

Backpropagation Algorithms and Reservoir Computing in Recurrent Neural Networks for the Forecasting of Complex Spatiotemporal Dynamics

P. R. Vlachas^a, J. Pathak^{b,c}, B. R. Hunt^{d,e}, T. P. Sapsis^f, M. Girvan^{b,c,d}, E. Ott^{b,c,g} and P. Koumoutsakos^{a,*}

^aComputational Science and Engineering Laboratory, ETH Zürich, Clausiusstrasse 33, Zürich CH-8092, Switzerland

^bInstitute for Research in Electronics and Applied Physics, University of Maryland, College Park, Maryland 20742, USA

^cDepartment of Physics, University of Maryland, College Park, Maryland 20742, USA

^dInstitute for Physical Science and Technology, University of Maryland, College Park, Maryland 20742, USA

^eDepartment of Mathematics, University of Maryland, College Park, Maryland 20742, USA

^fDepartment of Mechanical Engineering, Massachusetts Institute of Technology, 77 Massachusetts Avenue, Cambridge, Massachusetts 02139-4307, USA

^gDepartment of Electrical and Computer Engineering, University of Maryland, Maryland 20742, USA

ARTICLE INFO

Keywords:

Time-series forecasting
RNN, LSTM, GRU
Reservoir Computing
BPTT
Echo-state networks
Unitary RNNs
Kuramoto-Sivashinsky
complex systems

ABSTRACT

We examine the efficiency of Recurrent Neural Networks in forecasting the spatiotemporal dynamics of high dimensional and reduced order complex systems using Reservoir Computing (RC) and Backpropagation through time (BPTT) for gated network architectures. We highlight advantages and limitations of each method and discuss their implementation for parallel computing architectures. We quantify the relative prediction accuracy of these algorithms for the long-term forecasting of chaotic systems using as benchmarks the Lorenz-96 and the Kuramoto-Sivashinsky (KS) equations. We find that, when the full state dynamics are available for training, RC outperforms BPTT approaches in terms of predictive performance and in capturing of the long-term statistics, while at the same time requiring much less training time. However, in the case of reduced order data, large scale RC models can be unstable and more likely than the BPTT algorithms to diverge. In contrast, RNNs trained via BPTT show superior forecasting abilities and capture well the dynamics of reduced order systems. Furthermore, the present study quantifies for the first time the Lyapunov Spectrum of the KS equation with BPTT, achieving similar accuracy as RC. This study establishes that RNNs are a potent computational framework for the learning and forecasting of complex spatiotemporal systems.

1. Introduction

In recent years we have observed significant advances in the field of machine learning (ML) that rely on potent algorithms and their deployment on powerful computing architectures. Some of these advances have been materialized by deploying ML algorithms on dynamic environments such as video games (Ha and Schmidhuber, 2018; Schrittwieser et al., 2019) and simplified physical systems (AI gym) (Brockman et al., 2016; Mnih et al., 2015; Silver et al., 2016). Dynamic environments are often encountered across disciplines ranging from engineering and physics to finance and social sciences. They can serve as bridge for scientists and engineers to advances in machine learning and at the same time they present a fertile ground for the development and testing of advanced ML algorithms (Hassabis et al., 2017). The deployment of advanced machine learning algorithms to complex systems is in its infancy. We believe that it deserves further exploration as it may have far-reaching implications for societal and scientific challenges ranging from weather and climate prediction (Weyn et al., 2019; Gneiting and Raftery, 2005), to energy networks, medicine (Esteva et al., 2017; Kurth et al., 2018), and the dynamics of ocean dynamics and turbulent flows (Aksamit et al., 2019; Sünderhauf et al., 2018; Brunton et al., 2020).

*Corresponding author

✉ vlachas@collegium.ethz.ch (P.R. Vlachas); jpathak@umd.edu (J. Pathak); bhunt@umd.edu (B.R. Hunt); sapsis@mit.edu (T.P. Sapsis); girvan@umd.edu (M. Girvan); edott@umd.edu (E. Ott); petros@ethz.ch (P. Koumoutsakos)

ORCID(s): 0000-0002-3311-2100 (P.R. Vlachas)

Complex systems are characterized by multiple, interacting spatiotemporal scales that challenge classical numerical methods for their prediction and control. The dynamics of such systems are typically chaotic and difficult to predict, a critical issue in problems such as weather and climate prediction. Recurrent Neural Networks (RNNs), offer a potent method for addressing these challenges. RNNs were developed for processing of sequential data, such as time-series (Hochreiter and Schmidhuber, 1997), speech (Graves and Jaitly, 2014), and language (Dong et al., 2015; Cho et al., 2014). Unlike classical numerical methods that aim at discretizing existing equations of complex systems, RNN models are data driven. RNNs keep track of a hidden state, that encodes information about the history of the system dynamics. Such data-driven models are of great importance in applications to complex systems where equations based on first principles may not exist, or may be expensive to discretize and evaluate, let alone control, in real-time.

Early application of neural networks for modeling and prediction of dynamical systems can be traced to the work of Lapedes et. al. (Lapedes and Farber, 1987), where they demonstrated the efficiency of feedforward artificial neural networks (ANNs) to model deterministic chaos. As an alternative to ANNs, wavelet networks were proposed in (Cao et al., 1995) for chaotic time-series prediction. However, these works have been limited to intrinsically low-order systems, and they have been often deployed in conjunction with dimensionality reduction tools. As shown in this work, RNNs have the potential to overcome these scalability problems and be applied to high-dimensional spatio-temporal dynamics. The works of Takens (Takens, 1981) and Sauer, Yorke and Casdagli (Sauer et al., 1991) showed that the dynamics on a D-dimensional attractor of a dynamical system can be unfolded in a time delayed embedding of dimension greater than 2D. The identification of a useful embedding and the construction of a forecasting model have been the subject of life-long research efforts (Bradley and Kantz, 2015). More recently, in (Lusch et al., 2018), a data-driven method based on the Koopman operator formalism (Koopman, 1931) was proposed, using feed-forward ANNs to identify an embedding space with linear dynamics that is then amenable to theoretical analysis.

There is limited work at the interface of RNNs and nonlinear dynamical systems (Vlachas et al., 2018; Wan et al., 2018; Pathak et al., 2017, 2018a; Lu et al., 2018). Here we examine and compare two of the most prominent nonlinear techniques in the forecasting of dynamical systems, namely RNNs trained with backpropagation and Reservoir Computing (RC). We note that our RC implementation also uses a recurrent neural network, but according to the RC paradigm, it does not train the internal network parameters. We consider the cases of fully observed systems as well as the case of partially observed systems such as reduced order models of real world problems, where typically we do not have access to all the degrees-of-freedom of the dynamical system.

Reservoir Computing (RC) has shown significant success in modeling the full-order space dynamics of high dimensional chaotic systems. This success has sparked the interest of theoretical researchers that proved universal approximation properties of these models (Grigoryeva and Ortega, 2018; Gonon and Ortega, 2019). In (Pathak et al., 2017, 2018b) RC is utilized to build surrogate models for chaotic systems and compute their Lyapunov exponents based solely on data. A scalable approach to high-dimensional systems with local interactions is proposed in (Pathak et al., 2018a). In this case, an ensemble of RC networks is used in parallel. Each ensemble member is forecasting the evolution of a group of modes while all other modes interacting with this group is fed at the input of the network. The model takes advantage of the local interactions in the state-space to decouple the forecasting of each mode group and improve the scalability.

RNNs are architectures designed to capture long-term dependencies in sequential data (Pascanu et al., 2013; Bengio et al., 1994; Hochreiter, 1998; Goodfellow et al., 2016). The potential of RNNs for capturing temporal dynamics in physical systems was explored first using low dimensional RNNs (Elman, 1990) without gates to predict unsteady boundary-layer development, separation, dynamic stall, and dynamic reattachment back in 1997 (Faller and Schreck, 1997). The utility of RNNs was limited by the finding that during the learning process the gradients may vanish or explode. In turn, the recent success of RNNs is largely attributed to a cell architecture termed Long Short-Term Memory (LSTM). LSTMs employ gates that effectively remember and forget information thus alleviating the problem of vanishing gradients (Hochreiter, 1998). In recent years (Bianchi et al., 2017) RNN architectures have been bench-marked for short-term load forecasting of demand and consumption of resources in a supply network, while in (Laptev et al., 2017) they are utilized for extreme event detection in low dimensional time-series. In (Wan and Sapsis, 2018) LSTM networks are used as surrogates to model the kinematics of spherical particles in fluid flows. In (Vlachas et al., 2018) RNNs with LSTM cells were utilized in conjunction with a mean stochastic model to capture the temporal dependencies and long-term statistics in the reduced order space of a dynamical system and forecast its evolution. The method demonstrated better accuracy and scaling to high-dimensions and longer sequences than Gaussian Processes (GPs). In (Wan et al., 2018) the LSTM is deployed to model the residual dynamics in an imperfect Galerkin-based reduced order model derived from the system equations. RC and LSTM networks are applied in the long-term forecasting

of partially observable chaotic chimera states in (Neofotistos et al., 2019), where instead of a completely model-free approach, ground-truth measurements of currently observed states are helping to improve the long-term forecasting capability. RNNs are practical and efficient data-driven approximators of chaotic dynamical systems, due to their (1) universal approximation ability (Schäfer and Zimmermann, 2006; Siegelmann and Sontag, 1995) and (2) ability to capture temporal dependencies and implicitly identify the required embedding for forecasting.

Despite the rich literature on both methods there are limited comparative studies of the two frameworks. The present work aims to fill this gap by examining these two machine learning algorithms on challenging physical problems. We compare the accuracy, performance, and computational efficiency of the two methods on the full-order and reduced-order modeling of two prototype chaotic dynamical systems. We also examine the modeling capabilities of the two approaches for reproducing correct Lyapunov Exponents and frequency spectra. Moreover, we include in the present work some more recent RNN architectures, like Unitary (Arjovsky et al., 2016; Jing et al., 2017) and Gated Recurrent Units (GRUs) (Chung et al., 2014; Cho et al., 2014) that have shown superior performance over LSTMs for a wide variety of language, speech signal and polyphonic music modeling tasks.

We are interested in model-agnostic treatment of chaotic dynamical systems, where the time evolution of the full state or some observable is available, but we do not possess any knowledge about the underlying equations. In the latter case, we examine which method is more suitable for modeling temporal dependencies in the reduced order space (observable) of dynamical systems. Furthermore, we evaluate the efficiency of an ensemble of RNNs in predicting the full state dynamics of a high-dimensional dynamical system in parallel and compare it with that of RC. Finally, we discuss the advantages, implementation aspects (such as RAM requirements and training time) and limitations of each model. We remark that the comparison in terms of time and RAM memory consumption, does not aim to quantify advantages/drawback among models but rather provide information for the end users of the software.

We hope that the present study may open to the ML community a new arena with highly structured and complex environments for developing and testing advanced new algorithms (Hassabis et al., 2017). At the same time it may offer a bridge to the physics community to appreciate and explore the importance of advanced ML algorithms for solving challenging physical problems (Brunton et al., 2020).

The structure of the paper is as follows. Section 2 provides an introduction to the tasks and an outline of the architectures and training methods used in this work. Section 3 introduces the measures used to compare the efficiency of the models. In Section 4 the networks are compared in forecasting reduced order dynamics in the Lorenz-96 system. In Section 5, a parallel architecture leveraging local interactions in the state space is introduced and utilized to forecast the dynamics of the Lorenz-96 system (Lorenz, 1995) and the Kuramoto-Sivashinsky equation (Kuramoto, 1978). In Section 6 the GRU and RC networks are utilized to reproduce the Lyapunov spectrum of the Kuramoto-Sivashinsky equation, while Section 7 concludes the paper.

2. Methods - Sequence Modeling

We consider machine learning algorithms for time-series forecasting. The models are trained on time-series of an observable $\mathbf{o} \in \mathbb{R}^{d_o}$ sampled at a fixed rate $1/\Delta t$, $\{\mathbf{o}_1, \dots, \mathbf{o}_T\}$, where we eliminate Δt from the notation for simplicity. The models possess an internal high-dimensional hidden state denoted by $\mathbf{h}_t \in \mathbb{R}^{d_h}$ that enables the encoding of temporal dependencies on past state history. Given the current observable \mathbf{o}_t , the output of each model is a forecast $\hat{\mathbf{o}}_{t+1}$ for the observable at the next time instant \mathbf{o}_{t+1} . This forecast is a function of the hidden state. As a consequence, the general functional form of the models is given by

$$\mathbf{h}_t = f_h^h(\mathbf{o}_t, \mathbf{h}_{t-1}), \quad \hat{\mathbf{o}}_{t+1} = f_h^o(\mathbf{h}_t), \quad (1)$$

where f_h^h is the hidden-to-hidden mapping and f_h^o is the hidden-to-output mapping. All recurrent models analyzed in this work share this common architecture. They differ in the realizations of f_h^h and f_h^o and in the way the parameters or **weights** of these functions are learned from data, i.e., **trained**, to forecast the dynamics.

2.1. Long Short-Term Memory

In Elman RNNs (Elman, 1990), the vanishing or exploding gradients problem stems from the fact that the gradient is multiplied repeatedly during back-propagation through time (Werbos, 1988) with a recurrent weight matrix. As a consequence, when the spectral radius of the weight matrix is positive (negative), the gradients are prone to explode (shrink). The LSTM (Hochreiter and Schmidhuber, 1997) was introduced in order to alleviate the vanishing gradient

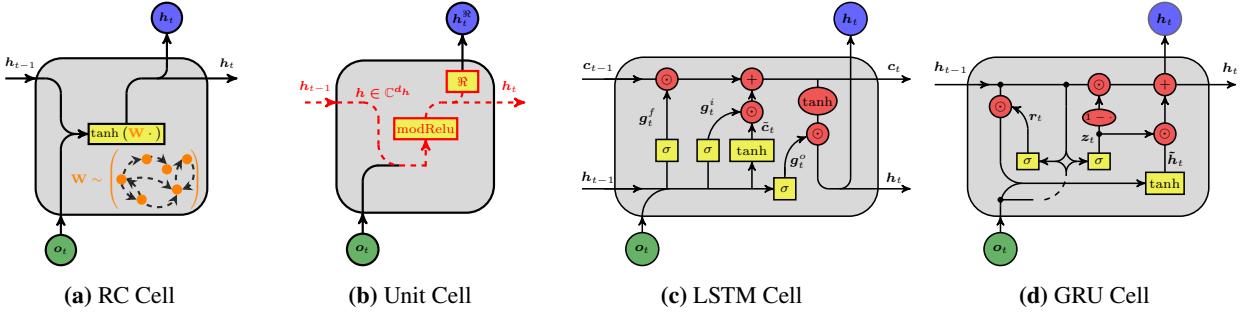


Figure 1: The information flow for a Reservoir Computing (RC) cell, a complex Unitary cell (Unit), a Long Short-Term Memory (LSTM) cell and a Gated Recurrent Unit (GRU) cell. The cells were conceptualized to tackle the vanishing gradients problem of Elman-RNNs. The cell used in RC is the standard architecture of the Elman-RNN. However, the weights of the recurrent connections are randomly picked to satisfy the echo state property and create a large reservoir of rich dynamics. Only the output weights are trained (e.g., with ridge regression). The Unitary RNN utilizes a complex unitary matrix to ensure that the gradients are not vanishing. LSTM and GRU cells employ gating mechanisms that allow forgetting and storing of information in the processing of the hidden state. Ellipses and circles denote entry-wise operations, while rectangles denote layer operations. The information flow of the complex hidden state in the Unitary RNN is illustrated with dashed red color, while the untrained randomly picked weights of the RC with orange.

problem of Elman RNNs (Hochreiter, 1998) by leveraging gating mechanisms that allow information to be forgotten. The equations that implicitly define the recurrent mapping f_h^h of the LSTM are given by

$$\begin{aligned} g_t^f &= \sigma_f(\mathbf{W}_f[\mathbf{h}_{t-1}, \mathbf{o}_t] + \mathbf{b}_f) & g_t^i &= \sigma_i(\mathbf{W}_i[\mathbf{h}_{t-1}, \mathbf{o}_t] + \mathbf{b}_i) \\ \tilde{c}_t &= \tanh(\mathbf{W}_c[\mathbf{h}_{t-1}, \mathbf{o}_t] + \mathbf{b}_c) & c_t &= g_t^f \odot c_{t-1} + g_t^i \odot \tilde{c}_t \\ g_t^o &= \sigma_h(\mathbf{W}_h[\mathbf{h}_{t-1}, \mathbf{o}_t] + \mathbf{b}_h) & \mathbf{h}_t &= g_t^o \odot \tanh(c_t), \end{aligned} \quad (2)$$

where $\mathbf{g}_t^f, \mathbf{g}_t^i, \mathbf{g}_t^o \in \mathbb{R}^{d_h}$, are the gate vector signals (forget, input and output gates), $\mathbf{o}_t \in \mathbb{R}^{d_o}$ is the observable input at time t , $\mathbf{h}_t \in \mathbb{R}^{d_h}$ is the hidden state, $c_t \in \mathbb{R}^{d_h}$ is the cell state, while $\mathbf{W}_f, \mathbf{W}_i, \mathbf{W}_c, \mathbf{W}_h \in \mathbb{R}^{d_h \times (d_h + d_o)}$, are weight matrices and $\mathbf{b}_f, \mathbf{b}_i, \mathbf{b}_c, \mathbf{b}_h \in \mathbb{R}^{d_h}$ biases. The symbol \odot denotes the element-wise product. The activation functions σ_f, σ_i and σ_h are sigmoids. For a more detailed explanation of the LSTM cell architecture refer to (Hochreiter and Schmidhuber, 1997). The dimension of the hidden state d_h (number of hidden units) controls the capacity of the cell to encode history information. The hidden-to-output functional form f_h^o is given by a linear layer

$$\hat{\mathbf{o}}_{t+1} = \mathbf{W}_o \mathbf{h}_t, \quad (3)$$

where $\mathbf{W}_o \in \mathbb{R}^{d_o \times d_h}$. The forget gate bias is initialized to one according to (Jozefowicz et al., 2015) to accelerate training. An illustration of the information flow in a LSTM cell is given in Figure 1c.

2.2. Gated Recurrent Unit

The Gated Recurrent Unit (GRU) (Cho et al., 2014) was proposed as a variation of LSTM utilizing a similar gating mechanism. Even though GRU lacks an output gate and thus has fewer parameters, it achieves comparable performance with LSTM in polyphonic music and speech signal datasets (Chung et al., 2014). The GRU equations are given by

$$\begin{aligned} z_t &= \sigma_g(\mathbf{W}_z[\mathbf{h}_{t-1}, \mathbf{o}_t] + \mathbf{b}_z) \\ r_t &= \sigma_g(\mathbf{W}_r[\mathbf{h}_{t-1}, \mathbf{o}_t] + \mathbf{b}_r) \\ \tilde{\mathbf{h}}_t &= \tanh(\mathbf{W}_h[r_t \odot \mathbf{h}_{t-1}, \mathbf{o}_t] + \mathbf{b}_h) \\ \mathbf{h}_t &= (1 - z_t) \odot \mathbf{h}_{t-1} + z_t \odot \tilde{\mathbf{h}}_t, \end{aligned} \quad (4)$$

where $\mathbf{o}_t \in \mathbb{R}^{d_o}$ is the observable at the input at time t , $\mathbf{z}_t \in \mathbb{R}^{d_h}$ is the update gate vector, $\mathbf{r}_t \in \mathbb{R}^{d_h}$ is the reset gate vector, $\tilde{\mathbf{h}}_t \in \mathbb{R}^{d_h}$, $\mathbf{h}_t \in \mathbb{R}^{d_h}$ is the hidden state, $\mathbf{W}_z, \mathbf{W}_r, \mathbf{W}_h \in \mathbb{R}^{d_h \times (d_h + d_o)}$ are weight matrices and $\mathbf{b}_z, \mathbf{b}_r, \mathbf{b}_h \in \mathbb{R}^{d_h}$

biases. The gating activation σ_g is a sigmoid. The output \hat{o}_{t+1} is given by the linear layer:

$$\hat{o}_{t+1} = \mathbf{W}_o \mathbf{h}_t, \quad (5)$$

where $\mathbf{W}_o \in \mathbb{R}^{d_o \times d_h}$. An illustration of the information flow in a GRU cell is given in Figure 1d.

2.3. Unitary Evolution

Unitary RNNs (Arjovsky et al., 2016; Jing et al., 2017), similar to LSTMs and GRUs, aim to alleviate the vanishing gradients problem of plain RNNs. Here, instead of employing sophisticated gating mechanisms, the effort is focused on the identification of a re-parametrization of the recurrent weight matrix, such that its spectral radius is a-priori set to one. This is achieved by optimizing the weights on the subspace of complex unitary matrices. The architecture of the Unitary RNN is given by

$$\begin{aligned} \mathbf{h}_t &= \text{modReLU}\left(\mathbf{W}_h \mathbf{h}_{t-1} + \mathbf{W}_o \mathbf{o}_t\right) \\ \hat{o}_{t+1} &= \mathbf{W}_o \Re(\mathbf{h}_t), \end{aligned} \quad (6)$$

where $\mathbf{W}_h \in \mathbb{C}^{d_h \times d_h}$ is the complex unitary **recurrent** weight matrix, $\mathbf{W}_o \in \mathbb{C}^{d_h \times d_o}$ is the complex **input** weight matrix, $\mathbf{h}_t \in \mathbb{C}^{d_h}$ is the complex state vector, $\Re(\cdot)$ denotes the real part of a complex number, $\mathbf{W}_o \in \mathbb{R}^{d_h \times d_h}$ is the real output matrix, and the modified ReLU non-linearity modReLU is given by

$$\left(\text{modReLU}(\mathbf{z}) \right)_i = \frac{z_i}{|z_i|} \odot \text{ReLU}(|z_i| + b_i), \quad (7)$$

where $|z_i|$ is the norm of the complex number z_i . The complex unitary matrix \mathbf{W}_h is parametrized as a product of a diagonal matrix and multiple rotational matrices. The reparametrization used in this work is the one proposed in (Jing et al., 2017). The complex input weight matrix $\mathbf{W}_o \in \mathbb{C}^{d_h \times d_o}$ is initialized with $\mathbf{W}_o^{re} + j \mathbf{W}_o^{im}$, with real matrices $\mathbf{W}_o^{re}, \mathbf{W}_o^{im} \in \mathbb{R}^{d_h \times d_o}$ whose values are drawn from a random uniform distribution $\mathcal{U}[-0.01, 0.01]$ according to (Jing et al., 2017). An illustration of the information flow in a Unitary RNN cell is given in Figure 1b.

In the original paper of (Jing et al., 2017) the architecture was evaluated on a speech spectrum prediction task, a copying memory task and a pixel permuted MNIST task demonstrating superior performance to LSTM either in terms of final testing accuracy or wall-clock training speed.

2.4. Back-Propagation Through Time

Backpropagation dates back to the works of (Dreyfus, 1962; Linnainmaa, 1976; Rumelhart et al., 1986), while its extension to RNNs termed Backpropagation through time (BPTT) was presented in (Werbos, 1988, 1990). A forward pass of the network is required to compute its output and compare it against the label (or target) from the training data based on an error metric (e.g. mean squared loss). Backpropagation amounts to the computation of the partial derivatives of this loss with respect to the network parameters by iteratively applying the chain rule, transversing backwards the network. These derivatives are computed analytically with automatic differentiation. Based on these partial derivatives the network parameters are updated using a first-order optimization method, e.g. stochastic gradient descent.

The power of BBTT lies in the fact that it can be deployed to learn the partial derivatives of the weights of any network architecture with differentiable activation functions, utilizing state-of-the-art automatic differentiation software, while (as the data are processed in small fragments called batches) it scales to large datasets and networks, and can be accelerated by employing Graphics Processing Units (GPUs). These factors made backpropagation the workhorse of state-of-the-art deep learning methods (Goodfellow et al., 2016).

In our study, we utilize BBTT to train the LSTM (Section 2.1), GRU (Section 2.2) and Unitary (Section 2.3) RNNs. There are three key parameters of this training method that can be tuned. The first hyperparameter κ_1 is the number of forward-pass timesteps performed to accumulate the error for back-propagation. The second parameter is the number of previous time steps for the back-propagation of the gradient κ_2 . This is also denoted as truncation length, or sequence length. This parameter has to be large enough to capture the temporal dependencies in the data. However, as κ_2 becomes larger, training becomes much slower, and may lead to vanishing gradients. In the following, we characterize as stateless, models whose hidden state before κ_2 is hard-coded to zero, i.e., $\mathbf{h}_{-\kappa_2} = 0$. Stateless models cannot learn dependencies that expand in a time horizon larger than κ_2 . However, in many practical cases stateless models are widely employed

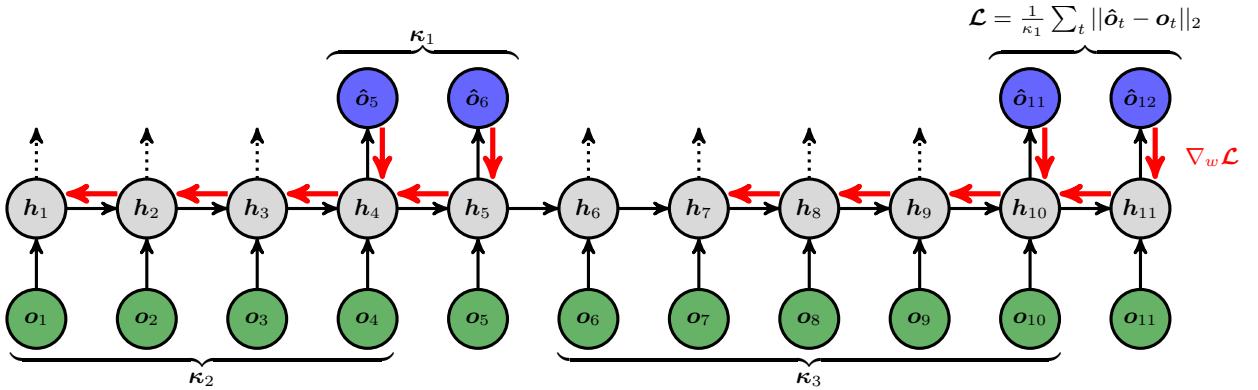


Figure 2: Illustration of an unfolded RNN. Time-series data o are provided at the input of the RNN. The RNN is forecasting the evolution of the observable at its outputs \hat{o} . The average difference (mean square error) between κ_1 iterative predictions (outputs) of the RNN \hat{o} and the targets o from the time-series data is computed every κ_3 steps. The gradient of this quantity, illustrated with red arrows, is back-propagated through time (BPTT) for κ_2 previous temporal time-steps, computing the gradients of the network parameters that are shared at each time layer. The output of intermediate steps illustrated with dashed lines is ignored. Stateless models initialize the hidden state before training at a specific fragment of the sequence of size κ_2 with zero (in this case $h_6 \hat{=} 0$) and cannot capture dependencies longer than κ_2 . In this way, consecutive training batches (sequence fragments) do not have to be temporally adjacent. In stateful models, the hidden state is never set to zero and in order to train at a specific fragment of the sequence, the initial hidden state has to be computed from the previously processed fragment. In order to eliminate the overlap between fragments, we teacher force the network with ground-truth data for $\kappa_3 \geq \kappa_2$ time-steps. In our study we pick $\kappa_3 = \kappa_2 + \kappa_1 - 1$ as illustrated in the figure.

assuming that only short-term temporal dependencies exist. In contrast, stateful models propagate the hidden state $h_{-\kappa_2} \neq 0$ between temporally consecutive batches. In our study, we consider only **stateful** networks.

Training stateful networks is challenging because the hidden state $h_{-\kappa_2}$ has to be available from a previous batch and the network has to be trained to learn temporal dependencies that may span many time-steps in the past. In order to avoid overlap between two subsequent data fragments and compute $h_{-\kappa_2}$ for the next batch update, the network is teacher-forced for κ_3 time-steps between two consecutive weight updates. That implies providing ground-truth values at the input and performing forward passing without any back-propagation. This parameter, has an influence on the training speed, as it determines how often the weights are updated. We pick $\kappa_3 = \kappa_2 + \kappa_1 - 1$ as illustrated in Figure 2, and optimize κ_1 as a hyperparameter.

The weights of the networks are initialized using the method of Xavier proposed in (Glorot and Bengio, 2010). We utilize a stochastic optimization method with adaptive learning rate called Adam (Kingma and Ba, 2015) to update the weights and biases. We add Zoneout (Krueger et al., 2017) regularization in the recurrent weights and variational dropout (Gal and Ghahramani, 2016) regularization at the output weights (with the same keep probability) to both GRU and LSTM networks to alleviate over-fitting. Furthermore, following (Vlachas et al., 2018) we add Gaussian noise sampled from $\mathcal{N}(0, \kappa_n \sigma)$ to the training data, where σ is the standard deviation of the data. The noise level κ_n is tuned. Moreover, we also vary the number of RNN layers by stacking **residual** layers (He et al., 2016) on top of each other. These deeper architectures may improve forecasting efficiency by learning more informative embedding at the cost of higher computing times.

In order to train the network on a data sequence of T time-steps, we pass the whole dataset in pieces (batches) for many iterations (epochs). An epoch is finished when the network has been trained on the whole dataset once. At the beginning of every epoch we sample uniformly $B = 32$ integers from the set $\mathcal{I} = \{1, \dots, T\}$, and remove them from it. Starting from these indexes we iteratively pass the data through the network till we reach the last (maximum) index in \mathcal{I} , training it with BBTT. Next, we remove all the intermediate indexes we trained on from \mathcal{I} . We repeat this process, until $\mathcal{I} = \emptyset$, proclaiming the end of the epoch. The batch-size is thus $B = 32$. We experimented with other batch-sizes $B \in \{8, 16, 64\}$ without significant improvement in performance of the methods used in this work.

As an additional over-fitting counter-measure we use validation-based early stopping, where 90% of the data is used for training and the rest 10% for validation. When the validation error stops decreasing for $N_{\text{patience}} = 20$ consecutive epochs, the training round is over. We train the network for $N_{\text{rounds}} = 10$ rounds decreasing the learning

rate geometrically by dividing with a factor of ten at each round to avoid tuning the learning rate of the Adam optimizer. When all rounds are finished, we pick the model with the lowest validation error among all epochs and rounds.

Preliminary work on tuning the hyperparameters of the Adam optimization algorithm apart from the learning rate, i.e. β_1 and β_2 in the original paper (Kingma and Ba, 2015), did not lead to important differences on the results. For this reason and due to our limited computational budget, we use the default values proposed in the paper (Kingma and Ba, 2015) ($\beta_1 = 0.9$ and $\beta_2 = 0.999$).

2.5. Reservoir Computing

Reservoir Computing (RC) aims to alleviate the difficulty in learning the recurrent connections of RNNs and reduce their training time (Lukoševičius and Jaeger, 2009; Lukoševičius, 2012). RC relies on **randomly selecting** the recurrent weights such that the hidden state captures the history of the evolution of the observable o_t and train the hidden-to-output weights. The evolution of the hidden state depends on the random initialization of the recurrent matrix and is driven by the input signal. The hidden state is termed reservoir state to denote the fact that it captures temporal features of the observed state history. This technique has been proposed in the context of Echo-State-Networks (ESNs) (Jaeger and Haas, 2004) and Liquid State Machines with spiking neurons (LSM) (Maass et al., 2002).

In this work, we consider reservoir computers with f_h^h given by the functional form

$$\mathbf{h}_t = \tanh(\mathbf{W}_{h,i} \mathbf{o}_t + \mathbf{W}_{h,h} \mathbf{h}_{t-1}), \quad (8)$$

where $\mathbf{W}_{h,i} \in \mathbb{R}^{d_h \times d_o}$, and $\mathbf{W}_{h,h} \in \mathbb{R}^{d_h \times d_h}$. Other choices of RC architectures are possible, including (Larger et al., 2012, 2017; Haynes et al., 2015; Antonik et al., 2017) Following (Jaeger and Haas, 2004), the entries of $\mathbf{W}_{h,i}$ are uniformly sampled from $[-\omega, \omega]$, where ω is a hyperparameter. The reservoir matrix $\mathbf{W}_{h,h}$ has to be selected in a way such that the network satisfies the “echo state property”. This property requires all of the conditional Lyapunov exponents of the evolution of \mathbf{h}_t conditioned on the input (observations \mathbf{o}_t) to be negative so that, for large t , the reservoir state \mathbf{h}_t does not depend on initial conditions. For this purpose, $\mathbf{W}_{h,h}$ is set to a large low-degree matrix, scaled appropriately to possess a spectral radius (absolute value of the largest eigenvalue) ρ whose value is a hyperparameter adjusted so that the echo state property holds¹. The effect of the spectral radius on the predictive performance of RC is analyzed in (Jiang and Lai, 2019). Following (Pathak et al., 2018a) the output coupling f_h^o is set to

$$\hat{\mathbf{o}}_{t+1} = \mathbf{W}_{o,h} \tilde{\mathbf{h}}_t, \quad (9)$$

where the augmented hidden state $\tilde{\mathbf{h}}_t$ is a d_h dimensional vector such that the i^{th} component of $\tilde{\mathbf{h}}_t$ is $\tilde{h}_t^i = h_t^i$ for half of the reservoir nodes and $\tilde{h}_t^i = (h_t^i)^2$ for the other half, enriching the dynamics with the square of the hidden state in half of the nodes. This was empirically shown to improve forecasting efficiency of RCs in the context of dynamical systems (Pathak et al., 2018a). The matrix $\mathbf{W}_{o,h} \in \mathbb{R}^{d_o \times d_h}$ is trained with regularized least-squares regression with Tikhonov regularization to alleviate overfitting (Tikhonov and Arsenin, 1977; Yan and Su, 2009) following the same recipe as in (Pathak et al., 2018a). The Tikhonov regularization η is optimized as a hyperparameter. Moreover, we further regularize the training procedure of RC by adding Gaussian noise in the training data. This was shown to be beneficial for both short-term performance and stabilizing the RC in long-term forecasting. For this reason, we add noise sampled from $\mathcal{N}(0, \kappa_n \sigma)$ to the training data, where σ is the standard deviation of the data and the noise level κ_n a tuned hyperparameter.

3. Comparison Metrics

The predictive performance of the models depends on the selection of model hyperparameters. For each model we perform an extensive grid search of optimal hyperparameters, reported in the Appendix. All model evaluations are mapped to a single Nvidia Tesla P100 GPU and are executed on the XC50 compute nodes of the Piz Daint supercomputer at the Swiss national supercomputing centre (CSCS). In the following we quantify the prediction accuracy of the methods in terms of the normalized root mean square error, given by

$$\text{NRMSE}(\hat{\mathbf{o}}) = \sqrt{\left\langle \frac{(\hat{\mathbf{o}} - \mathbf{o})^2}{\sigma^2} \right\rangle}, \quad (10)$$

¹Because of the nonlinearity of the tanh function, $\rho < 1$ is not necessarily required for the echo state property to hold true.

where $\hat{o} \in \mathbb{R}^{d_o}$ is the forecast at a single time-step, $o \in \mathbb{R}^{d_o}$ is the target value, and $\sigma \in \mathbb{R}^{d_o}$ is the standard deviation in time of each state component. In Equation (10), the notation $\langle \cdot \rangle$ denotes the state space average (average of all elements of a vector). To alleviate the dependency on the initial condition, we report the evolution of the NRMSE over time averaged over 100 initial conditions randomly sampled from the attractor.

Perhaps the most basic characterization of chaotic motion is through the concept of Lyapunov exponents (Ott, 2002): Considering two infinitesimally close initial conditions $U(t=0)$ and $U(t=0) + \delta U(t=0)$, their separation $|\delta U(t)|$ on average diverges exponentially in time, $|\delta U(t)|/|\delta U(t=0)| \sim \exp(\Lambda t)$, as $t \rightarrow \infty$. Note that the dimensionality of the vector displacement $\delta U(t)$ is that of the state space. In general, the Lyapunov exponent Λ depends on the orientation $(\delta U(t)/|\delta U(t)|)$ of the vector displacement $\delta U(t)$. In the $t \rightarrow \infty$ limit, the number of possible values of Λ is typically equal to the state space dimensionality. We denote these values $\Lambda_1 \geq \Lambda_2 \geq \Lambda_3 \geq \dots$ and collectively call them the Lyapunov exponent spectrum (LS) of the particular chaotic system. The Lyapunov exponent spectrum will be evaluated in Section 6.

However, we note that a special role is played by Λ_1 , and only Λ_1 , the largest Lyapunov exponent. We refer to the largest Lyapunov exponent as the Maximal Lyapunov exponent (MLE). Chaotic motion of a bounded trajectory is defined by the condition $\Lambda_1 > 0$. Importantly, if the orientation of $\delta U(t=0)$ is chosen randomly, the exponential rate at which the orbits separate is Λ_1 with probability one. This is because in order for any of the other exponents ($\Lambda_2, \Lambda_3, \dots$) to be realized, $\delta U(t=0)$ must be chosen to lie on a subspace of lower dimensionality than that of the state space; i.e., the orientation of $\delta U(t=0)$ must be chosen in an absolutely precise way never realized by random choice. Hence, the rate at which typical pairs of nearby orbits separate is Λ_1 , and $T^{\Lambda_1} = \Lambda_1^{-1}$, the “Lyapunov time”, provides a characteristic time scale for judging the quality of predictions based on the observed prediction error growth.

In order to obtain a single metric of the predictive performance of the models we compute the valid prediction time (VPT) in terms of the MLE of the system Λ_1 as

$$\text{VPT} = \frac{1}{\Lambda_1} \underset{t_f}{\operatorname{argmax}} \{ t_f \mid \text{NRMSE}(o_t) < \epsilon, \forall t \leq t_f \} \quad (11)$$

which is the largest time t_f the model forecasts the dynamics with a NRMSE error smaller than ϵ normalized with respect to Λ_1 . In the following, we set $\epsilon = 0.5$.

In order to evaluate the efficiency of the methods in capturing the long-term statistics of the dynamical system, we evaluate the mean power spectral density (power spectrum) of the state evolution over all $i \in \{1, \dots, d_o\}$ elements o_i^j of the state (since the state o_i is a vector). The power spectrum of the evolution of o_i^j is given by $PSD(f) = 20 \log_{10} (2|U(f)|) \text{ dB}$, where $U(f) = FFT(o_i^j)$ is the complex Fourier spectrum of the state evolution.

4. Forecasting Reduced Order Observable Dynamics in the Lorenz-96

The accurate long-term forecasting of the state of a deterministic chaotic dynamical system is challenging as even a minor initial error can be propagated exponentially in time due to the system dynamics even if the model predictions are perfect. A characteristic time-scale of this propagation is the Maximal Lyapunov Exponent (MLE) of the system as elaborated in Section 3. In practice, we are often interested in forecasting the evolution of an observable (that we can measure and obtain data from), which does not contain the full state information of the system. The observable dynamics are more irregular and challenging to model and forecast because of the additional loss of information.

Classical approaches to forecast the observable dynamics based on Takens seminal work (Takens, 1981), rely on reconstructing the full dynamics in a high-dimensional **phase space**. The state of the phase space is constructed by stacking delayed versions of the observed state. Assume that the state of the dynamical system is x_t , but we only have access to the less informative observable o_t . The phase space state, i.e., the **embedding state**, is given by $z_t = [o_t, o_{t-\tau}, \dots, o_{t-(d-1)\tau}]$, where the time-lag τ and the embedding dimension d are the embedding parameters. For d large enough, and in the case of deterministic nonlinear dynamical chaotic systems, there is generally a one-to-one mapping between a point in the phase space and the full state of the system and vice versa. This implies that the dynamics of the system are deterministically reconstructed in the phase space (Kantz and Schreiber, 1997) and that there exists a phase space forecasting rule $z_{t+1} = \mathcal{F}^z(z_t)$, and thus an observable forecasting rule $\hat{o}_{t+1} = \mathcal{F}^o(o_t, o_{t-\tau}, \dots, o_{t-(d-1)\tau})$.

The recurrent architectures presented in Section 2 fit to this framework, as the embedding state information can be captured in the high-dimensional hidden state h_t of the networks by processing the observable time series o_t , without having to tune the embedding parameters τ and d .

In the following, we introduce a high-dimensional dynamical system, the Lorenz-96 model and evaluate the efficiency of the methods to forecast the evolution of a reduced order observable of the state of this system. Here the observable is not the full state of the system, and the networks need to capture temporal dependencies to efficiently forecast the dynamics.

4.1. Lorenz-96 Model

The Lorenz-96 model was introduced by Edward Lorenz (Lorenz, 1995) to model the large-scale behavior of the mid-latitude atmosphere. The model describes the time evolution of an atmospheric variable that is discretized spatially over a single latitude circle modelled in the high-dimensional state $\mathbf{x} = [\mathbf{x}_0, \dots, \mathbf{x}_{J-1}]$, and is defined by the equations

$$\frac{d\mathbf{x}_j}{dt} = (\mathbf{x}_{j+1} - \mathbf{x}_{j-2})\mathbf{x}_{j-1} - \mathbf{x}_j + F, \quad (12)$$

for $j \in \{0, 1, \dots, J-1\}$, where we assume periodic boundary conditions $\mathbf{x}_{-1} = \mathbf{x}_{J-1}$, $\mathbf{x}_{-2} = \mathbf{x}_{J-2}$. In the following we consider a grid-size $J = 40$ and two different forcing regimes, $F = 8$ and $F = 10$.

We solve Equation (12) starting from a random initial condition with a Fourth Order Runge-Kutta scheme and a time-step of $\delta t = 0.01$. We run the solver up to $T = 2000$ after ensuring that transient effects are discarded ($T_{trans} = 1000$). The first half 10^5 samples are used for training and the rest for testing. For the forecasting test in the reduced order space, we construct observables of dimension $d_o \in \{35, 40\}$ by performing Singular Value Decomposition (SVD) and keeping the most energetic d_o components. The complete procedure is described in the Appendix. The 35 most energetic modes taken into account in the reduced order observable, explain approximately 98% of the total energy of the system in both $F \in \{8, 10\}$.

As a reference timescale that characterizes the chaoticity of the system we use the Lyapunov time, which is the inverse of the MLE, i.e., $T^{\Lambda_1} = 1/\Lambda_1$. The Lyapunov spectrum of the Lorenz-96 system is calculated using a standard technique based on QR decomposition (Abarbanel, 2012). This leads to $\Lambda_1 \approx 1.68$ for $F = 8$ and $\Lambda_1 \approx 2.27$ for $F = 10$.

4.2. Results on the Lorenz-96 Model

The evolution of the NRMSE of the model with the largest VPT of each architecture for $F \in \{8, 10\}$ is plotted in Figure 3 for two values of the dimension of the observable $d_o \in \{35, 40\}$, where $d_o = 40$ corresponds to full state information. Note that the observable is given by first transforming the state to its SVD modes and then keeping the d_o most energetic ones. As indicated by the slopes of the curves, models predicting the observable containing full state information ($d_o = 40$) exhibit a slightly slower NRMSE increase compared to models predicting in the reduced order state, as expected.

When the full state of the system is observed, the predictive performance of RC is superior to that of all other models. Unitary networks diverge from the attractor in both reduced order and full space in both forcing regimes $F \in \{8, 10\}$. This divergence (inability to reproduce the long-term climate of the dynamics) stems from the iterative propagation of the forecasting error. The issue has been also demonstrated in previous studies in both RC (Pathak et al., 2018b; Lu et al., 2018) and RNNs (Vlachas et al., 2018). This is because the accuracy of the network for long-term climate modeling, depends not only on how well it approximates the dynamics on the attractor locally, but also on how it behaves near the attractor, where we do not have data. As noted in Ref. (Lu et al., 2018), assuming that the network has a full Lyapunov spectrum near the attractor, if any of the Lyapunov exponents that correspond to infinitesimal perturbations transverse to the attractor phase space is positive, then the predictions of the network will eventually diverge from the attractor. Empirically, the divergence effect can also be attributed to insufficient network size (model expressiveness) and training, or attractor regions in the state space that are underrepresented in the training data (poor sampling). Even with a densely sampled attractor, during iterative forecasting in the test data, the model is propagating its own predictions, which might lead to a region near (but not on) the attractor where any positive Lyapunov exponent corresponding to infinitesimal perturbations transverse to the attractor will cause divergence.

In this work, we use 10^5 samples to densely capture the attractor. Still, RC suffers from the iterative propagation of errors leading to divergence especially in the reduced order forecasting scenario. In order to alleviate the problem, a parallel scheme for RC is proposed in (Pathak et al., 2018b) that enables training of many reservoirs locally forecasting the state. However, this method is limited to systems with local interactions in their state space. In the case we discuss here the observable obtained by singular value decomposition does not fulfill this assumption. In many systems the assumption of local interaction may not hold. GRU and LSTM show superior forecasting performance in the reduced

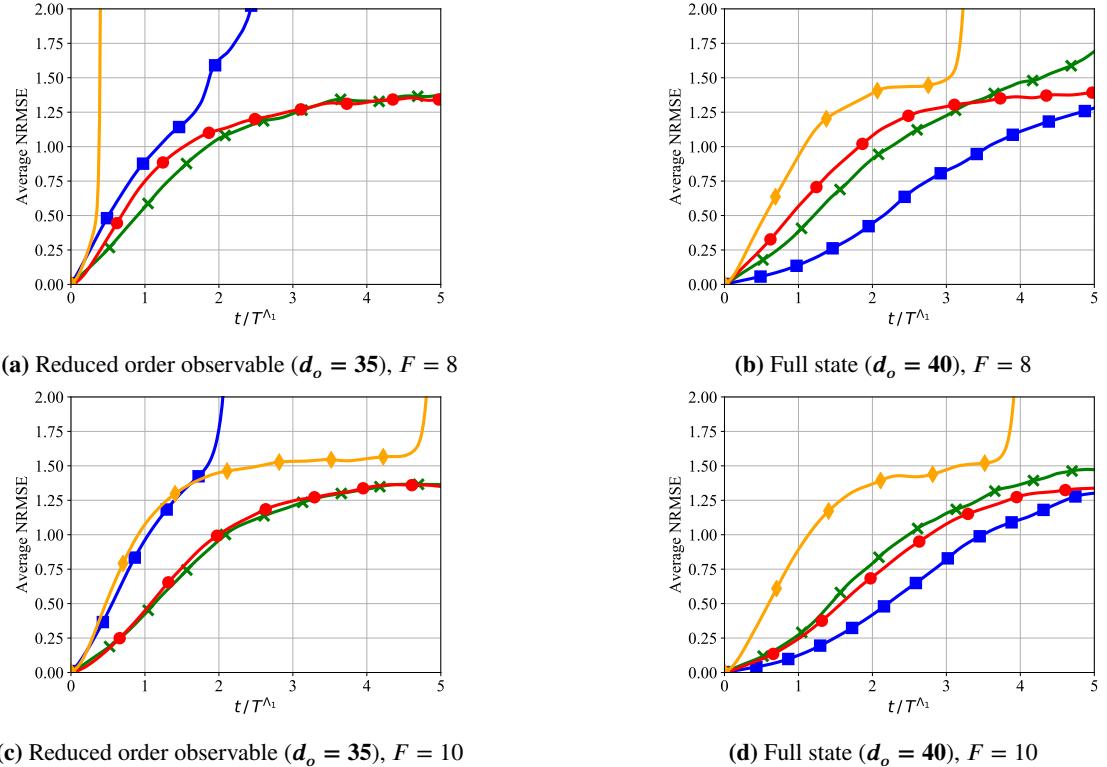


Figure 3: The evolution of the NRMSE error (average over 100 initial conditions) of the model with the highest VPT for each architecture in the Lorenz-96 with $F \in \{8, 10\}$ and $d_o \in \{35, 40\}$. Reservoir computers show remarkable predictive capabilities when the full state is observed, surpassing all other models (plots (b) and (d)). Predictions of Unitary networks diverge from the attractor in all scenarios, while iterative forecasts of RC suffer from instabilities when only partial information of a reduced order observable is available. In contrast, GRU and LSTM show stable behavior and superior performance in the reduced order scenario (plots (a) and (c)).

RC ■; GRU ✕; LSTM ●; Unitary ▲;

order scenario setting in Lorenz-96 as depicted in Figure 3a-Figure 3c. Especially in the case of $F = 10$, the LSTM and GRU models are able to predict up to 2 Lyapunov times ahead before reaching an NRMSE of $\epsilon = 1$, compared to RC and Unitary RNNs that reach this error threshold in 1 Lyapunov time. However, it should be noted that the predictive utility of all models (considering an error threshold of $\epsilon = 0.5$) is limited to one Lyapunov time when applied to reduced order data and up to two Lyapunov times in the full state.

In order to analyze the sensitivity of the VPT to the hyperparameter selection, we present violin plots in Figure 4, showing a smoothed kernel density estimate of the VPT values of all tested hyperparameter sets for $d_o = 35$ and $d_o = 45$ and $F = 8$ and $F = 10$. The horizontal markers denote the maximum, average and minimum value. Quantitative results for both $F \in \{8, 10\}$ are provided on Table 1.

In the full state scenario ($d_o = 40$) and forcing regime $F = 8$, RC shows a remarkable performance with a maximum VPT ≈ 2.31 , while GRU exhibits a max VPT of ≈ 1.34 . The LSTM has a max VPT of ≈ 0.97 , while Unitary RNNs show the lowest forecasting ability with a max VPT of ≈ 0.58 . From the violin plots in Figure 4 we notice that densities are wider at the lower part, corresponding to many models (hyperparameter sets) having much lower VPT than the maximum, emphasizing the importance of tuning the hyperparameters. Similar results are obtained for the forcing regime $F = 10$. One noticeable difference is that the LSTM exhibits a max VPT of ≈ 1.73 which is higher than that of GRU which is ≈ 1.59 . Still, the VPT of RC in the full state scenario is ≈ 2.35 which is the highest among all models.

In contrast, in the case of $d_o = 35$ where the models are forecasting on the reduced order space in the forcing regime $F = 8$, GRU is superior to all other models with a maximum VPT ≈ 0.98 compared to LSTM showing a max VPT ≈ 0.74 . LSTM shows inferior performance to GRU which we speculate may be due to insufficient hyperparameter optimization. Observing the results on $F = 10$ justifies our claim, as indeed both the GRU and the LSTM show the

Table 1

Maximum and average Valid Prediction Time (VPT) over all hyperparameter sets averaged over 100 initial conditions sampled from the testing data for each model.

Model \ Scenario	$F = 8$				$F = 10$			
	$d_o = 35$		$d_o = 40$		$d_o = 35$		$d_o = 40$	
	MAX	AVG	MAX	AVG	MAX	AVG	MAX	AVG
Unit	0.43	0.34	0.58	0.34	0.49	0.39	0.63	0.41
LSTM	0.74	0.37	0.97	0.45	1.17	0.47	1.73	0.66
GRU	0.98	0.37	1.34	0.55	1.22	0.43	1.59	0.71
RC	0.55	0.32	2.31	0.79	0.60	0.36	2.35	0.83

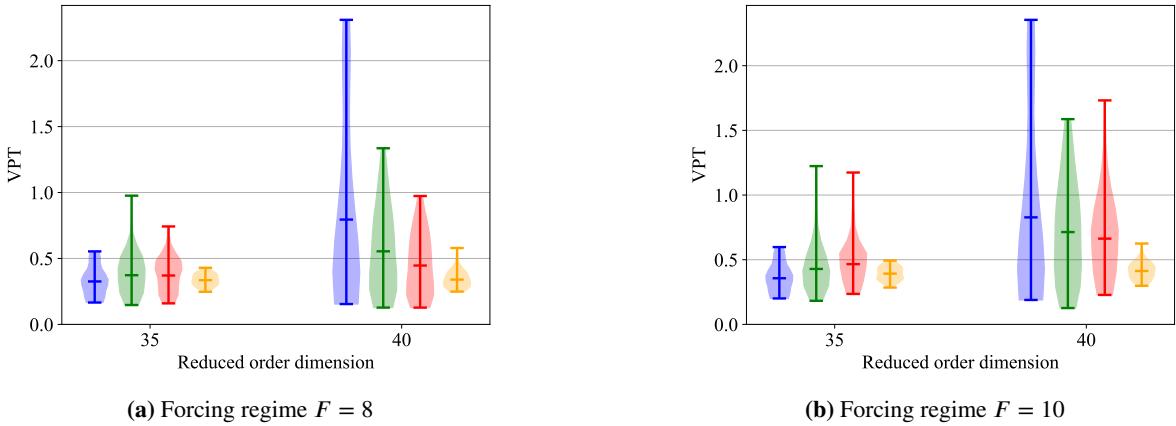


Figure 4: Violin plot showing the probability density of the VPT of all hyperparameter sets for each model for observable dimension $d_o = 35$ and $d_o = 40$ and forcing regimes (a) $F = 8$ and (b) $F = 10$ in the Lorenz-96.
RC ; GRU ; LSTM ; Unit

highest VPT values of ≈ 1.22 and ≈ 1.17 respectively. In both scenarios $F = 8$ and $F = 10$, when forecasting the reduced order space $d_o = 35$, RC shows inferior performance compared to both GRU and LSTM networks with max VPT ≈ 0.55 for $F = 8$ and ≈ 0.60 for $F = 10$. Last but not least, we observe that Unitary RNNs show the lowest forecasting ability among all models. This may not be attributed to the expressiveness of Unitary networks, but rather to the difficulty on identifying the right hyperparameters (Greff et al., 2016). In Figure 4 we observe that the violin plots in the reduced order state are much thinner at the top compared to the ones in the full state. This implies that the identification of hyperparameter sets that achieve a high VPT in the reduced order space is more challenging. This emphasizes that forecasting on the reduced order state is a more difficult task compared to the full state scenario.

In the following, we evaluate the ability of the trained networks to forecast the long-term statistics of the dynamical system. In almost all scenarios and all cases considered in this work, forecasts of Unitary RNN networks fail to remain close to the attractor and diverge. For this reason, we omit the results on these networks.

We quantify the long-term behavior in terms of the power spectrum of the predicted dynamics and its difference with the true spectrum of the testing data. In Figure 5, we plot the power spectrum of the predicted dynamics from the model (hyperparameter set) with the lowest power spectrum error for each architecture for $d_o \in \{35, 40\}$ and $F \in \{8, 10\}$ against the ground-truth spectrum computed from the testing data (dashed black line). In the full state scenario in both forcing regimes (Figure 5b, Figure 5d), all models match the true statistics in the test dataset, as the predicted power spectra match the ground-truth. These results imply that RC is a powerful predictive tool in the full order state scenario, as RC models both capture the long-term statistics and have the highest VPT among all other models analyzed in this work. However, in the case of a reduced order observable, the RC cannot match the statistics. In contrast, GRU and LSTM networks achieve superior forecasting performance while matching the long-term statistics, even at this challenging setting of a chaotic system with reduced order information.

An important aspect of machine learning models is their scalability to high-dimensional systems and their re-

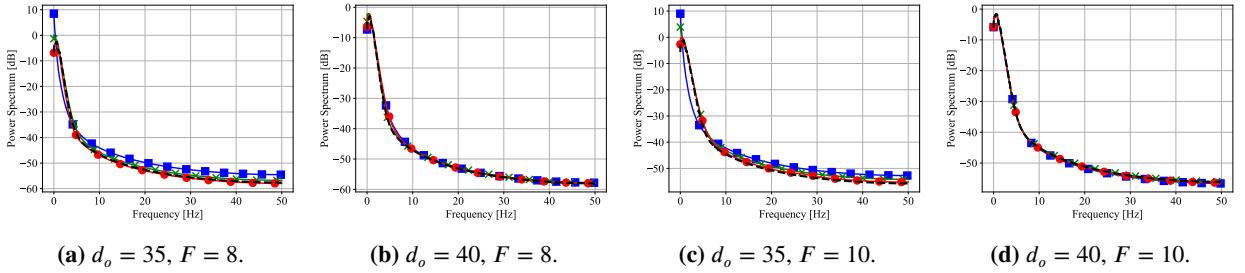


Figure 5: Predicted power spectrum of the RC, GRU, and LSTM networks with the lowest spectrum error forecasting the dynamics of an observable consisting of the SVD modes of the Lorenz-96 system with forcing $F \in \{8, 10\}$. The observable consists of the $d_o = 35$ most energetic modes or full state information $d_o = 40$. **(a)** Reduced order observable at forcing $F = 8$. **(b)** Full state observable at forcing $F = 8$. **(c)** Reduced order observable at forcing $F = 10$. **(d)** Full state observable at forcing $F = 10$.

RC ■; GRU ✖; LSTM ●; Groundtruth ---;

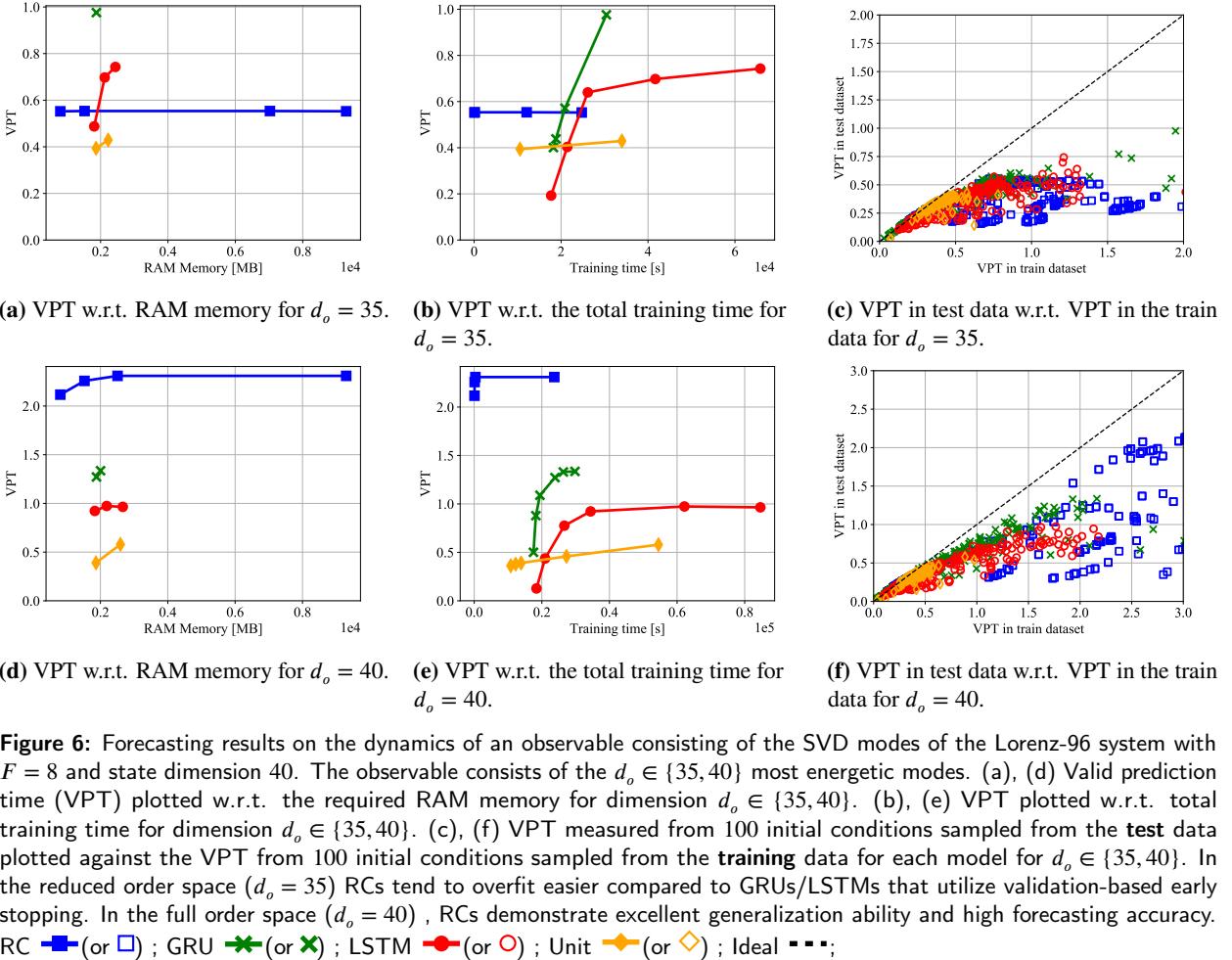
uirements in terms of training time and memory utilization. Large memory requirements and/or high training times might hinder the application of the models in challenging scenarios, like high-performance applications in climate forecasting (Kurth et al., 2018). In Figure 6a and Figure 6d, we present a Pareto front of the VPT with respect to the CPU RAM memory utilized to train the models with the highest VPT for each architecture for an input dimensions of $d_o = 35$ (reduced order) and $d_o = 40$ (full dimension) respectively. Figure 6b and Figure 6e, show the corresponding Pareto fronts of the VPT with respect to the training time. In case of the full state space ($d_o = 40$), the RC is able to achieve superior VPT with smaller memory usage and vastly smaller training time than the other methods. However, in the case of reduced order information ($d_o = 35$), the BPTT algorithms (GRU and LSTM) are superior to the RC even when the latter is provided with one order of magnitude more memory.

Due to the fact that the RNN models are **learning** the recurrent connections, they are able to reach higher VPT when forecasting in the reduced order space without the need for large models. In contrast, in RC the maximum reservoir size (imposed by computer memory limitations) may not be sufficient to capture the dynamics of high-dimensional systems with reduced order information and non-local interactions. We argue that this is the reason why the RC models do not reach the performance of GRU/LSTM trained with Back-propagation (see Figure 6a).

At the same time, letting memory limitations aside, training of RC models requires the solution of a linear system of equations $\mathbf{H}\mathbf{W}_{out}^T = \mathbf{Y}$, with $\mathbf{H} \in \mathbb{R}^{d_N \times d_h}$, $\mathbf{W}_{out}^T \in \mathbb{R}^{d_h \times d_o}$ and $\mathbf{Y} \in \mathbb{R}^{d_N \times d_o}$ (see Appendix A). The Moore-Penrose method of solving this system, scales cubically with the reservoir size as it requires the inversion of a matrix with dimensions $d_h \times d_h$. We also tried an approximate iterative method termed LSQR based on diagonalization, without any significant influence on the training time. In contrast, the training time of an RNN is very difficult to estimate a priori, as convergence of the training method depends on initialization and various other hyperparameters and are not necessarily dependent on the size. That is why we observe a greater variation of the training time of RNN models. Similar results are obtained for $F = 10$, the interested reader is referred to the appendix.

In the following, we evaluate to which extend the trained models overfit to the training data. For this reason, we measure the VPT in the training dataset and plot it against the VPT in the test dataset for every model we trained. This plot provides insight on the **generalization** error of the models. The results are shown in Figure 6c, and Figure 6f for $d_o = 35$ and $d_o = 40$. Ideally a model architecture that guards effectively against overfitting, exhibits a low generalization error, and should be represented by a point in this plot that is close to the identity line (zero generalization error). As the expressive power of a model increases, the model may fit better to the training data, but bigger models are more prone to memorizing the training dataset and overfitting (high generalization error). Such models would be represented by points on the right side of the plot. In the reduced order scenario, GRU and LSTM models lie closer to the identity line than RC models, exhibiting lower generalization errors. This is due to the validation-based early stopping routine utilized in the RNNs that guards effectively against overfitting.

We may alleviate the overfitting in RC by tuning the Tikhonov regularization parameter (η). However, this requires to rerun the training for every other combination of hyperparameters. For the four tested values $\eta \in \{10^{-3}, 10^{-4}, 10^{-5}, 10^{-6}\}$ of the Tikhonov regularization parameter the RC models tend to exhibit higher generalization error compared to the RNNs trained with BBTT. We also tested more values $\eta \in \{10^{-1}, 10^{-2}, 10^{-3}, 10^{-4}, 10^{-5}, 10^{-6}, 10^{-7}, 10^{-8}\}$,



while keeping fixed the other hyperparameters, without any observable differences in the results.

However, in the full-order scenario, the RC models achieve superior forecasting accuracy and generalization ability as clearly depicted in Figure 6f. Especially the additional regularization of the training procedure introduced by adding Gaussian noise in the data was decisive to achieve this result.

An example of an iterative forecast in the test dataset, is illustrated in Figure 7 for $F = 8$ and $d_o \in \{35, 40\}$.

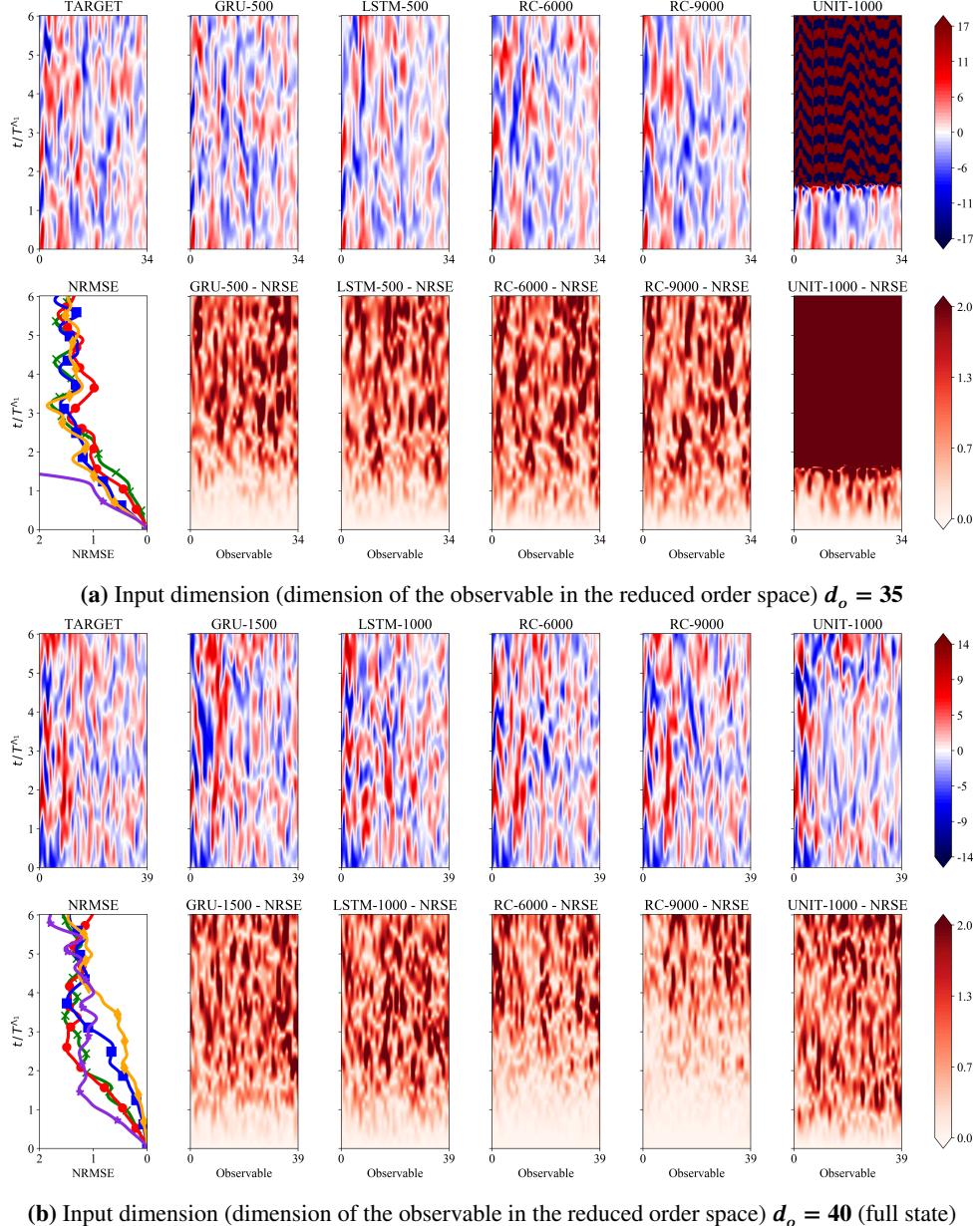
(b) Input dimension (dimension of the observable in the reduced order space) $d_o = 40$ (full state)

Figure 7: Contour plots of a spatio-temporal forecast on the SVD modes of the Lorenz-96 system with $F = 8$ in the testing dataset with GRU, LSTM, RC and a Unitary network along with the true (target) evolution and the associated NRSE contours for the reduced order observable **(a)** $d_o = 35$ and the full state **(b)** $d_o = 40$. The evolution of the component average NRSE (NMRSE) is plotted to facilitate comparison. Unitary networks suffer from propagation of forecasting error and eventually their forecasts diverge from the attractor. Forecasts in the case of an observable dimension $d_o = 40$ diverge slower as the dynamics are deterministic. In contrast, forecasting the observable with $d_o = 35$ is challenging due to both (1) sensitivity to initial condition and (2) incomplete state information that requires the capturing of temporal dependencies. In the full-state setting, RC models achieve superior forecasting accuracy compared to all other models. In the challenging reduced order scenario, LSTM and GRU networks demonstrate a stable behavior in iterative prediction and reproduce the long-term statistics of the attractor. In contrast, in the reduced order scenario RC suffer from frequent divergence (refer to the appendix).

GRU ; LSTM ; RC-6000 ; RC-9000 ; Unit

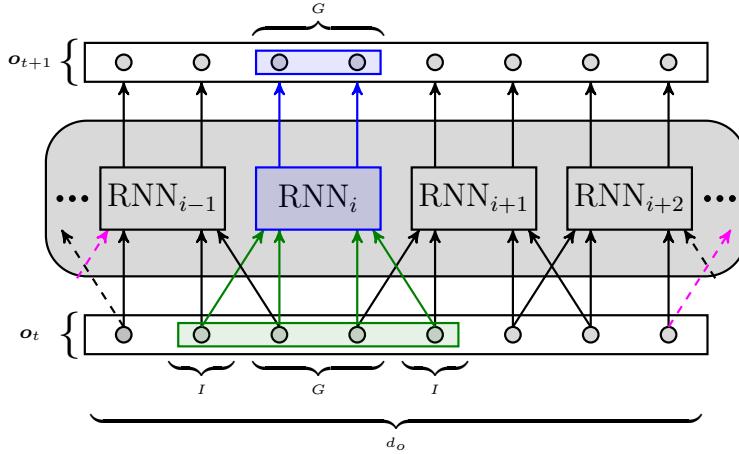


Figure 8: Illustration of the parallel architecture for a group size of $G = 2$ and an interaction length of $I = 1$. The network consists of multiple RNNs with different parameters. Each RNN is trained to forecast the evolution of G elements of the observable. Additional information of I elements from each neighboring network (left and right) are provided as additional input to capture local correlations.

5. Parallel Forecasting Leveraging Local Interactions

In spatially extended dynamical systems the state space (e.g., vorticity, velocity field, etc.) is high-dimensional (or even infinite dimensional), since an adequately fine grid is needed to resolve the relevant spatio-temporal scales of the dynamics. Even though RC and RNNs can be utilized for modeling and forecasting of these systems in the short-term, the RC and RNN methods described in Section 2 do not scale efficiently with the input dimension, i.e., as the dimensionality of the observable $o_t \in \mathbb{R}^{d_o}$ increases. Two limiting factors are the required time and RAM memory to train the model. As d_o increases, the size d_h of the reservoir network required to predict the system using only a single reservoir rises. This implies higher training times and more computational resources (RAM memory), which render the problem intractable for large values of d_o . The same applies for RNNs. More limiting factors arise by taking the process of identification of optimal model hyperparameters into account, since loading, storing and processing a very large number of large models can be computationally infeasible. However, these scaling problems for large systems can be alleviated in case the system is characterized by local state interactions or translationally invariant dynamics. In the first case, as shown in Figure 8 the modeling and forecasting task can be parallelized by employing multiple individually trained networks forecasting locally in parallel exploiting the local interactions, while, if translation invariance also applies, the individual parallel networks can be identical and training of only one will be sufficient. This parallelization concept is utilized in RC in (Pathak et al., 2018a; Parlitz and Merkowirth, 2000). The idea dates back to local delay coordinates (Parlitz and Merkowirth, 2000). The model shares ideas from convolutional RNN architectures (Sainath et al., 2015; Shi et al., 2015) designed to capture local features that are translationally invariant in image and video processing tasks. In this section, we extend this parallelization scheme to RNNs and compare the efficiency of parallel RNNs and RCs in forecasting the state dynamics of the Lorenz-96 model and Kuramoto-Sivashinsky equation discretized in a fine grid.

5.1. Parallel Architecture

Assume that the observable is $o_t \in \mathbb{R}^{d_o}$ and each element of the observable is denoted by $o_t^i \in \mathbb{R}$, $\forall i \in \{1, \dots, d_o\}$. In case of local interactions, the evolution of each element is affected by its spatially neighboring grid points. The elements o^i are split into N_g groups, each of which consisting of G spatially neighboring elements such that $d_o = GN_g$. The parallel model employs N_g RNNs, each of which is utilized to predict a spatially local region of the system observable indicated by the G group elements o^i . Each of the N_g RNNs receives G inputs o^i from the elements i it forecasts in addition to I inputs from neighboring elements on the left and on the right, where I is the interaction length. An example with $G = 2$ and $I = 1$ is illustrated in Figure 8.

During the training process, the networks can be trained independently. However, for long-term forecasting, a communication protocol has to be utilized as each network requires the predictions of neighboring networks to infer. In

the case of a homogeneous system, where the dynamics are translation invariant, the training process can be drastically reduced by utilizing one single RNN and training it on data from all groups. The weights of this RNN are then copied to all other members of the network. In the following we assume that we have no knowledge of the underlying data generating mechanism and its properties, so we assume the data is not homogeneous.

The elements of the parallel architecture are trained independently, while the MPI (Dalcín et al., 2011, 2008; Walker and Dongarra, 1996) communication protocol is utilized to communicate the elements of the interaction for long-term forecasting.

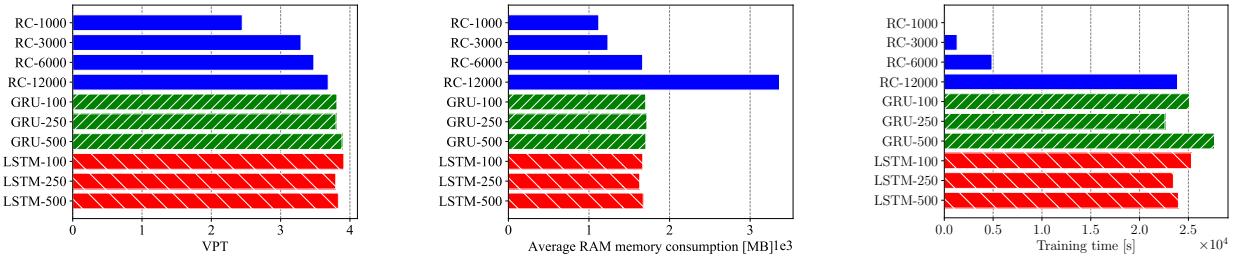
5.2. Results on the Lorenz-96

In this section, we employ the parallel architecture to forecast the state dynamics of the Lorenz-96 system explained in Section 4.1 with a state dimension of $d_o = 40$. Note that in contrast to Section 4.2, we do not construct an observable and then forecast the reduced order dynamics. Instead, we leverage the local interactions in the state space and employ an ensemble of networks forecasting the local dynamics.

Instead of a single RNN model forecasting the $d_o = 40$ dimensional global state (composed of the values of the state in the 40 grid nodes), we consider $N_g = 20$ separate RNN models, each forecasting the evolution of a $G = 2$ dimensional local state (composed of the values of the state in 2 grid nodes). In order to forecast the evolution of the local state, we take into account its interaction with $I = 4$ grid nodes on its left and on its right. The group size of the parallel models is thus $G = 2$, while the interaction length is $I = 4$. As a consequence, each model receives at its input an $2I + G = 10$ dimensional state and forecasts the evolution of a local state composed from 2 grid nodes. The size of the hidden state in RC is $d_h \in \{1000, 3000, 6000, 12000\}$. Smaller networks of size $d_h \in \{100, 250, 500\}$ are selected for GRU and LSTM. The rest of the hyperparameters are given in the appendix. Results for Unitary networks are omitted, as the identification of hyperparameters leading to stable iterative forecasting was computationally heavy and all trained models led to unstable systems that diverged after a few iterations.

In Figure 9a, we plot the VPT time of the RC and the BPTT networks. We find that RNN trained by BPTT achieve comparable predictions with RC, albeit using much smaller number hidden nodes (between 100 and 500 for BPTT vs 6000 to 12000 for RC). We remark that RC with 3000 and 6000 nodes have slightly lower VPT than GRU and LSTM but require significantly lower training times as shown in Figure 9c. At the same time, using 12000 nodes for RC implies high RAM requirements, more than 3 GB per rank, as depicted in Figure 9b.

As elaborated in Section 4.2 and depicted in Figure 3a, the VPT reached by large nonparallelized models that are forecasting the 40 SVD modes of the system is approximately 1.4. We also verified that the nonparallelized models of Section 4.1 when forecasting the 40 dimensional state containing local interactions instead of the 40 modes of SVD, reach the same predictive performance. Consequently, as expected the VPT remains the same whether we are forecasting the state or the SVD modes as the system is deterministic. By exploiting the local interactions and employing the parallel networks, the VPT is increased from ≈ 1.4 to ≈ 3.9 as shown in Figure 9a. The NRMSE error of the best performing hyperparameters is given in Figure 10a. All models are able to reproduce the climate as the reconstructed power spectrum plotted in Figure 10b matches the true one. An example of an iterative prediction with LSTM, GRU and RC models starting from an initial condition in the test dataset is provided in Figure 11.

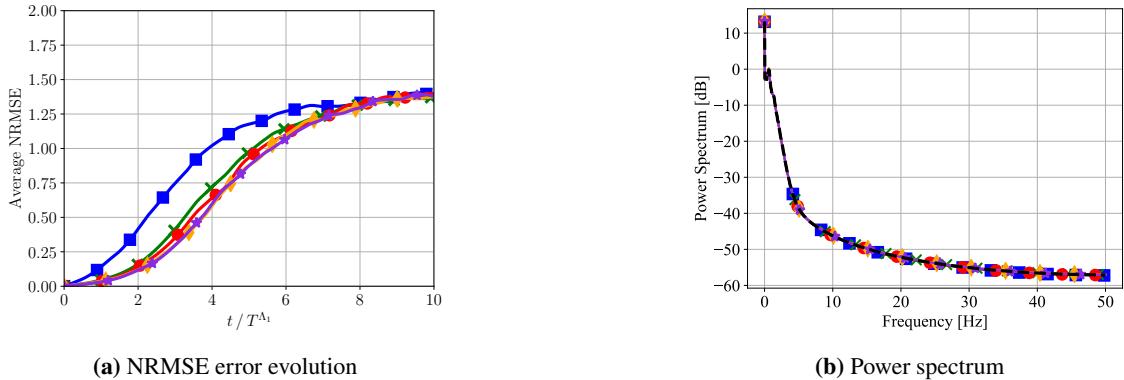


(a) Valid prediction time in the test dataset

(b) Average RAM memory requirement

(c) Training time

Figure 9: (a) Valid prediction time (VPT), (b) CPU memory utilization and (c) total training time of RNN parallel architectures with group size $G = 2$ and an interaction length $I = 4$ forecasting the dynamics of Lorenz-96 with state dimension $d_o = 40$ (full state). GRU and LSTM results do not depend significantly on network size. RC with 3000 or 6000 nodes have slightly lower VPT, but require much less training time. Increasing RC size to more than 12000 nodes was not feasible due to memory requirements.



(a) NRMSE error evolution

(b) Power spectrum

Figure 10: (a) The evolution of the NRMSE error (averaged over 100 initial conditions) of different parallel models in the Lorenz-96 with state dimension $d_o = 40$. (b) The reconstructed power spectrum. All models accurately capture the power spectrum. RCs with $d_h \in \{6000, 12000\}$ nodes are needed to match the predictive performance of an LSTM with 100 nodes.
RC-1000 ■; RC-6000 ▲; RC-12000 ●; GRU-500 ◇; LSTM-100 ▾; Groundtruth - - -;

5.3. Kuramoto-Sivashinsky

The Kuramoto-Sivashinsky (KS) equation is a nonlinear partial differential equation of fourth order that is used as a turbulence model for various phenomena. It was derived by Kuramoto in (Kuramoto, 1978) to model the chaotic behavior of the phase gradient of a slowly varying amplitude in a reaction-diffusion type medium with negative viscosity coefficient. Moreover, Sivashinsky (Sivashinsky, 1977) derived the same equations when studying the instantaneous instabilities in a laminar flame front. For our study, we restrict ourselves to the one dimensional K-S equation

$$\frac{\partial u}{\partial t} = -v \frac{\partial^4 u}{\partial x^4} - \frac{\partial^2 u}{\partial x^2} - u \frac{\partial u}{\partial x}, \quad (13)$$

on the domain $\Omega = [0, L]$ with periodic boundary conditions $u(0, t) = u(L, t)$. The dimensionless boundary size L directly affects the dimensionality of the attractor. For large values of L , the attractor dimension scales linearly with L (Manneville, 1984).

In order to spatially discretize Equation (13) we select a grid size Δx with $D = L/\Delta x + 1$ the number of nodes. Further, we denote with $u_i = u(i\Delta x)$ the value of u at node $i \in \{0, \dots, D-1\}$. In the following, we select $v = 1$, $L = 200$, $\delta t = 0.25$ and a grid of $d_o = 512$ nodes. We discretize Equation (13) and solve it using the fourth-order method for stiff PDEs introduced in (Kassam and Trefethen, 2005) up to $T = 6 \cdot 10^4$. This corresponds to $24 \cdot 10^4$ samples. The first $4 \cdot 10^4$ samples are truncated to avoid initial transients. The remaining data are divided to a training and a testing dataset of 10^5 samples each. The observable is considered to be the $d_o = 512$ dimensional state. The Lyapunov time T^A_1 of

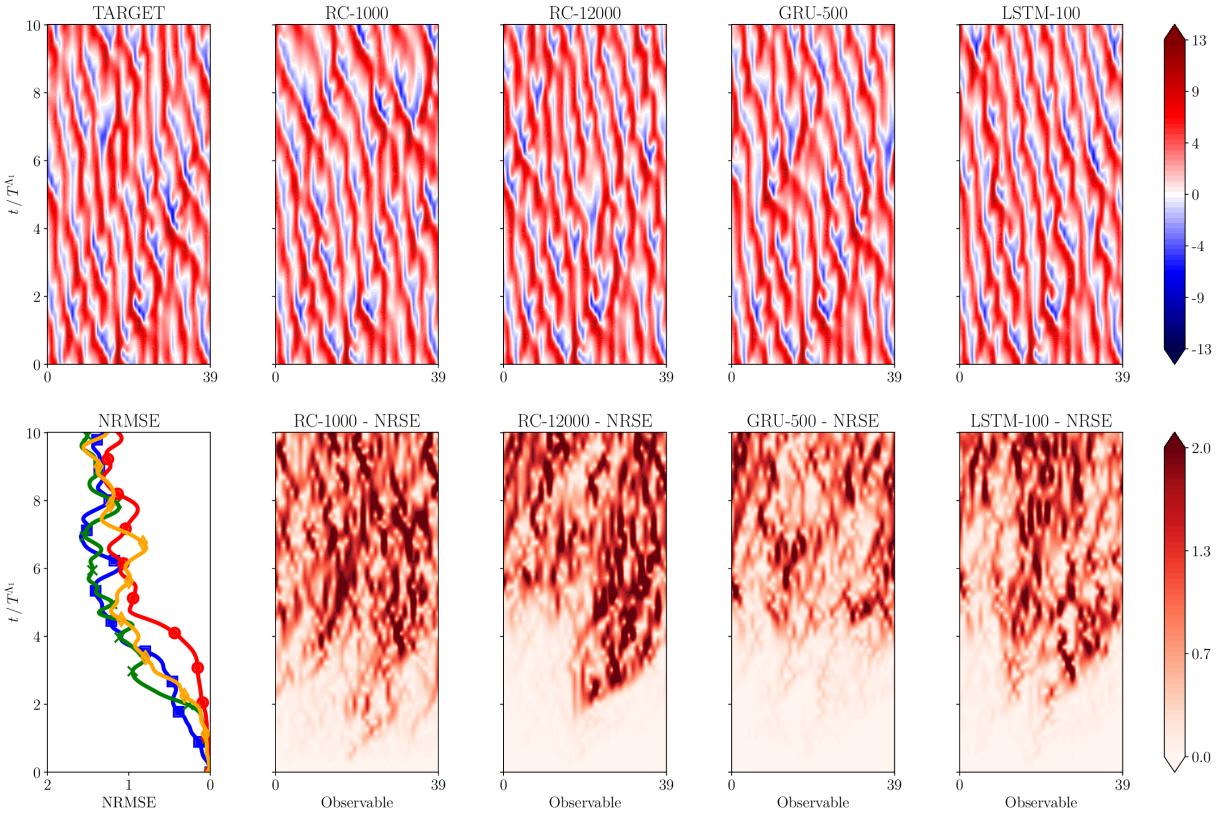


Figure 11: Contour plots of a spatio-temporal forecast in the testing dataset with parallel GRU, LSTM, and RC networks along with the true (target) evolution and the associated NRSE contours in the Lorenz-96 system with the full state as an observable $d_o = 40$. The evolution of the component average NRSE (NMRSE) is plotted to facilitate comparison.
 RC-1000 ■; RC-12000 ▲; GRU-500 ●; LSTM-100 □;

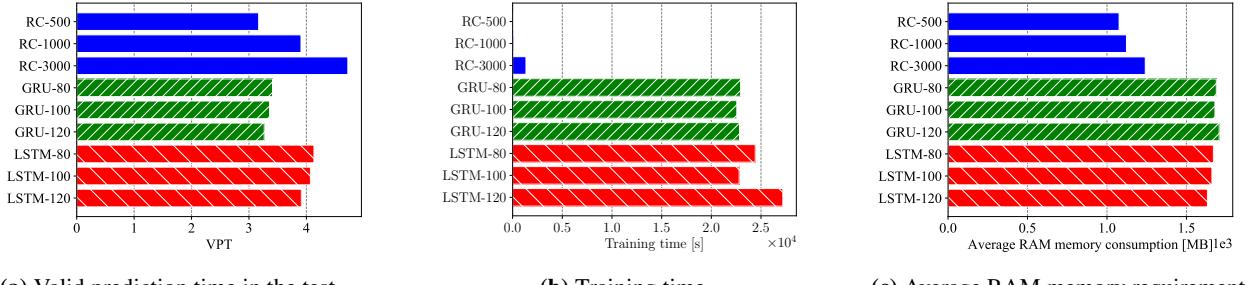
the system (see Section 3) is utilized as a reference timescale. We approximate it with the method of Pathak (Pathak et al., 2018a) for $L = 200$ and it is found to be $T^{\Lambda_1} \approx 0.094$.

5.4. Results on the Kuramoto-Sivashinsky Equation

In this section, we present the results of the parallel models in the Kuramoto-Sivashinsky equation. The full system state is used as an observable, i.e., $d_o = 512$. The group-size of the parallel models is set to $G = 8$, while the interaction length is $I = 8$. The total number of groups is $N_g = 64$. Each member forecasts the evolution of 8 state components, receiving at the input 24 components in total. The size of the reservoir in RC is $d_h \in \{500, 1000, 3000\}$. For GRU and LSTM networks we vary $d_h \in \{100, 250, 500\}$. The rest of the hyperparameters are given in the appendix. Results on Unitary networks are omitted, as the configurations tried in this work led to unstable models diverging after a few time-steps in the iterative forecasting procedure.

The results are summed up in the bar-plots in Figure 12. In Figure 12a, we plot the VPT time of the models. LSTM models reach VPTs of ≈ 4 , while GRU show an inferior predictive performance with VPTs of ≈ 3.5 . An RC with $d_h = 500$ reaches a VPT of ≈ 3.2 , and an RC with 1000 modes reaches the VPT of LSTM models with a VPT of ≈ 3.9 . Increasing the reservoir capacity of the RC to $d_h = 3000$ leads to a model exhibiting a VPT of ≈ 4.8 . In this case, the large RC model shows slightly superior performance to GRU/LSTM. The low performance of GRU models can be attributed to the fact that in the parallel setting the probability that any RNN may converge to bad local minima rises, with a detrimental effect on the total predictive performance of the parallel ensemble. In case of spatially translational invariant systems, we could alleviate this problem by using one single network. Still, training the single network to data from all spatial locations would be expensive.

As depicted in Figure 12, the reservoir size of 3000 is enough for RC to reach and surpass the predictive performance



(a) Valid prediction time in the test dataset
(b) Training time
(c) Average RAM memory requirement

Figure 12: **(a)** Valid prediction time (VPT), **(b)** total training time, and **(c)** CPU memory utilization of parallel RNN architectures with group size $G = 8$ and an interaction length $I = 8$ forecasting the dynamics of Kuramoto-Sivashinsky equation with state dimension $d_o = 512$.

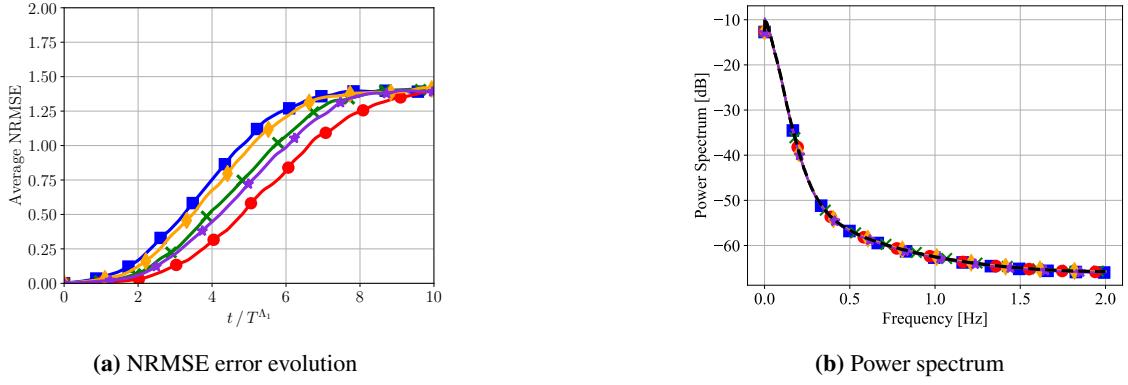


Figure 13: **(a)** The evolution of the NRMSE error (averaged over 100 initial conditions) of different parallel models in the Kuramoto-Sivashinsky equation with state dimension $d_o = 512$. **(b)** The power spectrum. All models capture the statistics of the system.

RC-500 ■; RC-1000 ✕; RC-3000 ●; GRU-80 ▲; LSTM-80 ▾; Groundtruth - - -;

of RNNs utilizing a similar amount of RAM memory and a much lower amount of training time as illustrated in Figure 12b.

The evolution of the NRMSE is given in Figure 13a. The predictive performance of a small LSTM network with 80 hidden units, matches that of a large RC with 1000 hidden units. In Figure 13b, the power spectrum of the predicted state dynamics of each model is plotted along with the true spectrum of the equations. The three models captured successfully the statistics of the system, as we observe a very good match. An example of an iterative prediction with LSTM, GRU and RC models starting from an initial condition in the test dataset is provided in Figure 14.

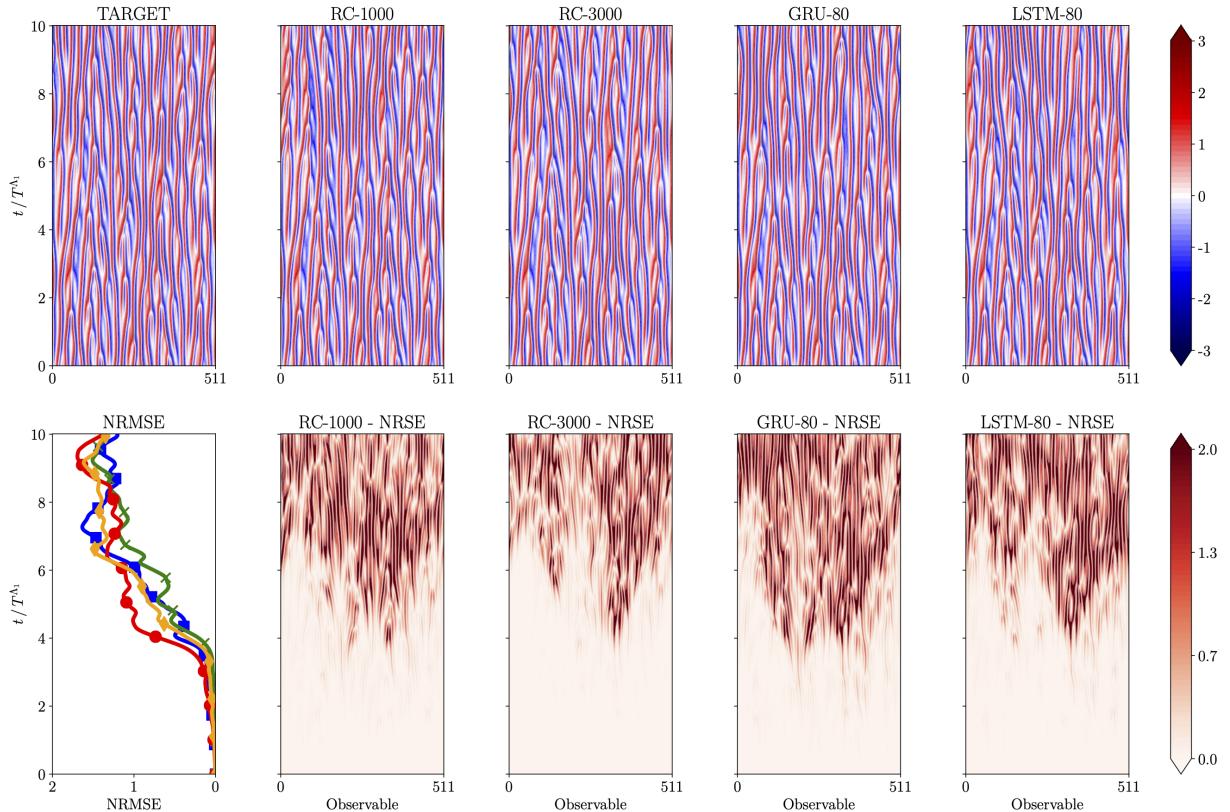


Figure 14: Contour plots of a spatio-temporal forecast starting from an initial condition in the testing dataset with parallel GRU, LSTM, and RC networks along with the true (target) evolution and the associated NRSE contours in the Kuramoto-Sivashinsky equation with the full state as an observable $d_o = 512$. The evolution of the component average NRSE (NMRSE) is plotted to facilitate comparison.

RC-1000 ; RC-3000 ; GRU-80 ; LSTM-80 ; Groundtruth .

6. Calculation of Lyapunov Exponents in the Kuramoto-Sivashinsky Equation

The recurrent models utilized in this study can be used as surrogate models to calculate the Lyapunov exponents (LEs) of a dynamical system relying only on experimental time-series data. The LEs characterize the rate of separation if positive (or convergence if negative) of trajectories that are initialized infinitesimally close in the phase space. They can provide an estimate of the attractor dimension according to the Kaplan-Yorke formula (Kaplan and Yorke, 1979). Early efforts to solve the challenging problem of data-driven LE identification led to local approaches (Wolf et al., 1985; Sano and Sawada, 1985) that are computationally inexpensive at the cost of requiring a large amount of data. Other approaches fit a global model to the data (Maus and Sprott, 2013) and calculate the LE spectrum using the Jacobian algorithm. These approaches were applied to low-order systems.

A recent machine learning approach utilizes deep convolutional neural networks for LE and chaos identification, without estimation of the dynamics (Makarenko, 2018). An RC-RNN approach capable of uncovering the whole LE spectrum in high-dimensional dynamical systems is proposed in (Pathak et al., 2018a). The method is based on the training of a surrogate RC model to forecast the evolution of the state dynamics, and the calculation of the Lyapunov spectrum of the hidden state of this surrogate model. The RC method demonstrates excellent agreement for all positive Lyapunov exponents and many of the negative exponents for the KS equation with $L = 60$ (Pathak et al., 2018a), alleviating the problem of spurious Lyapunov exponents of delay coordinate embeddings (Dechert and Gençay, 1996). We build on top of this work and demonstrate that a GRU trained with BPTT can reconstruct the Lyapunov spectrum accurately with lower error for all positive Lyapunov exponents at the cost of higher training times.

The Lyapunov spectrum of the KS equation is computed by solving the KS equations in the Fourier space with a fourth order time-stepping method called ETDRK4 (Kassam and Trefethen, 2005) and utilizing a QR decomposition approach as in (Pathak et al., 2018a). The Lyapunov spectrum of the RNN and RC surrogate models is computed based on the Jacobian of the hidden state dynamics along a reference trajectory, while Gram-Schmidt orthonormalization is utilized to alleviate numerical divergence. We employ a GRU-RNN over LSTM-RNN, due to the fact that the latter has two coupled hidden states, rendering the computation of the Lyapunov spectrum mathematically more involved and computationally more expensive. The interested reader can refer to the Appendix for the details of the method. The identified maximum LE is $\Lambda_1 \approx 0.08844$. In this work, a large RC with $d_h = 9000$ nodes is employed for LS calculation in the Kuramoto-Sivashinsky equation with parameter $L = 60$ and $D = 128$ grid points as in (Pathak et al., 2018a). The largest LE identified in this case is $\Lambda_1 \approx 0.08378$ leading to a relative error of 5.3%. In order to evaluate the efficiency of RNNs, we utilize a large GRU with $d_h = 2000$ hidden units. An iterative RNN roll-out of $N = 10^4$ total time-steps was needed to achieve convergence of the spectrum. The largest Lyapunov exponent identified by the GRU is $\Lambda_1 \approx 0.0849$ reducing the error to $\approx 4\%$. Both surrogate models identify the correct Kaplan-Yorke dimension $KY \approx 15$, which is the largest LE such that $\sum_i \Lambda_i > 0$.

The first 26 Lyapunov exponents computed the GRU, RC as well as using the true equations of the Kuramoto-Sivashinsky are plotted in Figure 15. We observe a good match between the positive Lyapunov exponents by both GRU and RC surrogates. The positive Lyapunov exponents are characteristic of chaotic behavior. However, the zero Lyapunov exponents Λ_7 and Λ_8 cannot be captured either with RC or with RNN surrogates. This is also observed in RC in (Pathak et al., 2018a), and apparently the GRU surrogate employed in this work does not alleviate the problem. In Figure 15b, we augment the RC and the GRU spectrum with these two additional exponents to illustrate that there is an excellent agreement between the true LE and the augmented LS identified by the surrogate models. The relative and absolute errors in the spectrum calculation are illustrated in Figure 16. After augmenting with the two zero LE, we get a mean absolute error of 0.012 for RC and 0.008 for GRU. The mean relative error is 0.23 for RC, and 0.22 for GRU. As a conclusion, GRU in par with RC networks can be used to replicate the chaotic behavior of a reference system and calculate the Lyapunov spectrum accurately.

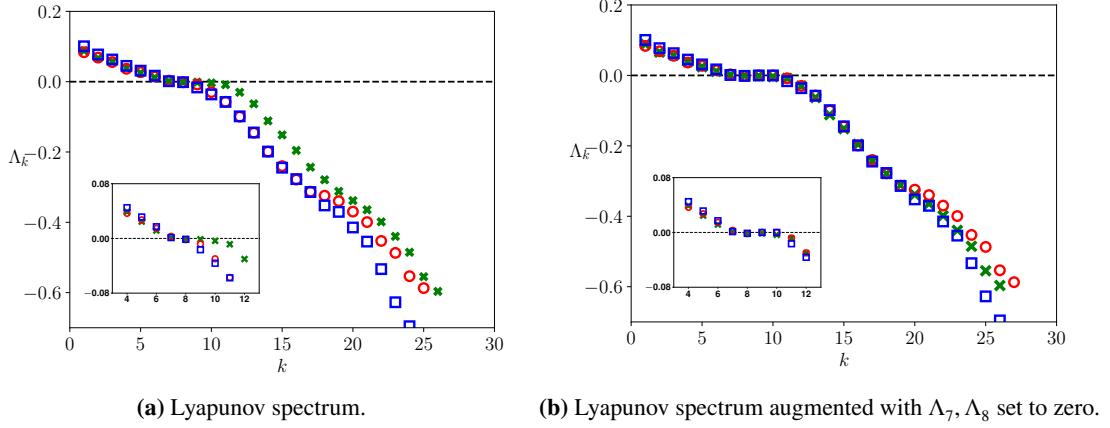


Figure 15: (a) Estimated Lyapunov exponents Λ_k of the KS equation with $L = 60$. The true Lyapunov exponents are illustrated with green crosses, red circles are calculated with the RC surrogate, while the blue rectangles with GRU. In (b) we augment the computed spectrums with the two zero Lyapunov exponents Λ_7, Λ_8 . Inset plots zoom in the zero crossing regions.

True X ; RC \circ ; GRU \square ;

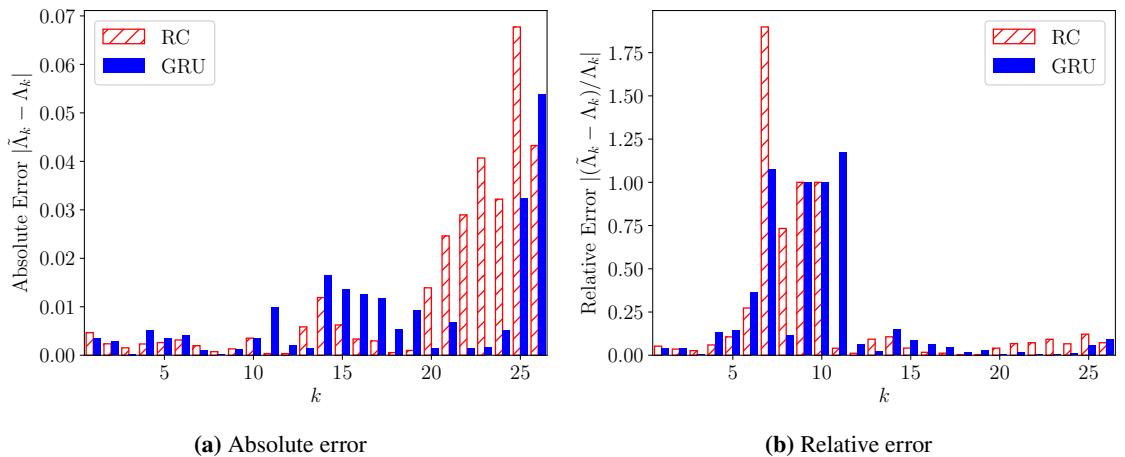


Figure 16: (a) Absolute and (b) Relative error of the LE spectrum of the KS equation with $L = 60$. The LE spectrum identified using the GRU shows a better agreement with the spectrum identified by the Kuramoto-Sivashinsky equations.

7. Conclusions

In this work, we employed several variants of recurrent neural networks and reservoir computing to forecast the dynamics of chaotic systems. We present a comparative study based on their efficiency in capturing temporal dependencies, evaluate how they scale to systems with high-dimensional state space, and how to guard against overfitting. We highlight the advantages and limitations of these methods and elucidate their applicability to forecasting spatiotemporal dynamics.

We considered three different types of RNN cells that alleviate the well-known vanishing and exploding gradient problem in Back-propagation through time training (BPTT), namely LSTM, GRU and Unitary cells. We benchmarked these networks against reservoir computers with random hidden to hidden connection weights, whose training procedure amounts to least square regression on the output weights.

The efficiency of the models in capturing temporal dependencies in the reduced order state space is evaluated on the Lorenz-96 system in two different forcing regimes $F = \{8, 10\}$, by constructing a reduced order observable using Singular Value Decomposition (SVD) and keeping the most energetic modes. Even though this forecasting task is challenging due to (1) chaotic dynamics and (2) reduced order information, LSTM and GRU show superior forecasting ability to RC utilizing similar amounts of memory at the cost of higher training times. GRU and LSTM models demonstrate stable behavior in the iterative forecasting procedure in the sense that the forecasting error usually does not diverge, in stark contrast to RC and Unitary forecasts. Large RC models tend to overfit easier than LSTM/GRU models, as the latter are utilizing validation-based early stopping and regularization techniques (e.g., Zoneout, Dropout) that guard against overfitting which are not directly applicable to RC. Validation in RC amounts to tuning an additional hyperparameter, the Tikhonov regularization. However, RC shows excellent forecasting efficiency when the full state of the system is observed, outperforming all other models by a wide margin, while also reproducing the frequency spectrum of the underlying dynamics.

RNNs and RC both suffer from scalability problems in high-dimensional systems, as the required hidden state size d_h to capture the high-dimensional dynamics can become prohibitively large especially with respect to the computational expense of training. In order to scale the models to high-dimensional systems we employ a parallelization scheme that exploits the local interactions in the state of a dynamical system. As a reference, we consider the Lorenz-96 system and the Kuramoto-Sivashinsky equation, and we train parallel RC, GRU, and LSTM models of various sizes. Iterative forecasting with parallel Unitary models diverged after a few timesteps in both systems. Parallel GRU, LSTM and RC networks reproduced the long-term attractor climate, as well as the power spectrum of the state of the Lorenz-96 and the Kuramoto-Sivashinsky equation matched with the predicted ones.

In the Lorenz-96 and the Kuramoto-Sivashinsky equation, the parallel LSTM and GRU models exhibited similar predictive performance compared to the parallel RC. The memory requirements of the models are comparable. RC networks require large reservoirs with 1000 – 6000 nodes per member to reach the predictive performance of parallel GRU/LSTM with a few hundred nodes, but their training time is significantly lower.

Last but not least, we evaluated and compared the efficiency of GRU and RC networks in capturing the Lyapunov spectrum of the KS equation. The positive Lyapunov exponents are captured accurately by both RC and GRU. Both networks cannot reproduce two zero LEs Λ_7 and Λ_8 . When these two are discarded from the spectrum, GRU and RC networks show comparable accuracy in terms of relative and absolute error of the Lyapunov spectrum.

Further investigation on the underlying reasons why the RNNs and RC cannot capture the zero Lyapunov exponents is a matter of ongoing work. Another interesting direction could include studying the memory capacity of the networks. This could offer more insight into which architecture and training method is appropriate for tasks with long-term dependencies. Moreover, we plan to investigate a coupling of the two training approaches to further improve their predictive performance, for example a network can utilize both RC and LSTM computers to identify the input to output mapping. While the weights of the RC are initialized randomly to satisfy the echo state property, the output weights alongside with the LSTM weights can be optimized by back-propagation. This approach, although more costly, might achieve higher efficiency, as the LSTM is used as a residual model correcting the error that a plain RC would have.

Although we considered a batched version of RC training to reduce the memory requirements, further research is needed to alleviate the memory burden associated with the matrix inversion (see Appendix A, Equation (15)) and the numerical problems associated with the eigenvalue decomposition of the sparse weight matrix.

Further directions could be the initialization of RNN weights with RC based heuristics based on the echo state property and fine-tuning with BPTT. This is possible for the plain cell RNN, where the heuristics are directly applicable. However, in more complex architectures like the LSTM or the GRU, more sophisticated initialization schemes that

ensure some form of echo state property have to be investigated. The computational cost of training networks of the size of RC with back-propagation is also challenging. Another interesting topic, is to analyze the influence of the amount of training data, and system size on the predictive efficiency of the methods, under the lens of the recently discovered “double descent” curve Belkin et al. (2019), supporting that over-parametrized networks (increasing model capacity beyond the point of interpolation) results in improved generalization.

One topic not covered in this work, is invertibility of the models, when forecasting the full state dynamics. Non-invertible models like the RNNs trained in this work, may suffer from spurious dynamics not present the training data and the underlying governing equations (Gicquel et al., 1998; Frouzakis et al., 1997). Invertible RNNs may constitute a promising alternative to further improve accurate short-term prediction and capturing of the long-term dynamics.

In conclusion, recurrent neural networks for data-driven surrogate modeling and forecasting of chaotic systems can efficiently be used to model high-dimensional dynamical systems, can be parallelized alleviating scaling problems and constitute a promising research subject that requires further analysis.

8. Data and Code

The code and data will be available upon publication in the following link <https://github.com/pvlachas/RNN-RC-Chaos> to assist reproducibility of the results. The software was written in Python utilizing Tensorflow (Abadi et al., 2016) and Pytorch (Paszke et al., 2017) for automatic differentiation and the design of the neural network architectures.

9. Acknowledgments

We thank Fabian Wermelinger and Ivica Kičić for helpful discussions and Guido Novati for valuable feedback on this manuscript. Moreover, we would like to acknowledge the time and effort of the Action Editor Herbert Jaeger along with two anonymous reviewers whose insightful and thorough feedback led to substantial improvements of the manuscript. TPS has been supported through the ARO-MURI grant W911NF-17-1-0306. This work has been supported at UMD by DARPA under grant number DMS51813027. We are also thankful to the Swiss National Supercomputing Centre (CSCS) providing the necessary computational resources under Projects s929.

CRediT authorship contribution statement

P. R. Vlachas: Conceptualization, Methodology, Software, Data Curation, Visualization, Writing - Original Draft. **J. Pathak:** Conceptualization, Methodology, Software, Data Curation, Writing - Review & Editing. **B. R. Hunt:** Conceptualization, Supervision, Writing - Review & Editing. **T. P. Sapsis:** Conceptualization, Supervision, Writing - Review & Editing. **M. Girvan:** Conceptualization, Supervision. **E. Ott:** Conceptualization, Supervision, Writing - Review & Editing. **P. Koumoutsakos:** Conceptualization, Methodology, Supervision, Writing - Review & Editing.

References

- Abadi, M., Barham, P., Chen, J., Chen, Z., Davis, A., Dean, J., Devin, M., Ghemawat, S., Irving, G., Isard, M., Kudlur, M., Levenberg, J., Monga, R., Moore, S., Murray, D.G., Steiner, B., Tucker, P., Vasudevan, V., Warden, P., Wicke, M., Yu, Y., Zheng, X., 2016. Tensorflow: A system for large-scale machine learning, in: 12th USENIX Symposium on Operating Systems Design and Implementation (OSDI 16), pp. 265–283. URL: <https://www.usenix.org/system/files/conference/osdi16/osdi16-abadi.pdf>.
- Abarbanel, H., 2012. Analysis of observed chaotic data. Springer Science & Business Media.
- Aksamit, N.O., Sapsis, T.P., Haller, G., 2019. Machine-learning ocean dynamics from lagrangian drifter trajectories. arXiv preprint arXiv:1909.12895
- Antonik, P., Haelterman, M., Massar, S., 2017. Brain-inspired photonic signal processor for generating periodic patterns and emulating chaotic systems. Phys. Rev. Applied 7, 054014. URL: <https://link.aps.org/doi/10.1103/PhysRevApplied.7.054014>, doi:10.1103/PhysRevApplied.7.054014.
- Arjovsky, M., Shah, A., Bengio, Y., 2016. Unitary evolution recurrent neural networks, in: Proceedings of the 33rd International Conference on International Conference on Machine Learning - Volume 48, JMLR.org. pp. 1120–1128. URL: <http://dl.acm.org/citation.cfm?id=3045390.3045509>.
- Belkin, M., Hsu, D., Ma, S., Mandal, S., 2019. Reconciling modern machine-learning practice and the classical bias–variance trade-off. Proceedings of the National Academy of Sciences 116, 15849–15854.
- Bengio, Y., Simard, P., Frasconi, P., 1994. Learning long-term dependencies with gradient descent is difficult. Trans. Neur. Netw. 5, 157–166. URL: <http://dx.doi.org/10.1109/72.279181>, doi:10.1109/72.279181.
- Bianchi, F.M., Maiorino, E., Kampffmeyer, M.C., Rizzi, A., Jenssen, R., 2017. An overview and comparative analysis of recurrent neural networks for short term load forecasting. CoRR abs/1705.04378. URL: <http://arxiv.org/abs/1705.04378>, arXiv:1705.04378.

- Bradley, E., Kantz, H., 2015. Nonlinear time-series analysis revisited. *Chaos: An Interdisciplinary Journal of Nonlinear Science* 25, 097610. doi:10.1063/1.4917289.
- Brockman, G., Cheung, V., Pettersson, L., Schneider, J., Schulman, J., Tang, J., Zaremba, W., 2016. Openai gym. arXiv preprint arXiv:1606.01540.
- Brunton, S.L., Noack, B.R., Koumoutsakos, P., 2020. Machine learning for fluid mechanics. *Annual Review of Fluid Mechanics* 52, 477–508.
- Cao, L., Hong, Y., Fang, H., He, G., 1995. Predicting chaotic time series with wavelet networks. *Physica D: Nonlinear Phenomena* 85, 225 – 238. URL: <http://www.sciencedirect.com/science/article/pii/0167278995001190>, doi:[https://doi.org/10.1016/0167-2789\(95\)00119-0](https://doi.org/10.1016/0167-2789(95)00119-0).
- Cho, K., van Merriënboer, B., Gulcehre, C., Bahdanau, D., Bougares, F., Schwenk, H., Bengio, Y., 2014. Learning phrase representations using rnn encoder–decoder for statistical machine translation, in: Proceedings of the 2014 Conference on Empirical Methods in Natural Language Processing (EMNLP), Association for Computational Linguistics. pp. 1724–1734. URL: <http://aclweb.org/anthology/D14-1179>, doi:10.3115/v1/D14-1179.
- Chung, J., Gulcehre, C., Cho, K., Bengio, Y., 2014. Empirical evaluation of gated recurrent neural networks on sequence modeling, in: NIPS 2014 Workshop on Deep Learning, December 2014.
- Dalcín, L., Paz, R., Storti, M., D'Elía, J., 2008. Mpi for python: Performance improvements and mpi-2 extensions. *Journal of Parallel and Distributed Computing* 68, 655–662.
- Dalcín, L.D., Paz, R.R., Kler, P.A., Cosimo, A., 2011. Parallel distributed computing using python. *Advances in Water Resources* 34, 1124–1139.
- Dechert, W.D., Gençay, R., 1996. The topological invariance of Lyapunov exponents in embedded dynamics. *Physica D Nonlinear Phenomena* 90, 40–55. doi:10.1016/0167-2789(95)00225-1.
- Dong, D., Wu, H., He, W., Yu, D., Wang, H., 2015. Multi-task learning for multiple language translation, in: Proceedings of the 53rd Annual Meeting of the Association for Computational Linguistics and the 7th International Joint Conference on Natural Language Processing (Volume 1: Long Papers), pp. 1723–1732.
- Dreyfus, S., 1962. The numerical solution of variational problems. *Journal of Mathematical Analysis and Applications* 5, 30 – 45.
- Elman, J.L., 1990. Finding structure in time. *COGNITIVE SCIENCE* 14, 179–211.
- Esteva, A., Kuprel, B., Novoa, R.A., Ko, J., Swetter, S.M., Blau, H.M., Thrun, S., 2017. Dermatologist-level classification of skin cancer with deep neural networks. *Nature* 542, 115.
- Faller, W.E., Schreck, S.J., 1997. Unsteady fluid mechanics applications of neural networks. *Journal of Aircraft* 34, 48–55. doi:10.2514/2.2134.
- Frouzakis, C.E., Gardini, L., Kevrekidis, I.G., Milleroux, G., Mira, C., 1997. On some properties of invariant sets of two-dimensional noninvertible maps. *International Journal of Bifurcation and Chaos* 7, 1167–1194.
- Gal, Y., Ghahramani, Z., 2016. A theoretically grounded application of dropout in recurrent neural networks, in: Lee, D.D., Sugiyama, M., Luxburg, U.V., Guyon, I., Garnett, R. (Eds.), *Advances in Neural Information Processing Systems 29*. Curran Associates, Inc., pp. 1019–1027. URL: <http://papers.nips.cc/paper/6241-a-theoretically-grounded-application-of-dropout-in-recurrent-neural-networks.pdf>.
- Gicquel, N., Anderson, J., Kevrekidis, I., 1998. Noninvertibility and resonance in discrete-time neural networks for time-series processing. *Physics Letters A* 238, 8–18.
- Glorot, X., Bengio, Y., 2010. Understanding the difficulty of training deep feedforward neural networks, in: Proceedings of the thirteenth international conference on artificial intelligence and statistics, pp. 249–256.
- Gneiting, T., Raftery, A.E., 2005. Weather forecasting with ensemble methods. *Science* 310, 248–249.
- Gonon, L., Ortega, J.P., 2019. Reservoir computing universality with stochastic inputs. *IEEE transactions on neural networks and learning systems*.
- Goodfellow, I., Bengio, Y., Courville, A., 2016. Deep Learning. MIT Press.
- Graves, A., Jaitly, N., 2014. Towards end-to-end speech recognition with recurrent neural networks, in: International conference on machine learning, pp. 1764–1772.
- Greff, K., Srivastava, R.K., Koutník, J., Steunebrink, B.R., Schmidhuber, J., 2016. Lstm: A search space odyssey. *IEEE transactions on neural networks and learning systems* 28, 2222–2232.
- Grigoryeva, L., Ortega, J.P., 2018. Echo state networks are universal. *Neural Networks* 108, 495–508.
- Ha, D., Schmidhuber, J., 2018. World models. arXiv preprint arXiv:1803.10122.
- Hassabis, D., Kumaran, D., Summerfield, C., Botvinick, M., 2017. Neuroscience-inspired artificial intelligence. *Neuron* 95, 245–258.
- Haynes, N.D., Soriano, M.C., Rosin, D.P., Fischer, I., Gauthier, D.J., 2015. Reservoir computing with a single time-delay autonomous boolean node. *Phys. Rev. E* 91, 020801. URL: <https://link.aps.org/doi/10.1103/PhysRevE.91.020801>, doi:10.1103/PhysRevE.91.020801.
- He, K., Zhang, X., Ren, S., Sun, J., 2016. Deep residual learning for image recognition. *2016 IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, 770–778.
- Hochreiter, S., 1998. The vanishing gradient problem during learning recurrent neural nets and problem solutions. *International Journal of Uncertainty, Fuzziness and Knowledge-Based Systems* 6, 107–116. doi:10.1142/S0218488598000094.
- Hochreiter, S., Schmidhuber, J., 1997. Long short-term memory. *Neural Comput.* 9, 1735–1780. URL: <http://dx.doi.org/10.1162/neco.1997.9.8.1735>, doi:10.1162/neco.1997.9.8.1735.
- Jaeger, H., Haas, H., 2004. Harnessing nonlinearity: Predicting chaotic systems and saving energy in wireless communication. *Science* 304, 78–80. URL: <http://science.sciencemag.org/content/304/5667/78>, doi:10.1126/science.1091277, arXiv:<http://science.sciencemag.org/content/304/5667/78.full.pdf>.
- Jiang, J., Lai, Y.C., 2019. Model-free prediction of spatiotemporal dynamical systems with recurrent neural networks: Role of network spectral radius. *Phys. Rev. Research* 1, 033056. URL: <https://link.aps.org/doi/10.1103/PhysRevResearch.1.033056>, doi:10.1103/PhysRevResearch.1.033056.
- Jing, L., Shen, Y., Dubcek, T., Peurifoy, J., Skirla, S.A., LeCun, Y., Tegmark, M., Soljacic, M., 2017. Tunable efficient unitary neural networks (EUNN) and their application to rnns, in: ICML, PMLR. pp. 1733–1741.
- Jozefowicz, R., Zaremba, W., Sutskever, I., 2015. An empirical exploration of recurrent network architectures, in: Proceedings of the 32Nd International Conference on International Conference on Machine Learning - Volume 37, JMLR.org. pp. 2342–2350. URL: <http://dl.acm.org>.

- [org/citation.cfm?id=3045118.3045367.](https://www.ncbi.nlm.nih.gov/citation.cfm?id=3045118.3045367)
- Kantz, H., Schreiber, T., 1997. Nonlinear Time Series Analysis. Cambridge University Press, New York, NY, USA.
- Kaplan, J.L., Yorke, J.A., 1979. Chaotic behavior of multidimensional difference equations, in: Peitgen, H.O., Walther, H.O. (Eds.), Functional Differential Equations and Approximation of Fixed Points, Springer Berlin Heidelberg, Berlin, Heidelberg, pp. 204–227.
- Kassam, A., Trefethen, L., 2005. Fourth-order time-stepping for stiff pdes. SIAM Journal on Scientific Computing 26, 1214–1233. doi:10.1137/S106482750410633.
- Kingma, D.P., Ba, J., 2015. Adam: A method for stochastic optimization, in: 3rd International Conference on Learning Representations, ICLR 2015, San Diego, CA, USA, May 7-9, 2015, Conference Track Proceedings.
- Koopman, B.O., 1931. Hamiltonian systems and transformation in hilbert space. Proceedings of the National Academy of Sciences of the United States of America 17, 315.
- Krueger, D., Maharaj, T., Kramár, J., Pezeshki, M., Ballas, N., Ke, N.R., Goyal, A., Bengio, Y., Courville, A.C., Pal, C.J., 2017. Zoneout: Regularizing rnns by randomly preserving hidden activations, in: 5th International Conference on Learning Representations, ICLR 2017, Toulon, France, April 24-26, 2017, Conference Track Proceedings. URL: <https://openreview.net/forum?id=rJqBEPcxe>.
- Kuramoto, Y., 1978. Diffusion-Induced Chaos in Reaction Systems. Progress of Theoretical Physics Supplement 64, 346–367. URL: <https://dx.doi.org/10.1143/PTPS.64.346>, doi:10.1143/PTPS.64.346, arXiv:<http://oup.prod.sis.lan/ptps/article-pdf/doi/10.1143/PTPS.64.346/5293041/64-346.pdf>.
- Kurth, T., Treichler, S., Romero, J., Mudigonda, M., Luehr, N., Phillips, E., Mahesh, A., Matheson, M., Deslippe, J., Fatica, M., et al., 2018. Exascale deep learning for climate analytics, in: SC18: International Conference for High Performance Computing, Networking, Storage and Analysis, IEEE, pp. 649–660.
- Lapedes, A., Farber, R., 1987. Nonlinear signal processing using neural networks: Prediction and system modelling. Technical Report.
- Laptev, N., Yosinski, J., Li, L.E., Smyl, S., 2017. Time-series extreme event forecasting with neural networks at uber.
- Larger, L., Baylón-Fuentes, A., Martinenghi, R., Udaltssov, V.S., Chembo, Y.K., Jacquot, M., 2017. High-speed photonic reservoir computing using a time-delay-based architecture: Million words per second classification. Phys. Rev. X 7, 011015. URL: <https://link.aps.org/doi/10.1103/PhysRevX.7.011015>, doi:10.1103/PhysRevX.7.011015.
- Larger, L., Soriano, M.C., Brunner, D., Appeltant, L., Gutierrez, J.M., Pesquera, L., Mirasso, C.R., Fischer, I., 2012. Photonic information processing beyond turing: an optoelectronic implementation of reservoir computing. Opt. Express 20, 3241–3249. URL: <http://www.opticsexpress.org/abstract.cfm?URI=oe-20-3-3241>, doi:10.1364/OE.20.003241.
- Linnainmaa, S., 1976. Taylor expansion of the accumulated rounding error. BIT Numerical Mathematics 16, 146–160. doi:10.1007/BF01931367.
- Lorenz, E., 1995. Predictability: a problem partly solved, in: Seminar on Predictability, 4-8 September 1995, ECMWF, ECMWF, Shinfield Park, Reading, pp. 1–18. URL: <https://www.ecmwf.int/node/10829>.
- Lu, Z., Hunt, B.R., Ott, E., 2018. Attractor reconstruction by machine learning. Chaos: An Interdisciplinary Journal of Nonlinear Science 28, 061104.
- Lukoševičius, M., 2012. A practical guide to applying echo state networks, in: Neural Networks: Tricks of the Trade.
- Lukoševičius, M., Jaeger, H., 2009. Reservoir computing approaches to recurrent neural network training. Computer Science Review 3, 127 – 149. doi:<https://doi.org/10.1016/j.cosrev.2009.03.005>.
- Lusch, B., Kutz, J.N., Brunton, S.L., 2018. Deep learning for universal linear embeddings of nonlinear dynamics. Nature communications 9, 4950.
- Maass, W., Natschläger, T., Markram, H., 2002. Real-time computing without stable states: A new framework for neural computation based on perturbations. Neural Comput. 14, 2531–2560. URL: <https://dx.doi.org/10.1162/089976602760407955>, doi:10.1162/089976602760407955.
- Makarenko, A.V., 2018. Deep convolutional neural networks for chaos identification in signal processing, in: 2018 26th European Signal Processing Conference (EUSIPCO), IEEE, pp. 1467–1471.
- Manneville, P., 1984. Macroscopic modelling of turbulent flows, proceedings of a workshop held at inria, sophia-antipolis, france, 1984. Lecture Notes in Physics 230, 319–326.
- Maus, A., Sprott, J.C., 2013. Evaluating lyapunov exponent spectra with neural networks. Chaos, Solitons and Fractals 51, 13–21. doi:10.1016/j.chaos.2013.03.001.
- Mnih, V., Kavukcuoglu, K., Silver, D., Rusu, A.A., Veness, J., Bellemare, M.G., Graves, A., Riedmiller, M., Fidjeland, A.K., Ostrovski, G., et al., 2015. Human-level control through deep reinforcement learning. Nature 518, 529.
- Neofotistos, G., Mattheakis, M., Barmparis, G.D., Hizanidis, J., Tsironis, G.P., Kaxiras, E., 2019. Machine learning with observers predicts complex spatiotemporal behavior. Frontiers in Physics 7, 24. URL: <https://www.frontiersin.org/article/10.3389/fphy.2019.00024>, doi:10.3389/fphy.2019.00024.
- Ott, E., 2002. Chaos in Dynamical Systems. 2 ed., Cambridge University Press. doi:10.1017/CBO9780511803260.
- Parlitz, U., Merkowirth, C., 2000. Prediction of spatiotemporal time series based on reconstructed local states. Phys. Rev. Lett. 84, 1890–1893. URL: <https://link.aps.org/doi/10.1103/PhysRevLett.84.1890>, doi:10.1103/PhysRevLett.84.1890.
- Pascanu, R., Mikolov, T., Bengio, Y., 2013. On the difficulty of training recurrent neural networks, in: Proceedings of the 30th International Conference on International Conference on Machine Learning - Volume 28, JMLR.org, pp. III–1310–III–1318. URL: <https://dl.acm.org/citation.cfm?id=3042817.3043083>.
- Paszke, A., Gross, S., Chintala, S., Chanan, G., Yang, E., DeVito, Z., Lin, Z., Desmaison, A., Antiga, L., Lerer, A., 2017. Automatic differentiation in pytorch.
- Pathak, J., Hunt, B., Girvan, M., Lu, Z., Ott, E., 2018a. Model-free prediction of large spatiotemporally chaotic systems from data: A reservoir computing approach. Phys. Rev. Lett. 120, 024102. URL: <https://link.aps.org/doi/10.1103/PhysRevLett.120.024102>, doi:10.1103/PhysRevLett.120.024102.
- Pathak, J., Lu, Z., Hunt, B.R., Girvan, M., Ott, E., 2017. Using machine learning to replicate chaotic attractors and calculate lyapunov exponents from data. Chaos: An Interdisciplinary Journal of Nonlinear Science 27, 121102. doi:10.1063/1.5010300.
- Pathak, J., Wikner, A., Fussell, R., Chandra, S., Hunt, B.R., Girvan, M., Ott, E., 2018b. Hybrid forecasting of chaotic processes: Using machine learning in conjunction with a knowledge-based model. Chaos: An Interdisciplinary Journal of Nonlinear Science 28,

041101. URL: <https://app.dimensions.ai/details/publication/pub.1103541484> and <http://arxiv.org/pdf/1803.04779.pdf>, doi:10.1063/1.5028373, exported from <https://app.dimensions.ai> on 2019/02/13.
- Rumelhart, D.E., Hinton, G.E., Williams, R.J., 1986. Learning Representations by Back-propagating Errors. *Nature* 323, 533–536. URL: <http://www.nature.com/articles/323533a0>, doi:10.1038/323533a0.
- Sainath, T.N., Vinyals, O., Senior, A., Sak, H., 2015. Convolutional, long short-term memory, fully connected deep neural networks, in: 2015 IEEE International Conference on Acoustics, Speech and Signal Processing (ICASSP), IEEE. pp. 4580–4584.
- Sano, M., Sawada, Y., 1985. Measurement of lyapunov spectrum from a chaotic time series. *Physical review letters* 55, 1082–1085. doi:10.1103/PhysRevLett.55.1082.
- Sauer, T., Yorke, J.A., Casdagli, M., 1991. Embedology. *Journal of Statistical Physics* 65, 579–616. URL: <https://doi.org/10.1007/BF01053745>, doi:10.1007/BF01053745.
- Schäfer, A.M., Zimmermann, H.G., 2006. Recurrent neural networks are universal approximators, in: Proceedings of the 16th International Conference on Artificial Neural Networks - Volume Part I, Springer-Verlag, Berlin, Heidelberg. pp. 632–640. URL: http://dx.doi.org/10.1007/11840817_66, doi:10.1007/11840817_66.
- Schriftwieser, J., Antonoglou, I., Hubert, T., Simonyan, K., Sifre, L., Schmitt, S., Guez, A., Lockhart, E., Hassabis, D., Graepel, T., et al., 2019. Mastering atari, go, chess and shogi by planning with a learned model. *arXiv preprint arXiv:1911.08265*.
- Shi, X., Chen, Z., Wang, H., Yeung, D.Y., Wong, W.K., chun Woo, W., 2015. Convolutional lstm network: A machine learning approach for precipitation nowcasting, in: NIPS.
- Siegelmann, H., Sontag, E., 1995. On the computational power of neural nets. *J. Comput. Syst. Sci.* 50, 132–150. URL: <http://dx.doi.org/10.1006/jcss.1995.1013>, doi:10.1006/jcss.1995.1013.
- Silver, D., Huang, A., Maddison, C.J., Guez, A., Sifre, L., Van Den Driessche, G., Schriftwieser, J., Antonoglou, I., Panneershelvam, V., Lanctot, M., et al., 2016. Mastering the game of go with deep neural networks and tree search. *nature* 529, 484.
- Sivashinsky, G.I., 1977. Nonlinear analysis of hydrodynamic instability in laminar flames — I. Derivation of basic equations. *Acta Astronautica* 4, 1177–1206. doi:10.1016/0094-5765(77)90096-0.
- Sünderhauf, N., Brock, O., Scheirer, W., Hadsell, R., Fox, D., Leitner, J., Upcroft, B., Abbeel, P., Burgard, W., Milford, M., et al., 2018. The limits and potentials of deep learning for robotics. *The International Journal of Robotics Research* 37, 405–420.
- Takens, F., 1981. Detecting strange attractors in turbulence, in: Rand, D., Young, L.S. (Eds.), *Dynamical Systems and Turbulence*, Warwick 1980, Springer Berlin Heidelberg, Berlin, Heidelberg. pp. 366–381.
- Tikhonov, A.N., Arsenin, V.Y., 1977. Solutions of Ill-posed problems. W.H. Winston.
- Vlachas, P.R., Byeon, W., Wan, Z.Y., Sapsis, T.P., Koumoutsakos, P., 2018. Data-driven forecasting of high-dimensional chaotic systems with long short-term memory networks. *Proceedings of the Royal Society A: Mathematical, Physical and Engineering Sciences* 474, 20170844. URL: <https://royalsocietypublishing.org/doi/abs/10.1098/rspa.2017.0844>, doi:10.1098/rspa.2017.0844, arXiv:<https://royalsocietypublishing.org/doi/10.1098/rspa.2017.0844>.
- Walker, D.W., Dongarra, J.J., 1996. Mpı: a standard message passing interface. *Supercomputer* 12, 56–68.
- Wan, Z.Y., Sapsis, T.P., 2018. Machine learning the kinematics of spherical particles in fluid flows. *Journal of Fluid Mechanics* 857. doi:10.1017/jfm.2018.797.
- Wan, Z.Y., Vlachas, P., Koumoutsakos, P., Sapsis, T., 2018. Data-assisted reduced-order modeling of extreme events in complex dynamical systems. *PLOS ONE* 13, 1–22. URL: <https://doi.org/10.1371/journal.pone.0197704>, doi:10.1371/journal.pone.0197704.
- Werbos, P.J., 1988. Generalization of backpropagation with application to a recurrent gas market model. *Neural Networks* 1, 339–356. URL: <http://www.sciencedirect.com/science/article/pii/089360808890007X>, doi:[https://doi.org/10.1016/0893-6080\(88\)90007-X](https://doi.org/10.1016/0893-6080(88)90007-X).
- Werbos, P.J., 1990. Backpropagation through time: what it does and how to do it. *Proceedings of the IEEE* 78, 1550–1560. doi:10.1109/5.58337.
- Weyn, J.A., Durran, D.R., Caruana, R., 2019. Can machines learn to predict weather? using deep learning to predict gridded 500-hpa geopotential height from historical weather data. *Journal of Advances in Modeling Earth Systems* 11, 2680–2693.
- Wolf, A., B. Swift, J., Swinney, H., A. Vastano, J., 1985. Determining lyapunov exponents from a time series. *Physica D: Nonlinear Phenomena* 16, 285–317. doi:10.1016/0167-2789(85)90011-9.
- Yan, X., Su, X.G., 2009. Linear Regression Analysis. WORLD SCIENTIFIC. URL: <https://www.worldscientific.com/doi/abs/10.1142/6986>, doi:10.1142/6986, arXiv:<https://www.worldscientific.com/doi/pdf/10.1142/6986>.

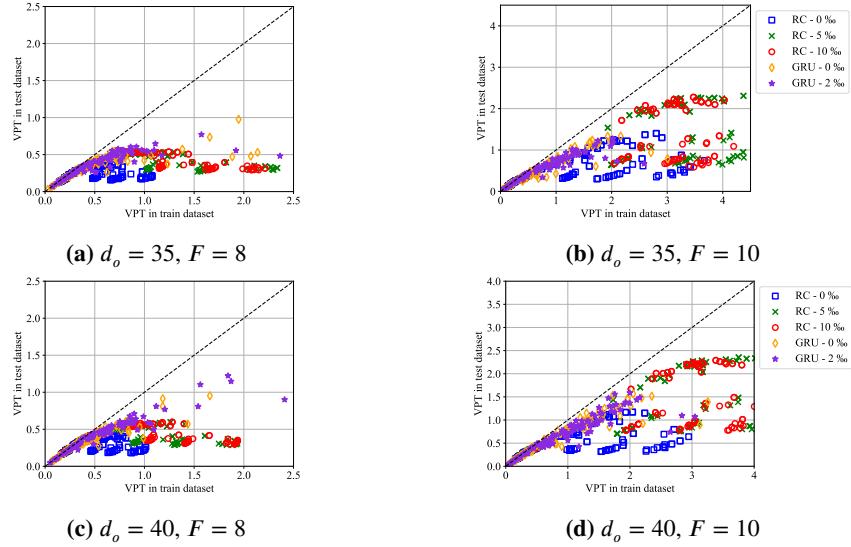


Figure 17: VPT in the testing data plotted against VPT in the training data for RC and GRU models trained with added noise of different levels in the data. Noise only slightly varies the forecasting efficiency in GRU networks. In contrast, the effectiveness of RC in forecasting the full-order system is increased as depicted in plots (b) and (d).

A. Memory Efficient Implementation of RC Training

In order to alleviate the RAM requirement for the computation of the RC weights we resort to a batched approach. Assuming the hidden reservoir size is given by $\mathbf{h} \in \mathbb{R}^{d_h}$, by teacher forcing the RC network with true data from the system for d_N time-steps and stacking the evolution of the hidden state in a single matrix we end up with matrix $\mathbf{H} \in \mathbb{R}^{d_N \times d_h}$. Moreover, by stacking the target values, which are the input data shifted by one time-step, we end up in the target matrix $\mathbf{Y} \in \mathbb{R}^{d_N \times d_o}$, where d_o is the dimension of the observable we are predicting. In order to identify the output weights $W_{out} \in \mathbb{R}^{d_o \times d_h}$, we need to solve the linear system of $d_N \cdot d_o$ equations

$$\mathbf{H}W_{out}^T = \mathbf{Y}. \quad (14)$$

A classical way to solve this system of equations is based on the Moore-Penrose inverse (pseudo-inverse) computed using

$$W_{out} = \underbrace{\mathbf{Y}^T \mathbf{H}}_{\bar{\mathbf{Y}}} \left(\underbrace{\mathbf{H}^T \mathbf{H}}_{\bar{\mathbf{H}}} + \eta \mathbf{I} \right)^{-1} \quad (15)$$

where η is the Tikhonov regularization parameter and \mathbf{I} the unit matrix. In our case d_N is of the order of 10^5 and $d_N \gg d_h$. To reduce the memory requirements of the training method, we compute the matrices $\bar{\mathbf{H}} = \mathbf{H}^T \mathbf{H} \in \mathbb{R}^{d_h \times d_h}$ and $\bar{\mathbf{Y}} = \mathbf{Y}^T \mathbf{H} \in \mathbb{R}^{d_o \times d_h}$ in a time-batched schedule.

Specifically, we initialize $\bar{\mathbf{Y}} = \mathbf{0}$ and $\bar{\mathbf{H}} = \mathbf{0}$. Then every d_n time-steps with $d_n \ll d_N$, we compute the batch matrix $\bar{\mathbf{H}}_b = \mathbf{H}_b^T \mathbf{H}_b \in \mathbb{R}^{d_h \times d_h}$, where $\mathbf{H}_b \in \mathbb{R}^{d_n \times d_h}$ is formed by the stacking the hidden state only for the last d_n time-steps. In the same way, we compute $\bar{\mathbf{Y}}_b = \mathbf{Y}_b^T \mathbf{H}_b \in \mathbb{R}^{d_o \times d_h}$, where $\mathbf{Y}_b \in \mathbb{R}^{d_n \times d_o}$ is formed by the stacking of the target data for the last d_n time-steps. After every batch computation we update our beliefs with $\bar{\mathbf{H}} \leftarrow \bar{\mathbf{H}} + \bar{\mathbf{H}}_b$ and $\bar{\mathbf{Y}} \leftarrow \bar{\mathbf{Y}} + \bar{\mathbf{Y}}_b$.

In addition, we also experimented with two alternative solvers for the linear system Equation (14) in the Lorenz-96. We tried a dedicated regularized least-squares routine utilizing an iterative procedure (`scipy.sparse.linalg.lsqr`) and a method based on stochastic gradient descent. We considered the solver as an additional hyperparameter of the RC models. After testing the solvers in Lorenz-96 systems, we found out that the method of pseudo-inverse provides the most accurate results. For this reason, and to spare computational resources, we used this method for the Kuramoto-Sivashinsky system.

B. Regularizing Training with Noise

In our study, we investigate the effect of noise to the training data. In Figure 17, we plot the Valid Prediction Time (VPT) in the testing data with respect to the VPT that each model achieves in the training data. We find out that RC models trained with additional noise of 5 – 10 % not only achieve better generalization, but their forecasting efficiency improves in both training and testing dataset. Moreover, the effect of divergent predictions by iterative forecasts is alleviated significantly. In contrast, adding noise does not seem to have an important impact on the performance of GRU models.

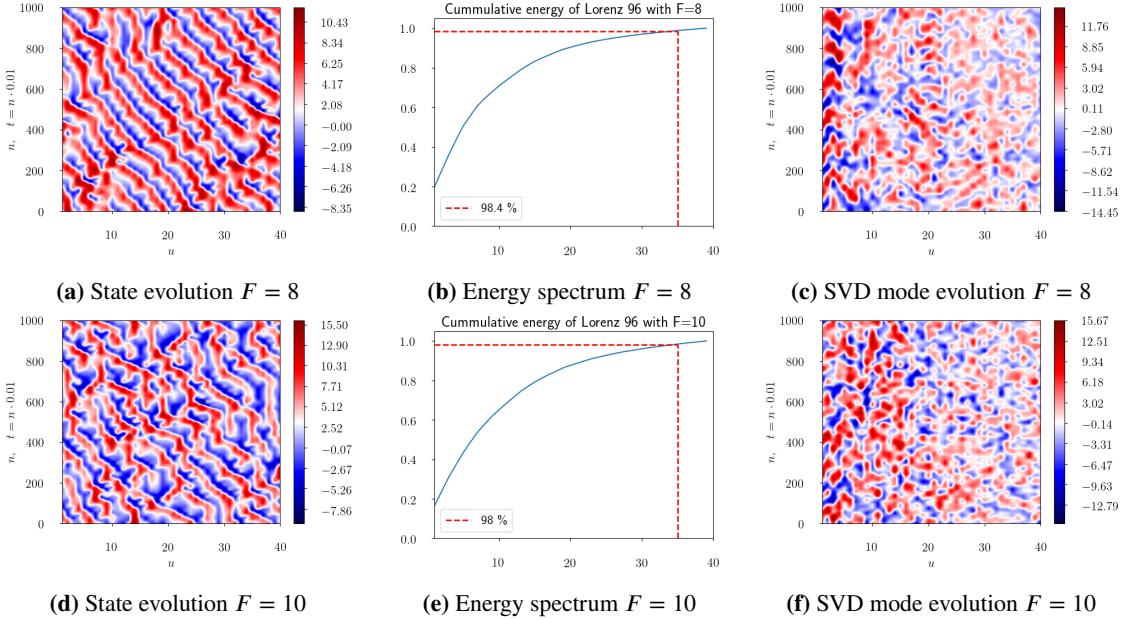


Figure 18: Energy spectrum of Lorenz-96

C. Dimensionality Reduction with Singular Value Decomposition

Singular Value Decomposition (SVD) can be utilized to perform dimensionality reduction in a dataset by identifying the modes that capture the highest variance in the data and then performing a projection on these modes. Assuming that a data matrix is given by stacking the time-evolution of a state $\mathbf{u} \in \mathbb{D}$ as $U = [\mathbf{u}_1, \mathbf{u}_2, \dots, \mathbf{u}_N]$, where the index N is the number of data samples. By subtracting the temporal mean $\bar{\mathbf{u}}$ and stacking the data, we end up with the data matrix $\mathbf{U} \in \mathbb{R}^{T \times D}$. Performing SVD on \mathbf{U} leads to

$$\mathbf{U} = \mathbf{M}\Sigma\mathbf{V}^T, \quad \mathbf{M} \in \mathbb{R}^{N \times N}, \quad \Sigma \in \mathbb{R}^{N \times D}, \quad \mathbf{V} \in \mathbb{R}^{D \times D}, \quad (16)$$

with Σ diagonal, with descending diagonal elements. The columns of matrix V are considered the modes of the SVD, while the square D singular values of Σ correspond to the data variance explained by these modes. This variance is also referred to as energy. In order to calculate the percentage of the total energy the square of the singular value of each mode has to be divided by the sum of squares of the singular values of all modes. In order to reduce the dimensionality of the dataset, we first have to decide on the reduced order dimension $r_{dim} < D$. Then we identify the eigenvectors corresponding to the most high-energetic eigenmodes. These are given by the first columns \mathbf{V}_r of \mathbf{V} , i.e., $\mathbf{V} = [\mathbf{V}_r, \mathbf{V}_{-r}]$. We discard the low-energetic modes \mathbf{V}_{-r} . The dimension of the truncated eigenvector matrix is $\mathbf{V}_r \in \mathbb{R}^{D \times r_{dim}}$. In order to reduce the dimensionality of the dataset, each vector $\mathbf{u} \in \mathbb{D}$ is projected to $\mathbf{u}_r \in r_{dim}$ by

$$\mathbf{c} = \mathbf{V}_r^T \mathbf{u}, \quad \mathbf{c} \in \mathbb{R}^{r_{dim}}. \quad (17)$$

In the Lorenz-96 system, we construct a reduced order observable with $d_o = 35$ modes of the system. The cumulative energy distribution along with a contour plot of the state and the mode evolution is illustrated in Figure 18.

D. Calculation of Lyapunov Spectrum

The true Lyapunov exponents of the KS equation are computed as in (Pathak et al., 2018a) by solving the KS equations in the Fourier space with a fourth order time-stepping method called ETDRK4 (Kassam and Trefethen, 2005) and utilizing a QR decomposition approach. The trained RNN model with GRU cell is used as a surrogate to compute the full Lyapunov spectrum of the Kuramoto-Sivashinsky system. Recall that the RNN dynamics are given by

$$\begin{aligned} \mathbf{h}_t &= f_h^h(\mathbf{o}_t, \mathbf{h}_{t-1}) \\ \mathbf{o}_{t+1} &= f_h^o(\mathbf{h}_t), \end{aligned} \quad (18)$$

where f_h^h is the hidden-to-hidden and f_h^o is the hidden-to-output mapping, $\mathbf{o} \in \mathbb{R}^{d_o}$ is an observable of the state, and $\mathbf{h}_t \in \mathbb{R}^{d_h}$ is the hidden state of the RNN. All models utilized in this work share this common architecture. They only differ in the forms of f_h^o and f_h^h . More importantly, the output mapping is linear, i.e.,

$$\mathbf{o}_{t+1} = f_h^o(\mathbf{h}_t) = \mathbf{W}_o \mathbf{h}_t. \quad (19)$$

The LEs are calculated based on the Jacobian $\mathbf{J} = \frac{dh_t}{dh_{t-1}}$ of the hidden state dynamics along the trajectory. In the following we compute the Jacobian using Equation (18). By writing down the equations for two consecutive time-steps, we get

$$\text{Timestep } t-1 : h_{t-1} = f_h^h(o_{t-1}, h_{t-2}) \quad (20)$$

$$o_t = f_h^o(h_{t-1}) = \mathbf{W}_o h_{t-1} \quad (21)$$

$$\text{Timestep } t : h_t = f_h^h(o_t, h_{t-1}). \quad (22)$$

The partial Jacobians needed to compute the total Jacobian are:

$$\frac{\partial f_h^h}{\partial o} = J_o^{hh} \in \mathbb{R}^{d_h \times d_o} \quad (23)$$

$$\frac{\partial f_h^h}{\partial h} = J_h^{hh} \in \mathbb{R}^{d_h \times d_h} \quad (24)$$

$$\frac{\partial f_h^o}{\partial h} = J_h^{oh} \in \mathbb{R}^{d_o \times d_h}. \quad (25)$$

In total we can write:

$$\frac{dh_t}{dh_{t-1}} = \frac{df_h^h(o_t, h_{t-1})}{dh_{t-1}} = \frac{\partial f_h^h(o_t, h_{t-1})}{\partial o_t} \frac{\partial o_t}{\partial h_{t-1}} + \frac{\partial f_h^h(o_t, h_{t-1})}{\partial h_{t-1}} \Rightarrow \quad (26)$$

$$\frac{dh_t}{dh_{t-1}} = \frac{\partial f_h^h(o_t, h_{t-1})}{\partial o_t} \frac{\partial f_h^o(h_{t-1})}{\partial h_{t-1}} + \frac{\partial f_h^h(o_t, h_{t-1})}{\partial h_{t-1}} \Rightarrow \quad (27)$$

$$\frac{dh_t}{dh_{t-1}} = \underbrace{J_o^{hh} \Big|_{(o_t, h_{t-1})}}_{\text{evaluated at } t} \cdot \underbrace{J_h^{oh} \Big|_{h_{t-1}}}_{\text{evaluated at } t-1} + \underbrace{J_h^{hh} \Big|_{(o_t, h_{t-1})}}_{\text{evaluated at } t} \quad (28)$$

A product of this Jacobian along the orbit δ is developed and iteratively orthonormalized every T_n steps using the Gram-Schmidt method to avoid numerical divergence and keep the columns of the matrix R independent. We check the convergence criterion by tracking the estimated LE values every T_c time-steps. The input provided to the algorithm is a short time-series of length T_w , to initialize the RNN and warm-up the hidden state $\tilde{o}_{1:T_w+1}$ (where the tilde denotes experimental or simulation data), the length of this warm-up time-series T_w , the number of the LE to calculate N , the maximum time to unroll the RNN T , a normalization time T_n and an additional threshold ϵ used as an additional termination criterion. The function ColumnSum(\cdot) computes the sum of each column of a matrix, i.e., $\text{sum}(\cdot, \text{axis} = 1)$. This method can be applied directly to RNNs with one hidden state like RC or GRUs. An adaptation to the LSTM is left for future research. The pseudocode of the algorithm to calculate the Lyapunov exponents of the RNN is given in Algorithm 1.

Algorithm 1 Algorithm to calculate Lyapunov Exponents of a trained surrogate RNN model

```

procedure LE_RNN( $\tilde{o}_{1:T_w+1}, T_w, N, T, T_n, \epsilon$ )
    Initialize  $\mathbf{h}_0 \leftarrow 0$ .
    for  $t = 1 : T_w$  do                                 $\triangleright$  Warming-up the hidden state of the RNN based on true data
         $|\quad \mathbf{h}_t \leftarrow f_h^h(\tilde{o}_t, \mathbf{h}_{t-1})$ 
    end for
     $\mathbf{h}_0 \leftarrow \mathbf{h}_{T_w}$ 
     $\mathbf{o}_1 \leftarrow \tilde{o}_{T_w+1}$ 
    Pick a random orthonormal matrix  $\delta \in \mathbb{R}^{d_h \times N_{LE}}$ .            $\triangleright$  Initializing  $N_{LE}$  deviation vectors
     $\tilde{T} \leftarrow T/T_n$ 
    Initialize  $\tilde{R} \leftarrow \mathbf{0} \in \mathbb{R}^{N \times \tilde{T}}$ .
     $l_{prev}, l \leftarrow \mathbf{0} \in \mathbb{R}^N$                                       $\triangleright$  Initializing the  $N$  LE to zero.
     $J_0 \leftarrow \nabla_{\mathbf{h}} f_h^o(\mathbf{h}_0)$ .                                $\triangleright$  Evolve the RNN dynamics
    for  $t = 1 : T$  do
         $|\quad \mathbf{h}_t \leftarrow f_h^h(\mathbf{o}_t, \mathbf{h}_{t-1})$ 
         $|\quad \mathbf{o}_{t+1} \leftarrow f_h^o(\mathbf{h}_t)$ 
         $|\quad J_1 \leftarrow \nabla_{\mathbf{h}} f_h^h(\mathbf{o}_{t+1}, \mathbf{h}_t)$ .            $\triangleright$  Calculating the partial Jacobians
         $|\quad J_2 \leftarrow \nabla_{\mathbf{o}} f_h^h(\mathbf{o}_{t+1}, \mathbf{h}_t)$ .
         $|\quad J \leftarrow J_1 + J_2 \cdot J_0$ .            $\triangleright$  Calculating the total Jacobian
         $|\quad \delta \leftarrow J \cdot \delta$                           $\triangleright$  Evolving the deviation vectors  $\delta$ 
        if mod( $t, T_{norm}$ ) = 0 then
             $|\quad Q, R \leftarrow QR(\delta)$ 
             $|\quad \delta \leftarrow Q[:, :N]$ 
             $|\quad \tilde{R}[:, t/T_{norm}] \leftarrow \log(\text{diag}(R[:, N, :N]))$             $\triangleright$  Re-orthonormalizing with QR-decomposition
            if mod( $t, T_c$ ) = 0 then
                 $|\quad l \leftarrow \text{Real}(\text{ColumnSum}(\tilde{R})) / (t * \delta t)$ 
                 $|\quad l \leftarrow \text{sort}(l)$ 
                 $|\quad d \leftarrow \|l - l_{prev}\|_2$ 
                if  $d < \epsilon$  then
                     $|\quad \quad \quad \text{break}$ 
                end if
            end if
        end if
         $|\quad J_0 \leftarrow \nabla_{\mathbf{h}} f_h^o(\mathbf{h}_t)$ 
    end for
    return  $l$ 
end procedure

```

Table 2Hyperparameters of **RC** for Lorenz-96

Hyperparameter	Explanation	Values
D_r	reservoir size	{6000, 9000, 12000, 18000}
N	training data samples	10^5
d	Solver degree of $W_{h,h}$	{3, 8}
ρ	radius of $W_{h,h}$	{0.4, 0.8, 0.9, 0.99}
ω	input scaling	{0.1, 0.5, 1.0, 1.5, 2.0}
η	regularization	{ 10^{-3} , 10^{-4} , 10^{-5} , 10^{-6} }
d_o	observed state dimension	{35, 40}
n_w	warm-up steps (testing)	2000
κ_n	noise level in data	{0, 0.5%, 1%}

Table 3Hyperparameters of **GRU/LSTM** for Lorenz-96

Hyperparameter	Explanation	Values
d_h	hidden state size	{1, 2, 3} layers of {500, 1000, 1500}
N	training data samples	10^5
B	batch-size	32
κ_1	BPTT forward time steps	{1, 8}
κ_2	BPTT truncated backprop. length	{8, 16}
κ_3	BPTT skip gradient parameter	= $\kappa_2 + \kappa_1 - 1$
η	initial learning rate	10^{-3}
p	zoneout probability	{0.99, 0.995 1.0}
d_o	observed state dimension	{35, 40}
n_w	warm-up steps (testing)	2000
κ_n	noise level in data	{0, 0.2%}

Table 4Hyperparameters of **Unitary Evolution** networks for Lorenz-96

Hyperparameter	Explanation	Values
d_h	hidden state size	{1, 2, 3} layers of {500, 1000, 1500}
N	training data samples	10^5
B	batch-size	32
κ_1	BPTT forward time steps	{1, 8}
κ_2	BPTT truncated backprop. length	{8, 16}
κ_3	BPTT skip gradient parameter	= $\kappa_2 + \kappa_1 - 1$
η	initial learning rate	10^{-3}
p	zoneout probability	1.0
d_o	observed state dimension	{35, 40}
n_w	warm-up steps (testing)	2000
κ_n	noise level in data	{0, 0.2%}

E. Model Hyperparameters

For the Lorenz-96 system space with $d_o \in \{35, 40\}$ (in the PCA mode), we used the hyperparameters reported on Table 2 for RC and Table 3 for GRU/LSTM models. For the parallel architectures in the state space of Lorenz-96 the hyperparameters are reported on Table 5 and Table 6 for the parallel RC and GRU/LSTM models respectively. For the parallel architectures in the state space of the Kuramoto-Sivashinsky architecture the hyperparameters are reported on Table 7 and Table 8 for the parallel RC and GRU/LSTM models respectively. We note here that in all RNN methods, the optimizer used to update the network can also be optimized. To alleviate the computational burden we stick to Adam.

Table 5

Hyperparameters of Parallel RC for Lorenz-96

Hyperparameter	Explanation	Values
D_r	reservoir size	{1000, 3000, 6000, 12000}
N_g	number of groups	20
G	group size	2
I	interaction length	4
N	training data samples	10^5
d	Solver	Pseudoinverse/LSQR/Gradient descent
ρ	degree of $W_{h,h}$	10
ω	radius of $W_{h,h}$	0.6
η	input scaling	0.5
d_o	regularization	10^{-6}
n_w	observed state dimension	40
	warm-up steps (testing)	2000

Table 6

Hyperparameters of Parallel GRU/LSTM for Lorenz-96

Hyperparameter	Explanation	Values
d_h	hidden state size	{100, 250, 500}
N_g	number of groups	20
G	group size	2
I	interaction length	4
N	training data samples	10^5
B	batch-size	32
κ_1	BPTT forward time steps	4
κ_2	BPTT truncated backprop. length	4
κ_3	BPTT skip gradient parameter	4
η	initial learning rate	10^{-3}
p	zoneout probability	{0.998, 1.0}
d_o	observed state dimension	40
n_w	warm-up steps (testing)	2000

Table 7

Hyperparameters of Parallel RC for Kuramoto-Sivashinsky

Hyperparameter	Explanation	Values
D_r	reservoir size	{500, 1000, 3000, 6000, 12000}
N_g	number of groups	64
G	group size	8
I	interaction length	8
N	training data samples	10^5
	Solver	Pseudoinverse
d	degree of $W_{h,h}$	10
ρ	radius of $W_{h,h}$	0.6
ω	input scaling	1.0
η	regularization	10^{-5}
d_o	observed state dimension	512
n_w	warm-up steps (testing)	2000

Table 8Hyperparameters of **Parallel GRU/LSTM** for Kuramoto-Sivashinsky

Hyperparameter	Explanation	Values
d_h	hidden state size	{80, 100, 120}
N_g	number of groups	64
G	group size	8
I	interaction length	8
N	training data samples	10^5
B	batch-size	32
κ_1	BPTT forward time steps	4
κ_2	BPTT truncated backprop. length	4
κ_3	BPTT skip gradient parameter	4
η	initial learning rate	10^{-3}
p	zoneout probability	{0.998, 1.0}
d_o	observed state dimension	512
n_w	warm-up steps (testing)	2000

Table 9Hyperparameters of **Parallel Unitary Evolution** networks for Kuramoto-Sivashinsky

Hyperparameter	Explanation	Values
d_h	hidden state size	{100, 200, 400}
N_g	number of groups	64
G	group size	8
I	interaction length	8
N	training data samples	10^5
B	batch-size	32
κ_1	BPTT forward time steps	4
κ_2	BPTT truncated backprop. length	4
κ_3	BPTT skip gradient parameter	4
η	initial learning rate	10^{-2}
p	zoneout probability	1.0
d_o	observed state dimension	512
n_w	warm-up steps (testing)	2000

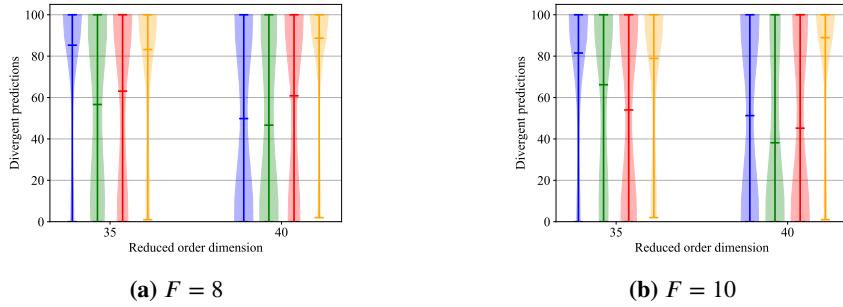


Figure 19: Violin plots with kernel density estimates of the number of divergent predictions over the 100 initial conditions from the test data, over all hyperparameter sets for RC, GRU, LSTM, and Unitary networks, for reduced order state $d_o = 35$ and full order state $d_o = 40$ in two forcing regimes (a) $F = 8$ and (b) $F = 10$ in the Lorenz-96 system. Most hyperparameter sets of Unitary networks, lead to models that diverge in iterative forecasting in both reduced order and full order scenario for both $F \in \{8, 10\}$. Although the divergence effect is a non-issue in RC in the full state scenario $d_o = 40$, indicated by the wide part in the lower end of the density plot, the effect is more prominent in the reduced order scenario compared to GRU and LSTM. Identification of hyperparameters for LSTM and GRU networks that show stable iterative forecasting behavior in the reduced order space is significantly easier compared to RC and Unitary networks, as indicated by the wide/thin lines in the lower part of the density plots of the first/latter.

RC ; GRU ; LSTM ; Unit ;

F. Additional Results - Lorenz-96 - Divergence of Unitary and RC RNNs

In this section, we try to quantify the divergence effect due to the accumulation of the forecasting error in the iterative prediction. In Figure 19 we present violin plots with fitted kernel density estimates for the number of divergent predictions of each hyperparameter set, computed based on all tested hyperparameter sets for forcing regimes $F \in \{8, 10\}$ and observable dimensions $d_o \in \{35, 40\}$. The annotated lines denote the minimum, mean and maximum number of divergent predictions over the 100 initial conditions of all hyperparameter sets. In the fully observable systems $d_o = 40$, in both forcing regimes $F \in \{8, 10\}$, there are many models (hyperparameter sets) with zero divergent predictions for RC, GRU and LSTM, as illustrated by the wide lower part of the violin plot. In contrast, most hyperparameter sets in Unitary networks lead to models whose iterative predictions diverge from the attractor, as illustrated by the wide upper part in the violin plot. In the reduced order scenario, this divergence effect seems to be more prominent in RC and Unitary networks, as indicated by the very thin lower part of their violin plots, for both forcing regimes. In contrast, many hyperparameter sets of GRU and LSTM models lead to stable iterative prediction. This indicates that hyperparameter tuning in RC and Unitary networks when the system state is not fully observed, is cumbersome compared to LSTM and GRU networks. One example of this divergence effect in an initial condition from the **test** dataset is illustrated in Figure 20. The RC and the Unitary networks diverge in the reduced order state predictions after approximately two Lyapunov times.

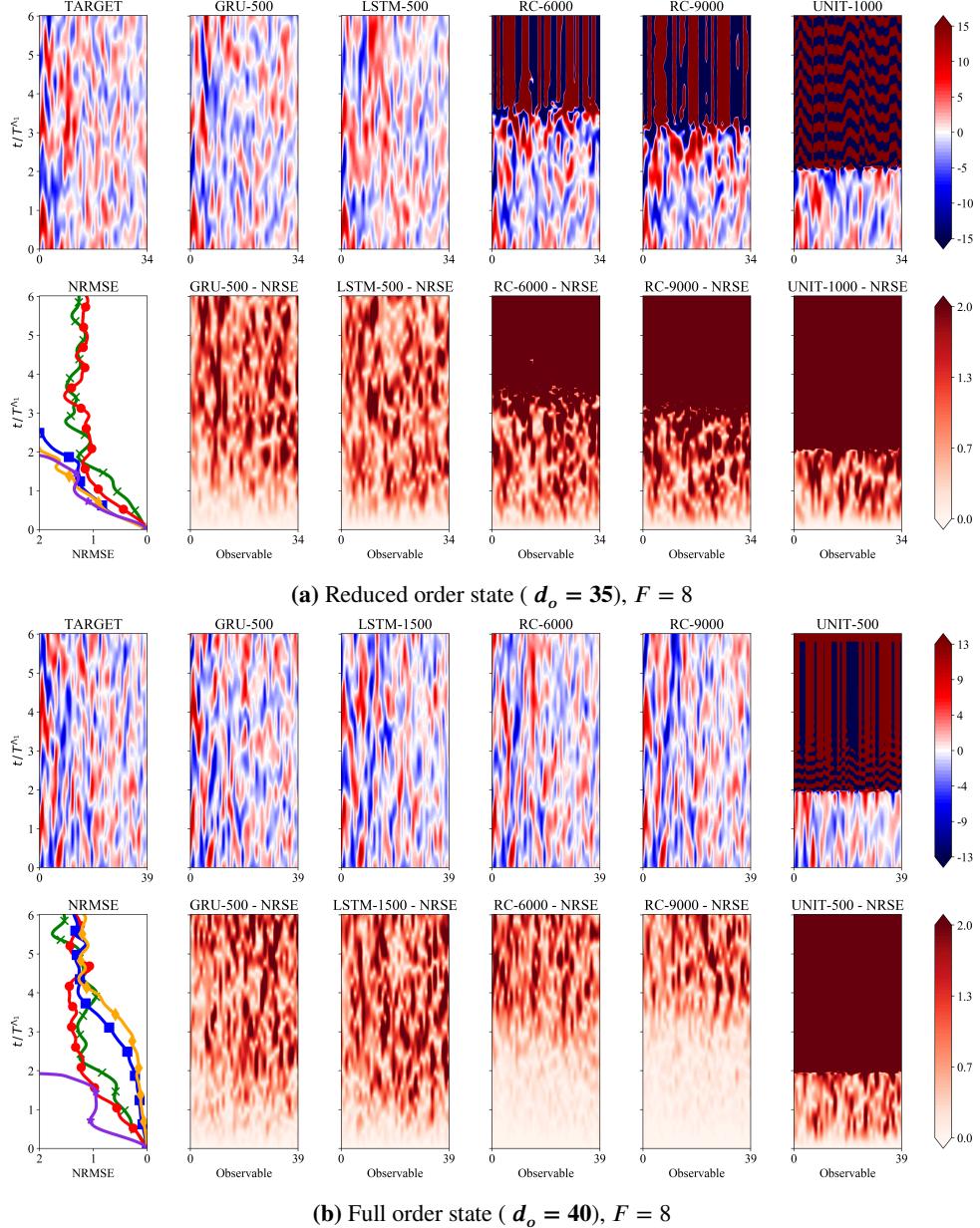
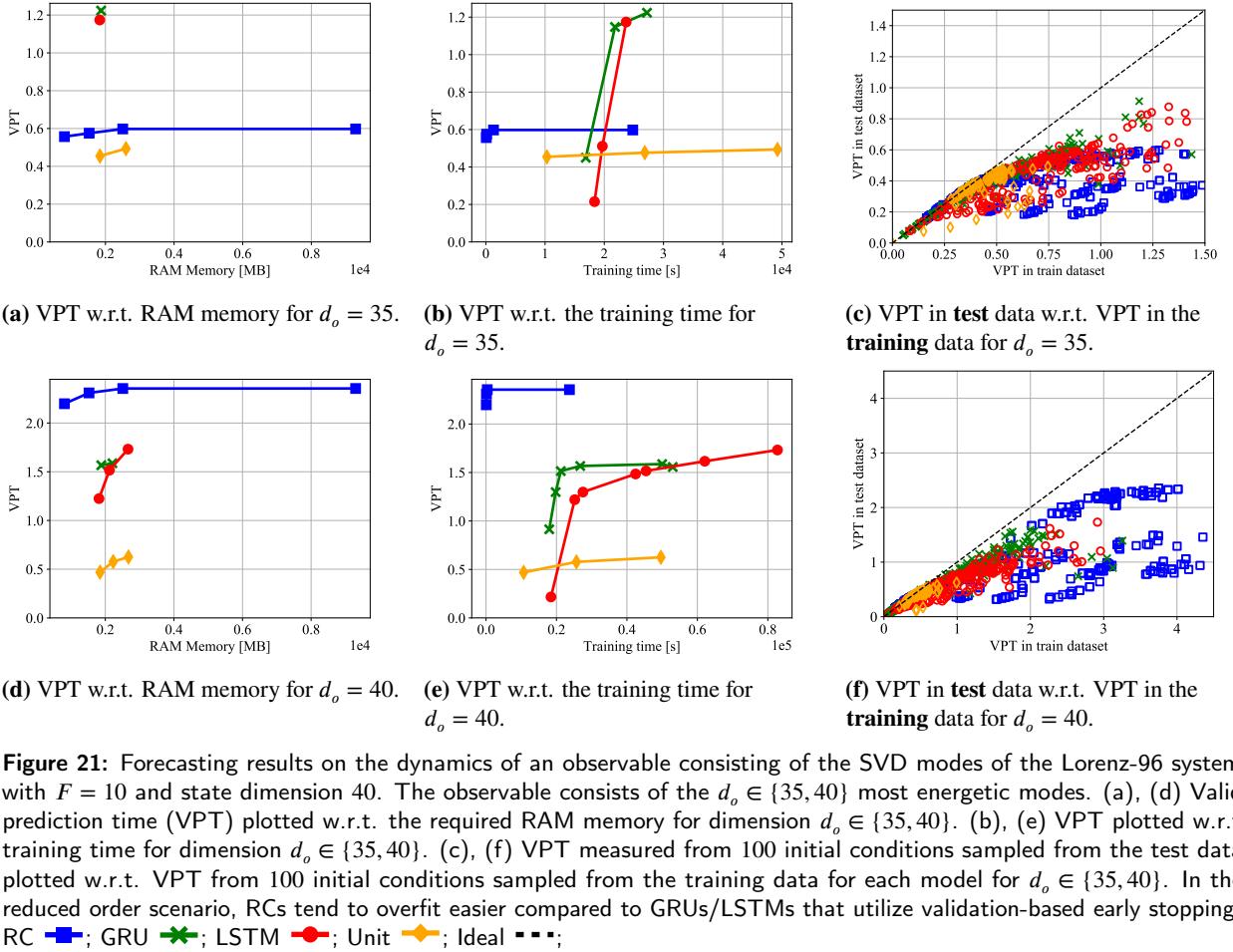


Figure 20: Contour plots of a spatio-temporal forecast on the SVD modes of the Lorenz-96 system with $F = 8$ in the testing dataset with GRU, LSTM, RC and a Unitary network along with the true (target) evolution and the associated NRSE contours for the reduced order observable (a) $d_o = 35$ and the full state (b) $d_o = 40$. The evolution of the component average NRSE (NMRSE) is plotted to facilitate comparison. Unitary networks suffer from propagation of forecasting error and eventually their forecasts diverge from the attractor. Forecasts in the case of an observable dimension $d_o = 40$ diverge slower as the dynamics are deterministic. In contrast, forecasting the observable with $d_o = 35$ is challenging due to both (1) sensitivity to initial condition and (2) incomplete state information that requires the capturing of temporal dependencies. In the full-state setting, RC models achieve superior forecasting accuracy compared to all other models. In the challenging reduced order scenario, LSTM and GRU networks demonstrate a stable behavior in iterative prediction and reproduce the long-term statistics of the attractor. In contrast, in the reduced order scenario RC suffer from frequent divergence. The divergence effect is illustrated in this chosen initial condition.

GRU ; LSTM ; RC-6000 ; RC-9000 ; Unit



G. Additional Results - Lorenz-96 - $F = 10$

In Figure 21, we provide additional results for the forcing regime $F = 10$ that are in agreement with the main conclusions drawn in the main manuscript for the forcing regime $F = 8$. An example of a single forecast of the models starting from an initial condition in the **test** dataset is given in Figure 22

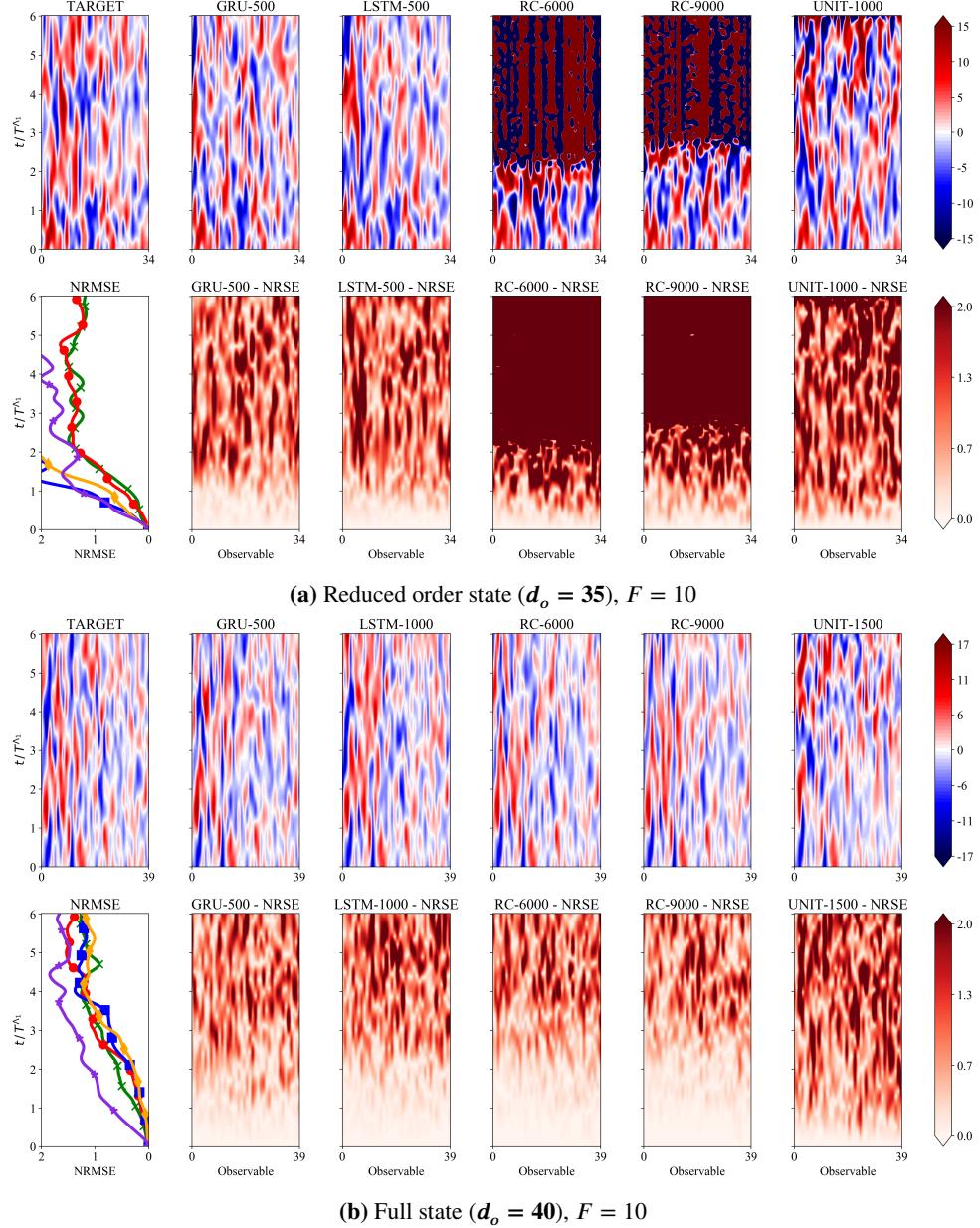


Figure 22: Contour plots of a spatio-temporal forecast on the SVD modes of the Lorenz-96 system with $F = 10$ in the testing dataset with GRU, LSTM, RC and a Unitary network along with the true (target) evolution and the associated NRSE contours for the reduced order observable (a) $d_o = 35$ and the full state (b) $d_o = 40$. The evolution of the component average NRSE (NMRSE) is plotted to facilitate comparison. Forecasting the observable with $d_o = 35$ is more challenging compared to the full state scenario ($d_o = 40$) due to both (1) sensitivity to initial condition and (2) incomplete state information that requires the capturing of temporal dependencies. For this reason, the iterative prediction error increases slower in $d_o = 40$. Even in the challenging scenario of $d_o = 35$, LSTM and GRU networks demonstrate a stable behavior in iterative prediction and reproduce the long-term statistics of the attractor. In the reduced order setting $d_o = 35$, accurate short-term predictions can be achieved with very large RC networks ($d_h = 9000$) at the cost of high memory requirements. However, even in this case, RC models may diverge from the attractor and do not reproduce the attractor climate.

GRU ; LSTM ; RC-6000 ; RC-9000 ; Unit

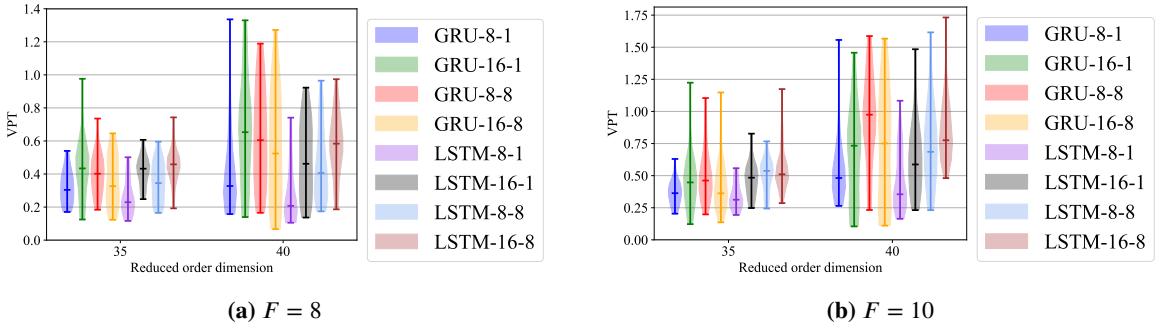


Figure 23: Violin plots of the VPT in the testing data for stateful LSTM and GRU models trained with different truncated Backpropagation through time parameters κ_1 and κ_2 in the (reduced) SVD mode observable of the Lorenz-96 system. The legend of each plot reports the models along with their $\kappa_1 - \kappa_2$ parameters used to train them. The three markers report the minimum, mean and maximum VPT. We observe that especially in the reduced order observable scenario ($d_0 = 35$), having a large truncated back-propagation parameter κ_1 (also referred to as sequence length) is vital to capture the temporal dependencies in the data and achieve high forecasting efficiency. In contrast in the full-state scenario ($d_0 = 40$) a model with a small back-propagation horizon suffices.

H. Temporal Dependencies and Backpropagation

In our study, in order to train the GRU and LSTM models with back-propagation through time (BPTT), we need to tune the parameters κ_1 and κ_2 . The first one denotes the truncated back-propagation length (also referred to as sequence length) and the second the number of future time-steps used to compute the loss and backpropagate the gradient during training at each batch. In the hyperparameter study considered in this work, we varied $\kappa_1 \in \{8, 16\}$ and $\kappa_2 \in \{1, 8\}$. For each of these hyperparameter sets, we varied all other hyperparameters according to the grid search reported in Appendix E.

In Figure 23 we present a violin plot, that illustrates the forecasting efficiency of LSTM and GRU models trained with the listed κ_1 and κ_2 (legend of the plot), while varying all other hyperparameters. The forecasting efficiency is quantified in terms of the Valid Prediction Time (VPT) in the test dataset (averaged over 100 initial conditions) on the Lorenz-96 system for $F \in \{8, 10\}$. The three bars in each violin plot, denote the minimum, average and maximum performance.

In the reduced order scenario case, we observe that models with a large sequence length κ_1 and a large prediction length κ_2 are pivotal in order to achieve a high forecasting efficiency. This implies that there are temporal correlations in the data that cannot be easily captured by other models with smaller horizons. In contrast, in the full order scenario, models with smaller κ_1 perform reasonably well, demonstrating that the need of capturing temporal correlations in the data in order to forecast the evolution is less prominent, since the full information of the state of the system is available.

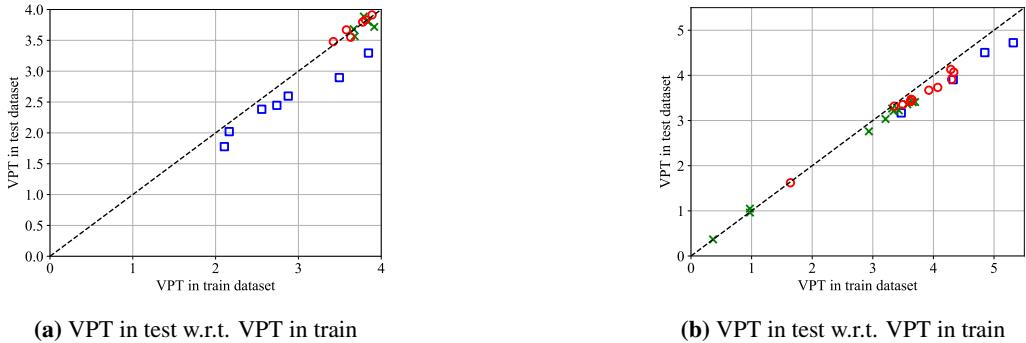


Figure 24: The average VPT measured from 100 initial conditions sampled from the **test** dataset is plotted against the average VPT measured from 100 initial conditions sampled from the **training** dataset for parallel models forecasting the dynamics of (a) the Lorenz-96 system with state dimension $d_o = 40$ and (b) the Kuramoto Sivashinsky equation with state dimension $d_o = 512$.
 RC \square ; GRU \times ; LSTM \circ ;

I. Over-fitting of Parallel Models

In this section, we provide results on the overfitting of the models in the parallel setting in the Lorenz-96 model in Figure 24a and the Kuramoto-Sivashinsky equation in Figure 24b. In both cases we do not observe overfitting issues as the predictive performance in terms of the VPT of the models in the test dataset is very close to that in the training dataset, emphasizing that the generalization ability of the models is good.