing tracking robustness. The results are compared against those of the KF.

314

315 Nonlinear dynamics: double pendulum and Lorenz96 For nonlinear dynamics problems, such 316 as the double pendulum and Lorenz96, DBF constructs a new latent space in addition to the original 317 physical space. Here, we took Strategy 2 in Sec.2.5 for the training: we simultaneously train NNs 318 for the IOO, nonlinear observation operator  $\phi$ , the dynamics matrix A, and the emission model's 319 standard deviation. We compare the performance of DBF with the classical DA algorithms (EnKF, 320 ETKF, PF), state-of-the-art assimilation methodologies (PFF Daum & Huang 2007; Hu & van 321 Leeuwen 2021, KalmanNet Revach et al. 2022), and DVAE-based approaches (deep Kalman Filter; DKF, Krishnan et al. 2015; 2016, variational recurrent neural network; VRNN, Chung et al. 322 2015, and stochastic recurrent neural network; SRNN, Fraccaro et al. 2016). DBF and other DVAEs 323 are trained by optimizing the evidence lower bound (ELBO), as described in Sec. 2.5.

324 For nonlinear dynamics experiments, we generate random initial conditions and evolve them using 325 the dynamics. Synthetic observations are produced by applying the observation operator with addi-326 tive noise. Noise levels, observation operators, and further details are given in Sec. C, C.1, C.2, and 327 C.3. Sec. C also provides computationally efficient parametrization of the latent dynamics matrix.

328

330

367

#### LINEAR DYNAMICS: OBJECT TRACKING 3.1

331 In a single-object tracking problem, a detector identifies a bounding box for the object in each 332 frame, and these boxes are then connected across frames. When the object is not fully visible or is obscured, the detector often fails to accurately determine its position. In such scenarios, the KF 333 aids by predicting and assimilating the object's true position. However, a key limitation of the KF 334 is its reliance on a fixed observation model throughout the tracking process. While low-confidence 335 observations can provide valuable approximate position information, they may also mislead the 336 tracker with inaccurate data, potentially degrading overall tracking performance. 337

338 We demonstrate that DBF can enhance tracking stability without requiring additional training. During the computation of the posterior  $p(z_t|o_{1:t})$  from  $p(z_t|o_{1:t-1})$ , the importance of the observation 339  $o_t$  is regulated via  $G_{\theta}(o_t)$ . This allows the observational confidence to be effectively incorporated 340 into the posterior estimation. We evaluate the tracking performance using the "airplane" category 341 from the LaSOT dataset (Fan et al., 2019; 2021). 342

343 We use the first 1,000 frames from 20 videos for evaluation. The first 10 videos serve as a validation 344 set for determining filter parameters (see Sec. C.1), while the performance is assessed using videos 11-20. Each set of 1,000 frames is divided into 20 subsets of 50 frames. Filters are initialized at the 345 ground truth coordinates of the bounding box in the first frame, after which each filter is responsible 346 for tracking the bounding box throughout the subset. We employ the YOLOv8n model (Jocher 347 et al., 2023) as the object detector. The detector outputs the bounding box position, X, along with 348 a confidence score, c. A detection threshold of 0.01 is applied. When multiple bounding boxes are 349 detected, the one with the highest posterior probability is selected. 350

The bounding box coordinates are used as  $f_{\theta}(o_t) = X$ . We experiment with linear confidence 351  $G_{\theta}(o_t) \propto c$  and squared confidence  $G_{\theta}(o_t) \propto c^2$  and find that the squared confidence  $G_{\theta}(o_t) \propto c^2$ 352 perform better. For further settings, see Sec. C.1. 353

354 Figure 2 presents the results. The left panel provides an illustrative example comparing the two 355 tracking algorithms. The KF tracker is visibly influenced by false detections, being pulled toward a 356 coordinate value of approximately 150 during frames 15–17. In contrast, the DBF tracker maintains stable predictions under the same conditions. The middle and the right panels offer a quantitative 357 comparison of KF and DBF in terms of intersection over union (IoU). Both filters perform well 358 in estimating bounding box positions in frames without detections. However, DBF demonstrates 359 a significant performance advantage in frames with low-confidence (c < 0.1) detections. This 360 improvement can be attributed to DBF's flexibility, allowing it to adaptively decide whether to trust 361 low-confidence observations or disregard them. 362



371 Figure 2: Left panel: x coordinate of bounding box center estimated with KF and DBF. Colored dots 372 show the coordinates of the bounding box reported by the YOLO model. The red band (frames 18 373 - 20) shows frames where the detection network reports no bounding boxes. Middle panel: fraction 374 of frames with IoU > 0.1 for each tracker. Detections with a confidence score greater than 0.1 are 375 categorized as "High", those below 0.1 as "Low", and those below 0.01 as "Missed". Right panel: 376 mean IoU for the three categories. The performance gain of DBF from KF is considerable in frames 377 with low-confidence detections.

#### 3.2 NONLINEAR DYNAMICS 1: DOUBLE PENDULUM



Figure 3: A schematic figure (panel a) and results for double pendulum experiments. Panel (b) shows the RMSE of angle velocities (averaged over  $\omega_1$  and  $\omega_2$ ) over time steps. Panels (c) and (d) show example histograms for normalized errors in DBF and ETKF samples compared against the unit Gaussian  $\mathcal{N}(x; \mu = 0, \sigma^2 = 1)$ . The small table compares the Jeffreys divergence of normalized errors and the unit Gaussian between DBF, EnKF, and ETKF predictions.

Table 1: RMSE at the final ten steps of assimilation in double pendulum experiments.

	$\sigma =$	: 0.1	$\sigma =$	: 0.3	$\sigma = 0.5$		
	$\theta \qquad \omega$		heta $ heta$		heta	$\omega$	
DBF	$\textbf{0.03} \pm \textbf{0.01}$	$\textbf{0.21} \pm \textbf{0.04}$	$\textbf{0.05} \pm \textbf{0.02}$	$\textbf{0.26} \pm \textbf{0.05}$	$\textbf{0.06} \pm \textbf{0.01}$	$\textbf{0.36} \pm \textbf{0.04}$	
EnKF	$0.05\pm0.00$	$0.33\pm0.07$	$0.14\pm0.01$	$0.71\pm0.09$	$0.24\pm0.01$	$1.17\pm0.22$	
ETKF	$0.05\pm0.01$	$0.46\pm0.08$	$0.22\pm0.05$	$1.41\pm0.41$	$0.36\pm0.08$	$2.70 \pm 1.25$	
PF	$0.05\pm0.00$	$0.63\pm0.24$	$0.21\pm0.14$	$1.41 \pm 1.30$	$0.32\pm0.08$	$2.36 \pm 2.29$	
PFF	$1.27\pm0.29$	$1.04\pm0.15$	NA	$5.99 \pm 1.09$	$5.88 \pm 0.67$	NA	
KNet	NA	NA	NA	NA	NA	NA	
VRNN	$0.04\pm0.01$	$0.44\pm0.19$	$0.06\pm0.02$	$0.35\pm0.14$	$0.08\pm0.04$	$0.40\pm0.16$	
SRNN	$0.05\pm0.02$	$0.52\pm0.18$	$0.06\pm0.02$	$0.44\pm0.08$	$0.08\pm0.03$	$0.52\pm0.22$	
DKF	$0.12\pm0.02$	$2.70\pm0.28$	$0.17\pm0.03$	$2.61\pm0.74$	$0.23\pm0.04$	$2.61\pm0.56$	

407 408

378

379 380

382

383

384

385

386

387

389 390

391

392

393

394

396 397

This section presents our experiments with a double pendulum system, selected for its nonlinear and
chaotic behavior. The pendulum consists of two 1 kg masses, P1 and P2, connected by two 1 meter
bars, B1 and B2. One end of the bar B1 is fixed at the origin ("O"), with the other end attached to
P1. Mass P2 is connected to P1 via bar B2. A schematic of the setup is shown in panel (a) of Fig. 3.

413 We use the angles  $\theta_1$  and  $\theta_2$ , and the two angular velocities,  $\omega_1$  and  $\omega_2$ , as target physical variables. 414 The latent dimension for DBF, VRNN, SRNN, and DKF is set to 50. For the choice of the latent 415 dimensions, refer to Sec. E.1. Observation data consists of the two-dimensional spatial positions of 416 masses  $P_1$  and  $P_2$ , corrupted by Gaussian noise. The observation operator combines trigonometric 417 functions for  $\theta_1$  and  $\theta_2$  which are highly nonlinear. Experiments are conducted with noise levels of 418  $\sigma = 0.1, 0.3,$  and 0.5 [m], with a time step of 0.03 [s] between observations. In the emission model 419  $p(z_t|h_t)$ , we assume von Mises distributions for  $\theta_1$  and  $\theta_2$ , while  $\omega_1$  and  $\omega_2$  follow Gaussians.

420 Table 1 presents the RMSE between the physical variables and the mean of the filtered distribution. 421 For both the angles  $\theta$  and angle velocities  $\omega$ , we compute the averages of the two variables across 422 two pendulums. Training for KalmanNet was unsuccessful under all conditions. For the DVAEs, 423 we exclude failed initial conditions (2/15 for VRNN and DKF, and 3/15 for SRNN) when calculating the RMSE. DBF outperforms both model-based and latent assimilation methods across all 424 settings, showing significant improvements in estimating  $\omega$ , which cannot be inferred from a single 425 observation. Fig. 3 (b) illustrates an example of RMSE evolution during assimilation, where DBF 426 consistently outperforms the other methods. The assimilation of  $\omega$  occurs within the first  $\sim 20$  steps, 427 maintaining an excellent estimation accuracy throughout the experiment. 428

429 A key feature of DBF is its ability to generate samples of  $z_t$  and assess the uncertainty in state 430 estimates. To evaluate this capability, we analyze the distributions of normalized errors defined as 431  $\epsilon_{norm,t,i} = (z_{t,sample,i} - z_{t,i})/\delta_i$ , where  $z_{t,i}$  represents the true value of dimension *i* at time *t*, and  $\delta_i$  is the standard deviation of  $z_{t,sample,i}$ . We collect  $\epsilon_{norm,t,i}$  across all time steps, focusing 432 on  $i = \omega_1$  and  $i = \omega_2$ , since  $\theta_1$  and  $\theta_2$  follow von Mises distributions. If the uncertainty esti-433 mates are accurate,  $\epsilon_{norm,t,i}$  should approximate a Gaussian distribution with a standard deviation 434 of one. To quantify the accuracy, we compute the symmetric KL divergence (Jeffreys divergence) 435  $KL_{sym}[p,q] = (KL[p||q] + KL[q||p])/2$  between the histogram of  $\epsilon_{norm,t,i}$  and a unit Gaussian. 436 DBF exhibits very low  $KL_{sym}$  values, indicating accurate error estimation. Panels (c) and (d) display example histograms of  $\epsilon_{norm,t,i}$  for DBF and ETKF. 437

438 439

440

448

449

#### 3.3 NONLINEAR DYNAMICS 2: LORENZ96

In the final experiment, we focus on state estimation in the 441 Lorenz96 model (Lorenz, 1995), a benchmark for testing data 442 assimilation algorithms on noisy, nonlinear observations. The 443 Lorenz96 model describes the evolution of a one-dimensional ar-444 ray of variables, each representing a physical quantity over a spatial 445 domain, like an equilatitude circle. The dynamics are governed by 446 the following coupled ordinary differential equations: 447

$$\frac{dz_i}{dt} = (z_{i+1} - z_{i-2})z_{i-1} - z_i + F, \quad i = 1, \dots, N,$$
(10)

where  $z_i$  is the value at grid i, N is the number of grid points, and F 450 is external forcing. For our experiments, we take (F, N) = (8, 40). 451

452 We consider two observation operators. The first adds Gaussian 453 noise to direct observations:  $o_{t,j} = z_{t,j} + \epsilon$ , with noise lev-454 els  $\sigma = 1, 3, 5$ . The second uses a nonlinear operator:  $o_{t,j} =$ 455  $\min(z_{t,i}^4, 10) + \epsilon$ , with the same noise levels. The dynamic range of  $z_{t,j}$  is around  $\pm 10$ , and observations are capped at 10 when  $z_{t,j}$ 456 exceeds 1.8. This makes it highly challenging for classical DA 457



Figure 4: A Hovmöller diagram for one of data in the test set. The observation operator is nonlinear,  $o_{t,i} =$  $min(z_{t,j}^4, 10) + \epsilon.$ 

methods, as each observation offers limited information. The filter must integrate data over long 458 timesteps, where nonlinear dynamics distort the probability distribution. Fig. 4 illustrates observa-459 tions and target values. All models use 80 observation steps with a 0.03 time interval. The latent 460 dimension for DBF, VRNN, SRNN, and DKF is set to 800 (for the choice of the latent dimension in 461 DBF, see Sec. E.2). For further details for the experiment, see Sec. C.3. 462

463 464

Table 2: RMSE at the final ten steps of assimilation in Lorenz96 experiments.

	d	irect observatio	n	nonlinear observation					
	$\sigma = 1$	$\sigma = 3$	$\sigma = 5$	$\sigma = 1$	$\sigma = 3$	$\sigma = 5$			
DBF	$0.53\pm0.04$	$\textbf{0.82} \pm \textbf{0.03}$	$\textbf{1.16} \pm \textbf{0.07}$	$1.08 \pm 0.15$	$\textbf{1.29} \pm \textbf{0.18}$	$\textbf{1.65} \pm \textbf{0.17}$			
EnKF	$0.31\pm0.01$	$0.83\pm0.10$	$1.73\pm0.12$	$4.69\pm0.14$	$3.93\pm0.08$	$3.81\pm0.07$			
ETKF	$\textbf{0.30} \pm \textbf{0.01}$	$1.06\pm0.15$	$2.42\pm0.11$	$4.57\pm0.25$	$4.28\pm0.04$	$4.23\pm0.07$			
PF	$2.80\pm0.04$	$3.12\pm0.06$	$3.62\pm0.13$	$6.05\pm0.16$	$4.95\pm0.12$	$4.58\pm0.14$			
PFF	$0.60\pm0.02$	$1.00\pm0.05$	$2.20\pm0.09$	$3.75\pm0.09$	$3.85\pm0.04$	$3.83\pm0.11$			
KNet	$0.60\pm0.02$	$1.81\pm0.05$	$3.02\pm0.09$	$2.97 \pm 0.21$	$3.47\pm0.17$	$3.99\pm0.25$			
VRNN	$3.67\pm0.06$	$3.67\pm0.06$	$3.67\pm0.06$	$3.69\pm0.04$	$2.51\pm0.79$	$3.67\pm0.06$			
SRNN	$3.08\pm0.56$	$3.63\pm0.05$	$3.40\pm0.29$	$3.30\pm0.81$	$3.62\pm0.41$	$2.96 \pm 0.32$			
DKF	3.70	NA	NA	NA	NA	NA			
	DBF EnKF ETKF PF KNet VRNN SRNN DKF	$\begin{array}{c} & \frac{\sigma=1}{\sigma=1} \\ \hline \\ DBF & 0.53 \pm 0.04 \\ EnKF & 0.31 \pm 0.01 \\ ETKF & \textbf{0.30} \pm \textbf{0.01} \\ PF & 2.80 \pm 0.04 \\ PFF & 0.60 \pm 0.02 \\ KNet & 0.60 \pm 0.02 \\ VRNN & 3.67 \pm 0.06 \\ SRNN & 3.08 \pm 0.56 \\ DKF & 3.70 \\ \end{array}$	$ \begin{array}{c} \hline \sigma = 1 & \sigma = 3 \\ \hline \sigma = 1 & \sigma = 3 \\ \hline \text{DBF} & 0.53 \pm 0.04 & \textbf{0.82} \pm \textbf{0.03} \\ \text{EnKF} & 0.31 \pm 0.01 & 0.83 \pm 0.10 \\ \text{ETKF} & \textbf{0.30} \pm \textbf{0.01} & 1.06 \pm 0.15 \\ \text{PF} & 2.80 \pm 0.04 & 3.12 \pm 0.06 \\ \text{PFF} & 0.60 \pm 0.02 & 1.00 \pm 0.05 \\ \text{KNet} & 0.60 \pm 0.02 & 1.81 \pm 0.05 \\ \text{VRNN} & 3.67 \pm 0.06 & 3.67 \pm 0.06 \\ \text{SRNN} & 3.08 \pm 0.56 & 3.63 \pm 0.05 \\ \text{DKF} & 3.70 & \text{NA} \\ \end{array} $	$\begin{tabular}{ c c c c c c c } \hline & & & & & & & & & & & & & & & & & & $	$\begin{tabular}{ c c c c c c c c c c c c c c c c c c c$	$ \begin{array}{c c c c c c c c c c c c c c c c c c c $			

475

476 Table 2 presents the assimilation performance across different noise levels and observation settings. 477 DBF outperforms existing methods in direct observations with  $\sigma = 3, 5$ , and across all noise levels for nonlinear observation cases. In the  $\sigma = 1$  setting with direct observation, traditional algorithms 478 like EnKF and ETKF outperform DBF. 479

480 The superior performance of EnKF and ETKF with direct observations at the lowest noise level 481 can be attributed to the minimal non-Gaussianity in the posteriors within physical space. Non-482 Gaussianity can originate from both the dynamics model (predict step) and the observation model (update step). In this setting, the linearity of the observation operator prevents non-Gaussianity from 483 being introduced during the update step, provided that the prior  $q(z_t|o_{1:t-1})$  is Gaussian. Addition-484 ally, state estimation from each observation is highly accurate due to small noise. As a result, the 485 prior  $q(z_t|o_{1:t-1})$  remains close to a Gaussian distribution, as the locally linear approximation of



Figure 5: RMSE results for Lorenz96 experiments. Panels (a), (b) show results for direct observation with  $\sigma = 1$  and  $\sigma = 5$ . Panel (c) shows results for nonlinear observation with  $\sigma = 1$ .

the dynamics adequately captures the time evolution of probability distributions. The poorer performance of EnKF and ETKF in the  $\sigma = 5$  experiment is attributed to the increased non-Gaussianity introduced during each predict step. Similarly, when the observation operator is nonlinear, each update step introduces substantial non-Gaussianity. This results in a significant drop in performance for traditional filtering methods across all noise levels. In these scenarios, DBF consistently maintains an advantage over classical DA algorithms.

504 We observe that training DVAE-based methods is 505 highly unstable, while that for DBF exhibits stability. 506 Dynamics in DVAEs are modeled by RNNs, which 507 often suffer from unstable training due to exploding 508 or vanishing gradients. In contrast, DBF employs 509 matrix multiplication for dynamics. If the eigenvalues of the matrix exceed one by a large margin, the 510 model predictions, and consequently the loss func-511 tion, would explode irrespective of inputs. Fig. 6 512 shows the histogram of the absolute values of eigen-513 values at the end of training, which are distributed 514 around or below one, indicating stable training. 515

	setting	max[abs(eig)]
	D, $\sigma = 1$	$1.016\pm0.002$
	D, $\sigma = 3$	$1.014\pm0.002$
	D, $\sigma = 5$	$1.011\pm0.001$
	N, $\sigma = 1$	$1.012\pm0.003$
	N, $\sigma = 3$	$1.008\pm0.004$
0.8 0.9 1.0 1.1 abs(eigenvalue)	N, $\sigma = 5$	$1.004\pm0.001$

Figure 6: Histogram of 800 eigenvalues of the dynamics matrix in Lorenz96. D for direct and N for nonlinear observations.

#### 4 LIMITATION

519 DBF's learning of IOO requires a training phase, unlike classical model-based data assimilation 520 methods. Specifically, when dealing with nonlinear dynamics, DBF requires either: (i) a pair of 521  $(z_t, o_t)$  generated from the original SSM, (ii) a pair of  $(z_t, o_t)$  obtained via, e.g., retrospective re-522 analysis (ERA5; Hersbach et al. 2020 in weather forecasting), or (iii) a pretrained Koopman operator 523 and observed data  $o_t$ .

In the Lorenz96 experiment, DBF's performance with direct observation with  $\sigma = 1$  falls short compared to EnKF and ETKF. In this setting, the non-Gaussianity of posteriors is weak, resulting in minor approximation errors due to Gaussian assumptions. Consequently, a model-based approach may be more advantageous in such situations, as it leverages complete SSM knowledge without introducing training biases.

529 530

531

516

517 518

495

496 497

#### 5 CONCLUSION

532 We propose DBF, a novel DA method. DBF is a NN-based extension of the KF designed to handle nonlinear observations. While constraining the test distributions to remain Gaussian, DBF enhances 533 their representational capacity by leveraging nonlinear transform expressed by a NN. DBF is the first 534 "Bayes-Faithful" amortized variational inference methodology, constructing test distributions that 535 mirror the inference structure of a SSM with the Markov property. This structured inference enables 536 analytical computation of test distributions, preventing the accumulation of Monte Carlo sampling 537 errors over time steps. DBF exhibits superior performance over existing methods in scenarios where 538 posterior distributions become highly non-Gaussian, such as in the presence of nonlinear observation operators or significant observation noise.

Reproducibility Statement We have provided the source code to reproduce the experiments for
double pendulum (Sec. 3.2) and the Lorenz96 (Sec. 3.3) in the supplementary material. The hyperparameters for the training are provided in Table 3 in the appendix. Generation method of the
training and test dataset, the dynamics model, the observation model, and the architectures are detailed in the appendix: Sec. C.1, C.2, and C.3.

546 REFERENCES

- Miguel Alfonzo and Dean S. Oliver. Seismic data assimilation with an imperfect model. *Computational Geosciences*, 24(2):889–905, 2020. Marine Environmental Monitoring and Prediction.
- Marcin Andrychowicz, Lasse Espeholt, Di Li, Samier Merchant, Alexander Merose, Fred Zyda, Shreya Agrawal, and Nal Kalchbrenner. Deep learning for day forecasts from sparse observations. ArXiv, abs/2306.06079, 2023. URL https://api.semanticscholar.org/
   CorpusID:259129311.
- Md Abdul Awal, Md Abu Rumman Refat, Feroza Naznin, and Md Zahidul Islam. A particle filter based visual object tracking: A systematic review of current trends and research challenges. *International Journal of Advanced Computer Science and Applications*, 14(11), 2023. doi: 10.14569/IJACSA.2023.01411131. URL http://dx.doi.org/10.14569/IJACSA.2023.01411131.
- Omri Azencot, N. Benjamin Erichson, Vanessa Lin, and Michael Mahoney. Forecasting sequential data using consistent koopman autoencoders. In Hal Daumé III and Aarti Singh (eds.), *Proceedings of the 37th International Conference on Machine Learning*, volume 119 of *Proceedings of Machine Learning Research*, pp. 475–485. PMLR, 13–18 Jul 2020. URL https://proceedings.mlr.press/v119/azencot20a.html.
- Eviatar Bach and Michael Ghil. A multi-model ensemble kalman filter for data assimilation and forecasting. Journal of Advances in Modeling Earth Systems, 15(1):e2022MS003123, 2023.
   doi: https://doi.org/10.1029/2022MS003123. URL https://agupubs.onlinelibrary.wiley.com/doi/abs/10.1029/2022MS003123. e2022MS003123 2022MS003123.
- Alexandros Beskos, Dan Crisan, and Ajay Jasra. On the stability of sequential Monte Carlo methods in high dimensions. *The Annals of Applied Probability*, 24(4):1396 – 1445, 2014. doi: 10.1214/ 13-AAP951. URL https://doi.org/10.1214/13-AAP951.
- 572
   573
   574
   574
   574
   575
   575
   576
   576
   576
   576
   576
   576
   576
   576
   576
   576
   577
   576
   577
   577
   576
   577
   577
   576
   577
   577
   577
   578
   579
   579
   570
   570
   571
   571
   572
   572
   573
   574
   575
   576
   576
   577
   577
   578
   579
   579
   570
   570
   571
   571
   572
   572
   573
   574
   574
   575
   575
   576
   577
   576
   577
   577
   578
   579
   579
   579
   570
   570
   571
   571
   572
   572
   572
   573
   574
   574
   575
   575
   576
   577
   578
   578
   579
   579
   570
   570
   571
   571
   572
   572
   573
   574
   574
   575
   575
   576
   576
   576
   577
   578
   579
   579
   570
   570
   571
   571
   572
   572
   574
   574
   575
   575
   576
   576
   576
   577
   578
   578
   579
- Yuming Chen, Daniel Sanz-Alonso, and Rebecca Willett. Autodifferentiable ensemble kalman filters. *SIAM Journal on Mathematics of Data Science*, 4(2):801–833, 2022. doi: 10.1137/21M1434477. URL https://doi.org/10.1137/21M1434477.
- Yuming Chen, Daniel Sanz-Alonso, and Rebecca Willett. Reduced-order autodifferentiable ensemble kalman filters. *Inverse Problems*, 39(12), 10 2023. doi: 10.1088/1361-6420/acff14.
- 583
   584
   585
   Nicolas Chopin and Omiros Papaspiliopoulos. An introduction to Sequential Monte Carlo. Springer series in statistics. Springer, 2020. URL https://ci.nii.ac.jp/ncid/BC03234800.
- Junyoung Chung, Kyle Kastner, Laurent Dinh, Kratarth Goel, Aaron C Courville, and Yoshua Bengio. A Recurrent Latent Variable Model for Sequential Data. In Advances in Neural Information Processing Systems, volume 28. Curran Associates, Inc., 2015. URL https://proceedings.neurips.cc/paper\_files/paper/2015/ hash/b618c3210e934362ac261db280128c22-Abstract.html.
- Fred Daum and Jim Huang. Nonlinear filters with log-homotopy. In Oliver E. Drummond and Richard D. Teichgraeber (eds.), *Signal and Data Processing of Small Targets 2007*, volume 6699, pp. 669918. International Society for Optics and Photonics, SPIE, 2007. doi: 10.1117/12.725684. URL https://doi.org/10.1117/12.725684.

631

- 594 Emmanuel de Bézenac, Syama Sundar Rangapuram, Konstantinos Benidis, Michael Bohlke-595 Schneider, Richard Kurle, Lorenzo Stella, Hilaf Hasson, Patrick Gallinari, and Tim 596 Januschowski. Normalizing kalman filters for multivariate time series analysis. In 597 H. Larochelle, M. Ranzato, R. Hadsell, M.F. Balcan, and H. Lin (eds.), Advances in Neu-598 ral Information Processing Systems, volume 33, pp. 2995-3007. Curran Associates, Inc., 2020. URL https://proceedings.neurips.cc/paper\_files/paper/2020/ file/1f47cef5e38c952f94c5d61726027439-Paper.pdf. 600
- 601 Geir Evensen. Sequential data assimilation with a nonlinear quasi-geostrophic model us-602 ing Monte Carlo methods to forecast error statistics. Journal of Geophysical Research: 603 Oceans, 99(C5):10143-10162, 1994. ISSN 2156-2202. doi: 10.1029/94JC00572. URL 604 https://onlinelibrary.wiley.com/doi/abs/10.1029/94JC00572. \_eprint: 605 https://onlinelibrary.wiley.com/doi/pdf/10.1029/94JC00572.
- Heng Fan, Liting Lin, Fan Yang, Peng Chu, Ge Deng, Sijia Yu, Hexin Bai, Yong Xu, Chunyuan 607 Liao, and Haibin Ling. Lasot: A high-quality benchmark for large-scale single object tracking, 608 2019. 609
- 610 Heng Fan, Hexin Bai, Liting Lin, Fan Yang, Peng Chu, Ge Deng, Sijia Yu, Harshit, Mingzhen 611 Huang, Juehuan Liu, Yong Xu, Chunyuan Liao, Lin Yuan, and Haibin Ling. Lasot: A high-612 quality large-scale single object tracking benchmark. Int. J. Comput. Vision, 129(2):439-461, feb 613 2021. ISSN 0920-5691. doi: 10.1007/s11263-020-01387-y. URL https://doi.org/10. 614 1007/s11263-020-01387-y.
- 615 Marco Fraccaro, Søren Kaae Sønderby, Ulrich Paquet, and Ole Winther. Sequential neural models 616 with stochastic layers. In Proceedings of the 30th International Conference on Neural Information 617 Processing Systems, NIPS'16, pp. 2207–2215, Red Hook, NY, USA, 2016. Curran Associates Inc. 618 ISBN 978-1-5108-3881-9. 619
- 620 Marco Fraccaro, Simon Kamronn, Ulrich Paquet, and Ole Winther. A disentangled recog-621 nition and nonlinear dynamics model for unsupervised learning. In I. Guyon, U. Von 622 Luxburg, S. Bengio, H. Wallach, R. Fergus, S. Vishwanathan, and R. Garnett (eds.), Advances in Neural Information Processing Systems, volume 30. Curran Associates, Inc., 623 URL https://proceedings.neurips.cc/paper\_files/paper/2017/ 2017. 624 file/7b7a53e239400a13bd6be6c91c4f6c4e-Paper.pdf. 625
- 626 Thomas Frerix, Dmitrii Kochkov, Jamie Smith, Daniel Cremers, Michael Brenner, and Stephan 627 Hoyer. Variational Data Assimilation with a Learned Inverse Observation Operator. In Proceed-628 ings of the 38th International Conference on Machine Learning, pp. 3449–3458. PMLR, July 629 2021. URL https://proceedings.mlr.press/v139/frerix21a.html. ISSN: 630 2640-3498.
- Laurent Girin, Simon Leglaive, Xiaoyu Bie, Julien Diard, Thomas Hueber, and Xavier Alameda-632 Pineda. Dynamical variational autoencoders: A comprehensive review. Foundations and Trends® 633 in Machine Learning, 15(1-2):1–175, 2021. ISSN 1935-8237. doi: 10.1561/2200000089. URL 634 http://dx.doi.org/10.1561/220000089. 635
- 636 Albert Gu and Tri Dao. Mamba: Linear-time sequence modeling with selective state spaces. arXiv 637 preprint arXiv:2312.00752, 2023.
- 638 Albert Gu, Karan Goel, and Christopher Re. Efficiently modeling long sequences with structured 639 state spaces. In International Conference on Learning Representations, 2022. URL https: 640 //openreview.net/forum?id=uYLFoz1vlAC.
- 642 Hans Hersbach, Bill Bell, Paul Berrisford, Shoji Hirahara, András Horányi, Joaquín Muñoz-643 Sabater, Julien Nicolas, Carole Peubey, Raluca Radu, Dinand Schepers, Adrian Simmons, Cor-644 nel Soci, Saleh Abdalla, Xavier Abellan, Gianpaolo Balsamo, Peter Bechtold, Gionata Biavati, Jean Bidlot, Massimo Bonavita, Giovanna De Chiara, Per Dahlgren, Dick Dee, Michail Dia-645 mantakis, Rossana Dragani, Johannes Flemming, Richard Forbes, Manuel Fuentes, Alan Geer, 646 Leo Haimberger, Sean Healy, Robin J. Hogan, Elías Hólm, Marta Janisková, Sarah Keeley, 647 Patrick Laloyaux, Philippe Lopez, Cristina Lupu, Gabor Radnoti, Patricia de Rosnay, Iryna

663

664

665

666

667

687

688 689

690

691

692

 Rozum, Freja Vamborg, Sebastien Villaume, and Jean-Noël Thépaut. The era5 global reanalysis. *Quarterly Journal of the Royal Meteorological Society*, 146(730):1999–2049, 2020. doi: https://doi.org/10.1002/qj.3803. URL https://rmets.onlinelibrary.wiley.com/ doi/abs/10.1002/qj.3803.

- Chih-Chi Hu and Peter Jan van Leeuwen. A particle flow filter for high-dimensional system applications. *Quarterly Journal of the Royal Meteorological Society*, 147(737):2352–2374, 2021. doi: https://doi.org/10.1002/qj.4028. URL https://rmets.onlinelibrary.wiley.com/doi/abs/10.1002/qj.4028.
- Brian R. Hunt, Eric J. Kostelich, and Istvan Szunyogh. Efficient data assimilation for spatiotemporal chaos: A local ensemble transform kalman filter. *Physica D: Nonlinear Phenomena*, 230(1):112–126, 2007. ISSN 0167-2789. doi: https://doi.org/10.1016/j.physd.2006.11.008. URL https://www.sciencedirect.com/science/article/pii/S0167278906004647. Data Assimilation.
  - Glenn Jocher, Ayush Chaurasia, and Jing Qiu. Ultralytics YOLO, January 2023. URL https://github.com/ultralytics/ultralytics.
  - B. O. Koopman. Hamiltonian systems and transformation in hilbert space. *Proceedings of the National Academy of Sciences*, 17(5):315–318, 1931. doi: 10.1073/pnas.17.5.315. URL https://www.pnas.org/doi/abs/10.1073/pnas.17.5.315.
- Rahul G. Krishnan, Uri Shalit, and David Sontag. Deep kalman filters, 2015. URL https://arxiv.org/abs/1511.05121.
- Rahul G. Krishnan, Uri Shalit, and David Sontag. Structured inference networks for nonlinear state space models, 2016. URL https://arxiv.org/abs/1609.09869.
- J. Larsen, J.L. Høyer, and J. She. Validation of a hybrid optimal interpolation and kalman filter
  scheme for sea surface temperature assimilation. *Journal of Marine Systems*, 65(1):122–133,
  2007. ISSN 0924-7963. doi: https://doi.org/10.1016/j.jmarsys.2005.09.013. URL https:
  //www.sciencedirect.com/science/article/pii/S0924796306002880. Marine Environmental Monitoring and Prediction.
- Phillip Lippe, Bastiaan S. Veeling, Paris Perdikaris, Richard E Turner, and Johannes Brandstetter. PDE-Refiner: Achieving Accurate Long Rollouts with Temporal Neural PDE Solvers. In *Thirty-seventh Conference on Neural Information Processing Systems*, 2023. URL https: //openreview.net/forum?id=Qv646811WS.
- Andrew C. Lorenc. Modelling of error covariances by 4d-var data assimilation. *Quarterly Journal of the Royal Meteorological Society*, 129(595):3167–3182, 2003. doi: https://doi.org/10.1256/qj.02.
   URL https://rmets.onlinelibrary.wiley.com/doi/abs/10.1256/qj. 02.131.
  - E.N. Lorenz. *Predictability: a problem partly solved*. PhD thesis, Shinfield Park, Reading, 1995 1995.
  - Bethany Lusch, J. Nathan Kutz, and Steven L. Brunton. Deep learning for universal linear embeddings of nonlinear dynamics. *Nature Communications*, 9:4950, November 2018. doi: 10.1038/s41467-018-07210-0.
- Matthias Müller, Adel Bibi, Silvio Giancola, Salman Alsubaihi, and Bernard Ghanem. Trackingnet: A large-scale dataset and benchmark for object tracking in the wild. In *Computer Vision – ECCV 2018: 15th European Conference, Munich, Germany, September 8-14, 2018, Proceedings, Part I*, pp. 310–327, Berlin, Heidelberg, 2018. Springer-Verlag. ISBN 978-3-030-01245-8. doi: 10.1007/978-3-030-01246-5\_19. URL https://doi.org/10.1007/ 978-3-030-01246-5\_19.
- Shun Ohishi, Takemasa Miyoshi, Takafusa Ando, Tomohiko Higashiuwatoko, Eri Yoshizawa, Hiroshi Murakami, and Misako Kachi. Letkf-based ocean research analysis (lora) version 1.0.
   *Geoscience Data Journal*, n/a(n/a), 2024. doi: https://doi.org/10.1002/gdj3.271. URL https://rmets.onlinelibrary.wiley.com/doi/abs/10.1002/gdj3.271.

Guy Revach, Nir Shlezinger, Xiaoyong Ni, Adrià López Escoriza, Ruud J. G. van Sloun, and Yon-ina C. Eldar. Kalmannet: Neural network aided kalman filtering for partially known dynamics. *IEEE Transactions on Signal Processing*, 70:1532–1547, 2022. doi: 10.1109/TSP.2022.3158588. Naoya Takeishi, Yoshinobu Kawahara, and Takehisa Yairi. Learning koopman invariant subspaces for dynamic mode decomposition. In I. Guyon, U. Von Luxburg, S. Bengio, H. Wallach, R. Fergus, S. Vishwanathan, and R. Garnett (eds.), Ad-vances in Neural Information Processing Systems, volume 30. Curran Associates, Inc., 2017. URL https://proceedings.neurips.cc/paper\_files/paper/2017/ file/3a835d3215755c435ef4fe9965a3f2a0-Paper.pdf. Liang Xu and Ruixin Niu. Ekfnet: Learning system noise covariance parameters for nonlinear tracking. IEEE Transactions on Signal Processing, 72:3139–3152, 2024. doi: 10.1109/TSP.2024. 3417350. Linqi Zhou, Michael Poli, Winnie Xu, Stefano Massaroli, and Stefano Ermon. Deep latent state space models for time-series generation. In Andreas Krause, Emma Brunskill, Kyunghyun Cho, Barbara Engelhardt, Sivan Sabato, and Jonathan Scarlett (eds.), International Confer-ence on Machine Learning, ICML 2023, 23-29 July 2023, Honolulu, Hawaii, USA, volume 202 of Proceedings of Machine Learning Research, pp. 42625-42643. PMLR, 2023. URL https://proceedings.mlr.press/v202/zhou23i.html. DERIVATION OF THE EVIDENCE LOWER-BOUND AND THE ASSOCIATED Α **MONTE-CARLO SAMPLING** A.1 LINEAR DYNAMICS CASE Following the definition of the probability density,  $p(o_t, h_t | o_{1:t-1}) = p(o_t | o_{1:t-1})p(h_t | o_{1:t})$ (11)

756 Using Eq. 11 at the third equality, 

 $\log p(o_{1:T}) = \sum_{t=1}^{T} \log p(o_t | o_{1:t-1})$  $= \sum_{t=1}^{I} \int q(h_t|o_{1:t}) \log p(o_t|o_{1:t-1}) dh_t$  $= \sum_{t=1}^{T} \int q(h_t|o_{1:t}) \log \frac{p(o_t, h_t|o_{1:t-1})}{p(h_t|o_{1:t})} dh_t$  $= \sum_{t=1}^{T} \int q(h_t|o_{1:t}) \log \left[ \frac{p(o_t, h_t|o_{1:t-1})}{q(h_t|o_{1:t})} \frac{q(h_t|o_{1:t})}{p(h_t|o_{1:t})} \right] dh_t$  $= \sum_{i=1}^{T} \int q(h_t|o_{1:t}) \log \left[ \frac{p(o_t, h_t|o_{1:t-1})}{q(h_t|o_{1:t})} \right] dh_t + KL[q(h_t|o_{1:t})||p(h_t|o_{1:t})]$  $= \sum_{t=1}^{T} \mathcal{L}_{ELBO,t} + KL[q(h_t|o_{1:t})||p(h_t|o_{1:t})]$  $\geq \sum^{1} \mathcal{L}_{ELBO,t}$ (12) $\mathcal{L}_{ELBO,t} = \int q(h_t|o_{1:t}) \log \left[\frac{p(o_t, h_t|o_{1:t-1})}{q(h_t|o_{1:t})}\right] dh_t$  $= \int q(h_t|o_{1:t}) \log \left[ \frac{p(h_t|o_{1:t-1})p(o_t|h_t)}{q(h_t|o_{1:t})} \right] dh_t$  $= \int q(h_t|o_{1:t}) \log p(o_t|h_t) dh_t + \int q(h_t|o_{1:t}) \frac{p(h_t|o_{1:t-1})}{q(h_t|o_{1:t})} dh_t$  $= \int q(h_t|o_{1:t}) \log p(o_t|h_t) dh_t - KL[q(h_t|o_{1:t})|p(h_t|o_{1:t-1})]$ (13)

The true prior at step t  $(p(h_t|o_{1:t-1}))$  on the right hand side of Eq. 13 could be replaced with the prior computed from the test distribution  $q(h_t|o_{1:t-1})$  when training.

#### A.2 NONLINEAR DYNAMICS CASE

 $p(o_t, z_t, h_t | o_{1:t-1}, z_{1:t-1}) = p(o_t, z_t | o_{1:t-1}, z_{1:t-1}) p(h_t | o_{1:t}, z_{1:t})$ (14)

The derivation proceeds parallel to the linear case. Using Eq. 14 at the third equality,  $\log p(o_{1:T}, z_{1:T}) = \sum \log p(o_t, z_t | o_{1:t-1}, z_{1:t-1})$  $= \sum_{i=1}^{T} \int q(h_t|o_{1:t}) \log p(o_t, z_t|o_{1:t-1}, z_{1:t-1}) dh_t$  $= \sum_{t=1}^{T} \int q(h_t|o_{1:t}) \log \frac{p(o_t, z_t, h_t|o_{1:t-1}, z_{1:t-1})}{p(h_t|o_{1:t}, z_{1:t})} dh_t$  $= \sum_{t=1}^{T} \int q(h_t|o_{1:t}) \log \left[ \frac{p(o_t, z_t, h_t|o_{1:t-1}, z_{1:t-1})}{q(h_t|o_{1:t})} \frac{q(h_t|o_{1:t})}{p(h_t|o_{1:t}, z_{1:t})} \right] dh_t$  $= \sum_{t=1}^{T} \int q(h_t|o_{1:t}) \log \left[ \frac{p(o_t, z_t, h_t|o_{1:t-1}, z_{1:t-1})}{q(h_t|o_{1:t})} \right] dh_t + KL[q(h_t|o_{1:t})||p(h_t|o_{1:t}, z_{1:t})]$  $= \sum_{i=1}^{T} \mathcal{L}_{ELBO,joint,t} + KL[q(h_t|o_{1:t})||p(h_t|o_{1:t}, z_{1:t})]$  $\geq \sum \mathcal{L}_{ELBO, joint, t}$ (15) $\mathcal{L}_{ELBO,joint,t} = \int q(h_t|o_{1:t}) \log \left[ \frac{p(o_t, z_t, h_t|o_{1:t-1}, z_{1:t-1})}{q(h_t|o_{1:t})} \right] dh_t$  $= \int q(h_t|o_{1:t}) \log \left[ \frac{p(h_t|o_{1:t-1}, z_{1:t-1})p(o_t, z_t|h_t)}{q(h_t|o_{1:t})} \right] dh_t$  $= \int q(h_t|o_{1:t})[\log p(z_t|h_t) + \log p(o_t|z_t)]dh_t + \int q(h_t|o_{1:t})\frac{p(h_t|o_{1:t-1}, z_{1:t-1})}{q(h_t|o_{1:t})}dh_t$  $= \int q(h_t|o_{1:t}) \log p(z_t|h_t) dh_t - KL[q(h_t|o_{1:t})|p(h_t|o_{1:t-1}, z_{1:t-1})] + \log p(o_t|z_t)$ (16) 

The true prior at step t  $(p(h_t|o_{1:t-1}, z_{1:t-1}))$  on the right hand side of Eq. 16 could be replaced with the prior computed from the test distribution  $q(h_t|o_{1:t-1})$  when training. The last term of the equation  $(\log p(o_t|z_t))$  can be neglected as it does not affect the new latent variables  $h_t$ .

#### A.3 COMPARISON TO OTHER DVAES IN TERMS OF MONTE-CARLO SAMPLING

The crucial difference from other DVAEs is that the Monte-Carlo samplings in DBF are not nested with each other. In DVAE, we need to evaluate an integral term  $\int q(h_{1:T}|o_{1:T}) \log p(o_{1:T}, h_{1:T}) dh_{1:T}$ , where  $q(h_{1:T}|o_{1:T}) = \prod_t q(h_t|h_{t-1}, o_t)$ . Although the log-term could be factorized as  $\sum_t \log p(o_t|h_t) + \log p(h_t|h_{t-1})$  thanks to the Markov property, we need MC (nested) sequential sampling over  $h_{1:T}$  if we want to evaluate the term at t = T. On the other hand, ELBO in DBF is  $\sum_t \int q(h_t|o_{1:t}) \log p(z_t|h_t) dh_t + KL[q(h_t|o_{1:t})|q(h_t|o_{1:t-1})]$  because DBF takes the lower limit of  $\sum_t \log p(o_t|o_{1:t-1})$ . Thanks to the analytic expressions of  $q(h_t|o_{1:t})$  and  $q(h_t|o_{1:t-1})$ , the KL term can be computed analytically. A MC sampling is needed to compute  $\int q(h_t|o_{1:t}) \log p(z_t|h_t) dh_t$  but this is independent from other timesteps.

#### B ADDITIONAL LINEAR DYNAMICS EXPERIMENT: TWO-BODY MOVING MNIST

This experiment demonstrates DBF's ability to handle linear dynamics where key parameters of the observation operator are unknown. The dataset consists of 2D figures containing two embedded images, each moving at a constant speed and bouncing off frame edges. The system's physical state is described by eight variables: the positions (x, y) and velocities  $(v_x, v_y)$  of the two embedded images. The dynamics matrix is block-diagonal, composed of four (two-body times two dimensions) translation matrices,  $A_{tr}$ :  $A_{tr} = \begin{pmatrix} 1 & 1 \\ 0 & 1 \end{pmatrix}$ ,  $z_t = \begin{pmatrix} x_t \\ v_{x_t} \end{pmatrix}$ . Observations are corrupted by additive Gaussian noise with a standard deviation of  $\sigma = 50$  per pixel, where the original pixel values range from 0 to 255 (see panel (a) of Fig. 7 for an example of the data provided).

870 The aim is to show that DBF can track the linear dynamics while estimating unknown system pa-871 rameters. DBF learns the pixel values of the embedded images from noisy observations, while 872 maintaining consistency with physical motion. The observation model contains 1,568 unknown pa-873 rameters, corresponding to the number of pixels in the images. In classical DA algorithms, it is not 874 possible to train unknown system parameters. However, it may be possible to infer these parameters 875 by incorporating them as new physical dimensions. We have adopted this strategy for classical DA 876 algorithms (EnKF, ETKF, and PF). For these, we tested at three different model noise levels ( $\sigma_{sys}$ of 1, 0.1, and 0.01) and chosen the best parameter. While DVAE generates latent variables, they are 877 different from the state variables of the original SSM: therefore, they cannot infer the position or 878 velocity from those images. We were unable to compare with KalmanNet as the high observation 879 dimensions of  $x_{dim}^2 = (44 \times 44)^2$  inhibits the training even with the batch size of one. 880

Fig. 7 summarizes the experiment. Panel (a) shows an example from the test set, illustrating the 882 challenges posed by strong noise and overlapping images. Panel (b) presents the DBF learning pro-883 cess. In the rightmost table, we compare the success rates of DBF against model-based approaches (EnKF, ETKF, PF). We define success as achieving a root-mean-square error (RMSE) of less than 884 1.0 for both position  $(x_1, y_1, x_2, y_2)$  and velocity  $(v_{x_1}, v_{y_1}, v_{x_2}, v_{y_2})$  of the two digits over the fi-885 nal ten steps. DBF successfully performs assimilation without explicit knowledge of the images, while all the other model-based approaches fail. The KF-inspired approaches (EnKF, ETKF) failed 887 because of very strong non-Gaussianity in the observation process and the high system dimension. 888 Similarly, PF underperformed because the number of particles (10,000) was insufficient for the prob-889 lem dimension ( $z_{\rm dim} = 8$  and two digits images  $2 \times 28 \times 28 = 1,568$ ). Figures for visualizing the 890 assimilation results for all the algorithms are given in the appendix (Fig. 13). 891

Panel (b) of Fig. 7 illustrates the evolution of the estimated figures. Initially, DBF assumes two random shapes. "Iterations" in panel (b) means the number of parameter update steps the DBF has undergone. As training progresses, it first identifies one of the numbers ("9") and subsequently detects the second shape ("5"). By the end of the training process, DBF nearly perfectly estimates the parameters of the observation model, including the positions of the figures, which is crucial for adjusting their reflective behavior.



Figure 7: Figures from the two-body Moving MNIST experiments. Panel (a) displays examples of the observation data. Panel (b) illustrates the evolution of the observation model parameters (the embedded images) during training. The table compares the success rates of four methodologies.

#### C SETTINGS AND ADDITIONAL RESULTS FOR EXPERIMENTS

ŀ

parametrization of the dynamics matrix We have parametrized the dynamics matrix A following Lusch et al. (2018): we consider that  $h_{dim}/2$  complex eigenvalues  $\lambda_i (0 \le i < h_{dim}/2)$  characterize A. Namely, A is a block-diagonal matrix of  $h_{dim}/2$  blocks. Each block consists of  $2 \times 2$ matrix, whose components are:

917

892

893

894

895

896

897

899

900

901

902

903 904 905

906

907

908 909 910

911

$$A_{block} = \exp(\rho_i) \begin{pmatrix} \cos(\omega_i) & -\sin(\omega_i) \\ \sin(\omega_i) & \cos(\omega_i) \end{pmatrix}, \tag{17}$$

where  $\rho_i = \operatorname{Re}[\lambda_i]$  and  $\omega_i = \operatorname{Im}[\lambda_i]$ . In contrast to Lusch et al. (2018), we apply the same dynamics matrix at any positions on the latent space. We consider that this representation is sufficiently expressive, as it can express any matrix on a complex number field that is diagonalizable.

One key advantage of DBF is that augmenting the latent dimension only results in a linear increase in computational demand. This scaling is due to the efficient parametrization of the dynamics matrix, where the block-diagonal structure allows operations to scale linearly with the latent dimension. In contrast, methods such as Sequential Monte Carlo (SMC) suffer from exponential increases in computational demand as the latent space grows, assuming that the same density of particles must be maintained to capture posterior distributions. This makes DBF particularly well-suited for highdimensional systems where traditional methods struggle with computational complexity.

929 Computational resources We conduct experiments on a cluster of V100 GPUs. Each GPU has930 memory of 32GB.

**hyperparameters for training** For all experiments, we have used Adam optimizer with default parameters. Table 3 shows hyperparameters employed in our experiments. Trainings for moving MNIST and double pendulum are conducted with one GPU, while that for Lorenz96 is with eight GPUs.

Table	3.	Hyper	narameters	for	training
raute	э.	riyper	parameters	101	uaming

	lr	batch size	$h_{dim}$	$N_{data,train}$	Epochs	train time per model
object tracking	-	-	8	-	-	-
double pendulum	$10^{-3}$	256	50	$1.0 \times 10^7$	1	6hr× 1GPU
Lorenz96	$3 \times 10^{-3}$	64	800	$2.6  imes 10^7$	1	15hr× 8GPUs
moving MNIST	$10^{-3}$	64	8	480,000	2	$3hr \times 1GPU$

#### C.1 OBJECT TRACKING

**Dataset:** "Airplane" movies in the LaSOT dataset (Fan et al., 2019; 2021). It contains 20 movies. Each movie has at least 1,000 frames. We chop the first 1,000 frames into 20 sets of 50 frames. Airplanes numbered one to ten are considered a validation set used to determine the model hyperparameters. We use the remaining data (airplane-11 to airplane-20) as a test set to evaluate the performance of the filters.

**Dynamics model:** Constant velocity model. The (x, y) coordinates and  $(v_x, v_y)$  velocities of the top left and bottom right edges are the latent (physical) variables.

$$h_{t+1} = Fh_t \tag{18}$$

$$F = \begin{pmatrix} 1 & 0 & 0 & 0 & dt & 0 & 0 & 0 \\ 0 & 1 & 0 & 0 & 0 & dt & 0 & 0 \\ 0 & 0 & 1 & 0 & 0 & 0 & dt & 0 \\ 0 & 0 & 0 & 1 & 0 & 0 & 0 & dt \\ 0 & 0 & 0 & 0 & 1 & 0 & 0 & 0 \\ 0 & 0 & 0 & 0 & 0 & 1 & 0 & 0 \\ 0 & 0 & 0 & 0 & 0 & 0 & 1 & 0 \\ 0 & 0 & 0 & 0 & 0 & 0 & 0 & 1 \end{pmatrix}, h_t = \begin{pmatrix} x_{1,t} \\ y_{1,t} \\ x_{2,t} \\ y_{2,t} \\ v_{y_{1,t}} \\ v_{y_{1,t}} \\ v_{y_{2,t}} \end{pmatrix}.$$
(19)

970 Here,  $x_{1,t}$  and  $y_{1,t}$  stand for the coordinates of the left top edge of the bounding box, and  $x_{2,t}$  and 971  $y_{2,t}$  are the right bottom edge of the box.  $v_{x_{1,t}}, v_{y_{1,t}}, v_{x_{2,t}}, v_{y_{2,t}}$  are velocities of box edges. dt is the time difference between frames, which we take as 1 (arbitrary).

973

984 985

986

987

988 989

990

991

992

993 994

995

996

997

998 999

974 Network architecture: We use a pre-trained detector YOLOv8n model (Jocher et al., 2023). The 975 detector yields the bounding box's position, X, and the box's confidence score, c. We set the de-976 tection threshold at 0.01. In cases where the detector reports multiple bounding boxes, we choose 977 the one with the highest posterior probability. We use the bounding box coordinates as  $f_{\theta}(o_t) = X$ . 978 Several choices for the relation between confidence score and  $G_{\theta}(o_t)$  are possible.

We experiment with linear confidence  $G_{\theta}(o_t) \propto c$  and squared confidence  $G_{\theta}(o_t) \propto c^2$ . We determine the system noise factor for either dependence with the validation set. We use normalized precision as the evaluation metric (Müller et al., 2018). Figure 8 shows the normalized precision score for the validation set for the system noise factor. The system noise factor of  $10^{-1}$  is chosen for KF. For DBF, squared confidence with the system noise factor of  $10^{-2}$  is employed.



Figure 8: Normalized precision scores for validation samples.

# C.2 DOUBLE PENDULUM

**Dataset:** The dataset consists of 2D coordinates representing the positions of two weights. The training set includes 10, 240, 000 initial conditions, while the test set contains 10 initial conditions. The number of training samples is sufficiently large to ensure that the training converges. During DVAE training, we observed that some initial conditions resulted in training failure due to instability; however, we maintained the total number of training samples since the training was successful for at least one initial condition. Both datasets comprise 80 time steps. Numerical integration is performed using the solve\_ivp function in SciPy, with relative tolerance (rtol) set to  $10^{-2}$  and absolute tolerance (atol) set to  $10^{-2}$ .

1010 A schematic figure explaining the problem setting is presented in panel (a) of Fig. 3 in the main text.

1011 Dynamics model is described in https://matplotlib.org/stable/gallery/animation/double\_pendulum.html. 1012 The length of the bars is 1 [m], and the positions of the two pendulum weights are observable with 1013 Gaussian noise of  $\sigma = 0.1, 0.3$ , or 0.5 [m]. The observation interval is 0.03 [s]. The task is to 1014 predict the positions of the two weights in the successive ten frames.

1015

1018

1019

1020 1021

1016 Network architecture:  $f_{\theta}$ : A sequence of ten "linear blocks" composed of fully connected layers, 1017 layer normalizations, and skip connections. Namely, each linear block has three components:

- fc: (input dimension) × (output dimension) linear layer,
- norm: layer normalization,
- skip: skip connection.

Taking four observation variables as input, the first linear block expands the dimensionality to 100. The intermediate linear blocks maintain these 100-dimensional variables. The final linear block reduces the 100-dimensional input to a 50-dimensional output, representing 50 latent space variables. The ReLU activation function is applied throughout the network. The structure of  $G_{\theta}$  mirrors that



Table 4: List of hyperparameters for double pendulum experiment.

Figure 9: PF results with 100,000 paticles for five example data in test set. Two left columns show evolution of  $\theta_1$  and  $\theta_2$  (rad) (, therefore, the values are cyclic with the period of  $2\pi \simeq 6.3$ , and we corrected for those periodic shifts) and the two right columns show  $\omega_1$  and  $\omega_2$  (rad/s).

of  $f_{\theta}$ , while  $\phi_{\theta}$  serves as the inverse of  $f_{\theta}$ . The initial eigenvalues are randomly sampled from the range between  $e^0$  and  $e^{0.01}$ . 

**Training:** All training variables (network weights for the IOO  $(f_{\theta}, G_{\theta})$ , the emission model op-erator  $\phi$ , eigenvalues  $\lambda$  for the dynamics matrix A, Gaussian noise parameter  $\sigma$  for angular velocity  $\omega$ , and the concentration parameter for Von Mises distribution used for angular coordinate  $\theta$ ) are trained together. 

**Examples:** Here, we show examples for assimilated  $\theta$  and  $\omega$  in Fig. 10. Also, we give an additional figure for the RMSE of  $\theta$  for various methods. 



Figure 11: Assimilation results for the angle variable  $\theta$ . All models successfully determine the angle coordinate in spite of the strong nonlinearity in the observation (trigonometric function). Among these, performance of DBF is the best.

Table 5: List of hyperparameters for Lorenz96 experiment.

parameter	value
$R_{init}$	diag[1]
Q	diag[ $e^{-8}$ ]

1138 1139

1134

1135 1136 1137

1140 1141

#### 1142 C.3 LORENZ96

**Dataset:** The dataset consists of physical and observed variables sampled at 40 grid points. The training set includes 25,600,000 initial conditions, while the test set contains 10 initial conditions. The number of training samples is sufficiently large to ensure that the training converges in most cases. The original datasets comprise 80 time steps. Numerical integration is performed using the solve\_ivp function in SciPy, with a relative tolerance rtol =  $10^{-2}$  and an absolute tolerance of atol =  $10^{-2}$ . Gaussian noise with standard deviations of  $\sigma = 1, 3$ , or 5 is added to all measurements.

For KalmanNet, we attempted to train with 25,600,000 and 400,000 initial conditions; however, the process was terminated due to memory limitations. Consequently, we report results using a dataset size of 120,000. For DKF, VRNN, and SRNN, we also tried training with 25,600,000 conditions, but all models encountered a RuntimeError due to instability during the backward computation. To obtain results, we reduced the number of training samples to 512,000. With this adjustment, both SRNN and VRNN successfully completed the training procedure for some initial conditions.

1156 1157 A physical quantity  $z_j$  is defined at each grid point  $j(1 \le j \le 40)$ . The time evolution of this quantity is described by the following set of differential equations:

1159 1160

1161

1165

1169

1170

1171

1172

$$\frac{dz(t)_j}{dt} = (z_{j+1} - z_{j-2})z_{j-1} - z_j + F, (1 \le j \le 40)$$
(20)

1162 In this equation, the driving term F is set to 8. The first term models the advection of the physical 1163 quantity, while the second term represents its diffusion along a fixed latitude. With these parameters, 1164 the evolution of the physical quantity exhibits chaotic behavior.

**Network architecture:** The NN  $f_{\theta}$  consists of ten convolutional blocks followed by a fully connected layer. Each convolutional block comprises a 1D convolution, layer normalization, and a skip connection:

- conv1d: nn.Conv1d( c<sub>in</sub>, c<sub>out</sub>, kernel\_size=5, padding=2, padding\_mode="circular", )
- norm: layer normalization,
- skip: skip connection.

The first convolutional block has  $c_{in} = 1$  and  $c_{out} = 20$ , expanding the input by a factor of 20 in the channel dimension. The subsequent eight layers maintain 20 channels. Finally, the 20 channels and 40 physical dimensions are flattened into 800-dimensional variables, which are then fed into a fully connected layer of size  $800 \times 800$ . For all layers, the activation function used is ReLU. The function  $G_{\theta}$  is structured identically to  $f_{\theta}$ , while  $\phi_{\theta}$  represents the inverse of  $f_{\theta}$ .

1178

1182

1184

1186

**Training:** All training variables, including the network weights for the inverse observation operator  $f_{\theta}$  and  $G_{\theta}$ , the emission model operator  $\phi$ , the eigenvalues  $\lambda$  for the dynamics matrix A, and the Gaussian noise parameter  $\sigma$ , are trained concurrently.

- **Examples:** We show an example figure for assimilation experiment with DBF in Fig. 12.
- 1185 C.4 MOVING MNIST (ADDITIONAL LINEAR EXPERIMENT)
- **Dataset:** The dataset consists of a series of 2D images, where each pixel has a dynamic range from 0 to 255. The training set contains 480,000 initial conditions, while the test set consists of ten initial



Figure 12: An example of assimilation output in the experiment with nonlinear observation operator. The observation is not very informative due to low threshold for saturation in the observation operator  $(o_{t,j} = min(z_{t,j}^4, 10) + \epsilon)$ , all cells with  $z_{t,j} > 1.8$  are just observed as  $10 + \epsilon$ ). In the first 20 steps, the model output resembles little with the target. However, as the step proceeds, the estimated state begins to capture features of the true state. Even with such a poor observation operator, DBF finds a latent space representation that captures the evolution of the true state.

conditions, with both datasets comprising 20 time steps each. The number of training samples and epochs is sufficiently large to ensure that the training converges effectively. A Gaussian noise with a standard deviation of  $\sigma = 50$  is added to all pixels. The MNIST images of the digits "9" (data point 5740) and "5" (data point 5742) move at constant speeds until they reach the edges, where reflection occurs.

**Training:** The network weights for  $G_{\theta}$  are fixed during the first epoch to facilitate the learning of  $f_{\theta}$  and the image tensor for the observation model. Subsequently,  $G_{\theta}$  is trained during the second epoch. In total, DBF undergoes training for two epochs. 

**Dynamics model:** Constant velocity model. The exact dynamics matrix we have used is:

 $F = \begin{pmatrix} 1 & 0 & 0 & 0 & ut & 0 & 0 & 0 \\ 0 & 1 & 0 & 0 & 0 & dt & 0 & 0 \\ 0 & 0 & 1 & 0 & 0 & 0 & dt & 0 \\ 0 & 0 & 0 & 1 & 0 & 0 & 0 & dt \\ 0 & 0 & 0 & 0 & 1 & 0 & 0 & 0 \\ 0 & 0 & 0 & 0 & 0 & 1 & 0 & 0 \\ 0 & 0 & 0 & 0 & 0 & 0 & 1 & 0 \\ 0 & 0 & 0 & 0 & 0 & 0 & 1 & 0 \\ 0 & 0 & 0 & 0 & 0 & 0 & 0 & 1 \\ \end{pmatrix}, z_t = \begin{pmatrix} u_{1,t} \\ y_{1,t} \\ v_{2,t} \\ v_{x_{1,t}} \\ v_{y_{1,t}} \\ v_{x_{2,t}} \end{pmatrix},$ 

$$z_{t+1} = F z_t \tag{21}$$

(22)

and true observation model:

 $O_{t}$ 

 $\tilde{x}_t$ 

$$=\begin{cases} (x_t \mod 16) & \text{if } x//16 \text{ is even} \\ 9 - (x_t \mod 16) & \text{if } x//16 \text{ is odd} \end{cases}, \text{ same for } y$$

$$= h(z_t), \dim(o_t) = 44 \times 44, \text{ a } 28 \times 28 \text{ image is embedded } \operatorname{at}(\tilde{x}_t, \tilde{y}_t).$$
(24)

The formulation above addresses image reflection through the observation operator, resulting in lin-ear dynamics while permitting multiple solutions for each observed figure. This approach presents significant challenges for the EnKF, which assumes a single-peak Gaussian distribution in the as-similating space. To ensure a fair comparison, we revise the dynamics and observation models to allow for a single solution for each figure. This adjustment notably enhances the performance of the EnKF if the image is provided. However, even with this modification, the EnKF fails to accurately estimate the position, velocity, and the embedded image.

**Network architecture:**  $f_{\theta}$ : Two-dimension convolutional NNs. Below is the list of layers.

- conv1: nn.Conv2d(1, 2, kernel\_size=3, stride=2, padding=1)
- conv2: nn.Conv2d(2, 4, kernel\_size=3, stride=2, padding=1)
- conv3: nn.Conv2d(4, 4, kernel\_size=3, stride=1, padding=1)
- conv4: nn.Conv2d(4, 4, kernel\_size=3, stride=1, padding=1)
- fc: nn.Linear $(11 \times 11 \times 4, 8)$

The input image, sized  $44 \times 44$ , is sequentially processed by convolutional layers (conv1, conv2, conv3, and conv4). The output is then flattened to serve as the input for the fully connected layer (fc). Ultimately, this process yields eight variables for  $f_{\theta}(o_t)$ . The network  $G_{\theta}$  follows the same architecture as  $f_{\theta}$ , but it produces only the diagonal components of  $G_{\theta}(o_t)$  through the NN. 

- **Example figures:** In Fig. 14, we show example images for observations and all the algorithms in image-informed setting.

1296						data	0				
1297		t=0	t=2	t=4	t=6	t=8	t=10	t=12	t=14	t=16	t=18
1298		S	12	C	Calls.	6	1		N.S. State	and the second	Sterry Street
1299		ê 🥱	ୁ 🤧	· 🥹	9	્છ	10	10	25.2	- 59	<b>1</b>
1201		214080020002	1200010000		and the second states	- MONE - 143			(2013) (A		
1202		њ	- &	- 6	- 5-	6	6	6	80		
1302				-	-	-	~	~	~	~	×
1303		щ	Edit.	Sec.	1000		678	21.000 TO 20.000 TO	1000	all and a	St. Antonio
1305		EnK				2.3	学会	1		1	1000
1306		statute.	1000	8.8.0-FA	10.110	16.5 9 <sup>3</sup> 8 <sup>24</sup>	8.80 A.87	R. R. A. P.	4,924.33		R. M. LAND
1307		ЧĂ		Sterios	100	2012		100	2012		100 A 2
1308		Ш	SCR1						- Carlor		
1309			-								
1310		H State	<b>1</b>	- <b>1</b>	<b>6</b>		<b>6</b> 7.				<b>1</b>
1311		3.1 11			1000 C						
1312											
1313						data	20				
1314		t=0	t=2	t=4	t=6	t=8	t=10	t=12	t=14	t=16	t=18
1315		v C	T.C.	1	1.2		100	人的标志		States -	
1316		do 🤤	9	: <u>@</u>	9	3	9	S 😪	S. 😌	S 😒	_ <b>**</b> *
1317		No. Providence and the second	10000000000	North China Conto	HORE HORE HE	2010/02/02/02/02/02/02/02/02/02/02/02/02/02	244 Carlo 2010		SECRET SECOND		200000000000000
1318		в 🭝	\$	8	8	6	<b>6</b>	<b>~</b>	9	9	9
1319				₽	<b>-</b>	· •	2	- 🔫	· •	-	
1320		₩ 2000	13763	44.83	13434	12.12	1212	1.43	1200	- 24	328
1321		Enl	225	300	35	184	2.2	33.50		32.24	1230
1323			5.M.M.	ALCON DO				522	(* 165)		51555.4
1324		TKF									1 Sec
1325		- Segur	25494	and the	4784835	204.54	14.95	25.59	TO SAC	- 13 <b>4</b> 68	
1326		ц 500				1.00	234			-	
1327		₽.		•							
1328											
1329											
1330						data	27				
1331		t=0	t=2	t=4	t=6	t=8	t=10	t=12	t=14	t=16	t=18
1332		sdo	- 🐨	8	୍ଷ 😓 🛛	8	50	- 9	9	9	6
1333		12年1日1月		當時的	CE STER	的体质		12-14-23	12 8 18	(2) (4) (4)	Stark:
1334		ш				6	50	50	5	6	- 5
1335		80 🛃	- <b>N</b>		· • • •	-	-	4	-		<b>2</b>
1337		10000-00001	128 20	1.3847750	1.12:0703	12:231	1002231	(茶店)	「検索型」	1285	194221
1338		INK	- 2 M	12.00	a de la	1 dec	100	12.00	S PR		101
1339		Ш. 1964 2.46 У		a case o	1.0001.001					18238	
1340		¥		300	25.3	3.5	202	1	28.4	100	100
1341		II SAS	1969840	0.963	1.1.1.1.1.1.1	13.6168	101950	· 31 86/3	25240	20085	194,0463
1342		201			and the second						
1343		H	-		-	1			No.		
1344		1997年1月									
1345											
1346											
1347	<b>D'</b> 1/		C	<b>c</b> ,	1 1		OHOT			• .1	<b>.</b>



1350					data	0				
1351	t=0	t=2	t=4	t=6	t=8	t=10	t=12	t=14	t=16	t=18
1352	S	- C		A Call		100		R.S. Stelly	Strenth -	
1353	e 🚱	୍ 🤧	9	9	્છ	10	10	35.2	49	S 2
1354		i sente		2.822.0012	1000000000					
1355	н 🔎	2	5	5	5	5	6	80		
1356	ē 🦻	9	9	9	- 10	-		20	1	<b>*</b>
1357	~	~	-	-						
1358	¥ 🔿	- 🕹	- 5-	- 5-	6	59	50	-	<b>\$</b>	<b>\$</b>
1359	ш	<i>•</i>	~	~	- <b>-</b>	~	-			- <b>-</b>
1360	щ		~		~		0		-	
1361	H 🤤	- <b>E</b>	8	S	- <del>22</del> -	8	7	1	- 🤶	<del>दि</del>
1362					•		•	•	•	
1363	ш 😒 👘			-			0		<u>_</u>	
1364	≏	2	2	- <b></b>	-	· 😎	- <b>27</b>	2		-
1365										
1366										
1367					data	20				
1368	t=0	t=2	t=4	t=6	t=8	t=10	t=12	t=14	t=16	t=18
1369	s s	200			是自己。	14 A	and the second		Sec. 1	
1370	e 😔		1 🔶	4	୍ 🥵	-	S 🚱	ୁ 🤤	S 😔 (	- <del>1</del>
1371				5162.0613						
1372	њ <u> </u>	5	2	~	~	-	~	0		
1373		9	9	8	8	- 😂	- <b>S</b>	- 🤝	- 🤝	
1374										
1375		95	<u>a</u>	8	- 65		<b>(</b>	9	<b>\$9</b>	<b>2</b>
1376	□ 🤛		2		· •	- <del></del>		<b>•</b>	3	
1377	யு 🥐	~		6		~	0	~	-	
1378	Ě 🥬	9	<b>27</b>	وچ ا	- 😎	6	2	5	85	9
1379			_							
1380	բ 😏	<u>_</u>	<u></u>	<u></u>	<u></u>	4			-	<u></u>
1381		-		-	-	-	-	•	-	-
1382										
1383										
1384					data	27				
1385	t=0	t=2	t=4	t=6	t=8	t=10	t=12	t=14	t=16	t=18
1386	S 🖘	<b>1</b>	0	2	Sea.	500	Sa	5	C.	6
1387	a series a	- <b></b>		a state	State of		Salar Sa			
1388					Contra al Orico e An		-	-		
1389	ВГ	· 😌 -	9	- 😒	5	50	50	5	5	- 🚑
1300		- <b>-</b> -						-		-
1301	щ					6	0	6	<b>_</b>	-
1301	Ye 🛃	- 😏	- <del>22</del>	5	- 😏	· è	<b>&gt;</b>	- 29	- 🌮 -	S 😂 🛛
1332	•									
1333	щ 🤪	<u></u>				0	0	9	0	0
1205	۳ <u>۲</u>	-	~	- 😕 -	95	5	5	5	5	5
1305										
1007	н 😌 🗌	<b>~</b>		<b>2</b>	<b>2</b>		4	9	<b>~</b>	- 65
1000						-	~			-
1000										
1399										



Table 6: List of hyperparameters for moving MNIST experiment. 1405 1406 parameter value 1407  $diag[e^6]$ R1408 Qdiag $[e^{-4}]$ 1409 1410 1411 Table 7: The success rates of different methodologies in the two-body moving MNIST problem. For 1412 the model-based approaches, we used the same dynamics and observation models that generated the 1413 data. For DBF, the model was initialized with random image tensors and trained solely on the data. 1414 1415 Method Success rate 1416 DBF 100% (50/50) 1417 EnKF 58% (29/50) 1418 ETKF 0% (0/50) 1419 PF 0% (0/50) 1420 1421 D TRAINING STABILITY 1422 1423 We observe that the training of our proposed method is stable compared to RNN-based models. 1424 Fig. 15 shows the evolution of the real parts of eigenvalues. Although we do not impose constraints 1425 on the real parts of eigenvalues, the values only marginally exceed one. Therefore, long-time dy-1426 namics is stable during training. 1427 1428 1429 1430 1431 1432 1433 1434 0.9 1.0 abs(eigenvalue) 0.9 1.0 abs(eigenvalue) 0.9 1.0 abs(eigenvalue) 0.8 0.9 1.0 nvalue) 0.8 1.1 0.8 1.1 0.8 1.1 0.8 ດ່ອ 10 0.8 0,0 1.1 1.1 10 1.1 1435 abs(eigenvalue) 1436 Figure 15: Evolution of histograms for the real parts of 800 complex eigenvalues in Lorenz96 exper-1437 iment. Initially, eigenvalues are taken as one. As the model learns the dynamics, eigenvalues lower 1438 than 1.0 appear. However, the largest eigenvalue  $\lambda_{max}$  mostly remains less than 1.02. 1439 1440 1441 1442 1443 Е HYPERPARAMETER STUDY ON THE LATENT DIMENSIONS 1444 1445 The dimension of the latent variables is a hyperparameter. We have tested the performance and 1446 computation (both training and inference) time for nonlinear problems. 1447 E.1 DOUBLE PENDULUM 1448 1449 1450 E.1.1 ACCURACY-COMPUTE TRADE-OFF IN DBF 1451 1452 For double pendulum problem, we test with the standard observation operator with the observation 1453 noise of  $\sigma = 0.1$ . Figs. 16, 17 show the relation between the RMSE and the latent dimensions of the 1454 system. Here, we show results with  $1.0 \times 10^7$  training data. For the double pendulum problem, we

1455 have tested with 4, 20, 80, and 200 latent dimensions. All the latent dimensions tested were too small 1456 to observe the impact of the compute-latent dimension relationship. To observe the slowdown, we 1457 need to test with higher dimensions. Please also refer to the results for Lorenz96. The performance (RMSE at the final 10 steps) for the angles  $\theta$  and angle velocities  $\omega$  are poor if the latent dimension



Figure 16: Left panel: RMSE as a function of the latent dimensionality of DBF. Right panel: the inference time as a function of the latent dimensionality of DBF.



1481Figure 17: Left panel: the training time for  $1.0 \times 10^7$  initial conditions as a function of the latent1482dimension. Right panel: RMSE as a function of the training time for five different numbers of latent1483dimensions.

1466

1467

1468

is four. By leveraging 20 latent dimensions, DBF achieves a very good assimilation performance. Further enhancing the latent dimensions to 80 and 200 did not improve performance. The training gradually gets slower when we use latent dimensions higher than 80. As can be seen from Fig. 16, The performance is rather insensitive to the latent dimensions in the range of [20, 200]: the RMSE for  $\theta$  is 0.036 at  $dim(h_t) = 20$ , 0.053 at  $dim(h_t) = 80$ , and 0.044 at  $dim(h_t) = 200$  and for  $\omega$  is 0.265 at  $dim(h_t) = 20$ , 0.375 at  $dim(h_t) = 80$ , and 0.302 at  $dim(h_t) = 200$ .

1492

### 1493 E.1.2 COMPARISON TO THE PF

The performance of PF depends on the number of particles used. We have tested with 20, 200, 2,000, 20,000, and 100,000 particles. The performance for  $\theta$  improves significantly if we use more than 200 particles. The RMSE for the angle velocities  $\omega$  almost saturates at RMSE  $\simeq 0.31$  when we use particles more than 20,000. To achieve that accuracy, the inference time required for PF is more than 200 seconds per initial condition. On the other hand, DBF achieves slightly better performance (RMSE  $\simeq 0.265$ ) with the latent dimensions of 20. The inference time for DBF is 0.1 seconds per batch.

1501 1502

1503 E.2 LORENZ96

1504 1505

# 1506 E.2.1 ACCURACY-COMPUTE TRADE-OFF IN DBF

For Lorenz96 problem, we test with the nonlinear observation operator with the observation noise of  $\sigma = 1$ . Figs. 19, 20 show the relation between the RMSE and the latent dimensions of the system. Here, we show results with  $1.0 \times 10^7$  training data. The dimensionality of the latent variables can be either larger or smaller than that of the physical variables, but there is a tradeoff: up to a certain latent dimensionality, increasing the dimension improves performance at the



Figure 19: Left panel: RMSE as a function of the latent dimensionality of DBF. Right panel: the inference time as a function of the latent dimensionality of DBF.

cost of longer computation time. Beyond that point, increasing the latent dimensionality no longer improves performance but only increases training time (although inference time remains relatively short compared to model-based approaches). Therefore, the optimal balance depends on the specific problem. For the Lorenz96 system, a dimensionality of 800 was a reasonable trade-off among 20, 80, 200, 800, and 2,000 dimensions. As shown in the figure, the RMSE changes by only 7 percent (1.31 vs 1.23) in the range from 200 to 2,000 dimensions, indicating that the impact is not critical in this range.

1561 1562

1564

#### 1563 E.2.2 COMPARISON TO THE PF

1565 The PF also has the trade-off. Although RMSE improves slowly as we increase the number of particles, the RMSE was poor (2.27) compared to the DBF results (RMSE  $\simeq 1.3$ ) even with mas-



An example is the performance of PF as a function of the particles used. Right panel: RMSE as a function of the inference time for the DBF and the PF. For the DBF, the latent dimensions are 20, 80, 200, 800, and 2,000. For the PF, the number of particles are 20, 200, 2,000, 20,000, 100,000.

sively large number of particles (100,000) with very long inference time (2,000 seconds per initial condition)



Figure 22: the performance of DBF for a low latent dimension case  $(dim(h_t) = 20)$  and a high latent dimension case  $(dim(h_t) = 800)$ . Even with the latent dimensions (20) smaller than that of the original state space (40), DBF shows the skillful assimilation. With higher latent dimensions (800), the performance further improves.