

Robust Optimization and Validation of Echo State Networks for learning chaotic dynamics

Alberto Racca^a, Luca Magri^{a,b,c,*}

^a*Department of Engineering, University of Cambridge, Trumpington Street, Cambridge CB2 1PZ, UK*

^b*The Alan Turing Institute, 96 Euston Road, London, England, NW1 2DB, UK*

^c*Institute for Advanced Study, Technical University of Munich, Lichtenbergstrasse 2a, 85748 Garching, Germany (visiting fellowship)*

Abstract

The time-accurate prediction of a chaotic system is challenging because its evolution becomes unpredictable after the predictability time. This is because infinitesimal errors in a chaotic system increase exponentially, i.e., two nearby time series diverge from each other. An approach to the time-accurate prediction of chaotic solutions is by learning temporal patterns from data. Echo State Networks (ESNs), which are a class of Reservoir Computing, can accurately predict the chaotic dynamics well beyond the predictability time. Existing studies, however, also showed that small changes in the hyperparameters may markedly affect the network’s performance. The overarching aim of this paper is to assess and improve the robustness of Echo State Networks for the time-accurate prediction of chaotic solutions. The goal is three-fold. First, we investigate the robustness of routinely used validation strategies. Second, we propose the *Recycle Validation*, and the *chaotic versions* of existing validation strategies, to specifically tackle the forecasting of chaotic systems. Third, we compare Bayesian optimization with the traditional Grid Search for optimal hyperparameter selection. Numerical tests are performed on two prototypical nonlinear systems that have both chaotic and quasiperiodic solutions. Both model-free and model-informed Echo State Networks are analysed. By comparing the network’s robustness in learning chaotic (unpredictable) versus quasiperiodic (predictable) solutions, we highlight fundamental challenges in learning chaotic solutions.

The proposed validation strategies, which are based on the dynamical systems properties of chaotic time series, are shown to outperform the state-of-the-art validation strategies. Because the strategies are principled—they are based on chaos theory such as the Lyapunov time—they can be applied to other Recurrent Neural Networks architectures with little modification. This work opens up new possibilities for the robust design and application of Echo State Networks, and Recurrent Neural Networks, to the time-accurate prediction of chaotic systems.

Keywords: Chaotic dynamical systems, Reservoir Computing, Robustness

1. Introduction

Chaotic systems naturally appear in many branches of science and engineering, from turbulent flows [e.g., 1, 2, 3], through vibrations [4], electronics and telecommunications [5], quantum mechanics [6], reacting flows [7, 8], to epidemic modelling [9], to name only a few. The time-accurate computation of chaotic systems is hindered by the “butterfly effect” [10]: an error in the system’s knowledge—e.g. initial conditions and parameters—grows exponentially until nonlinear saturation. Practically, it is not possible to time-accurately predict chaotic solutions after a time scale, known as the predictability time. The predictability

*Corresponding author

Email address: `lm547@cam.ac.uk` (Luca Magri)

time scales with the inverse of the dominant Lyapunov exponent, which is typically a small characteristic scale of the system under investigation [2].

An approach to the prediction of chaotic dynamics is data-driven. Given a time series (data), we wish to learn the underlying chaotic dynamics to predict the future evolution. The data-driven approach, also known as model-free, traces back to the delay coordinate embedding by Takens [11], which is widely used in time series analysis, in particular, in low-dimensional systems [12]. An alternative data-driven approach to inferring (or, equivalently, learning) chaotic dynamics from data is machine learning. Machine learning is establishing itself as a paradigm that is complementary to first-principles modelling of nonlinear systems in computational science and engineering [13]. In the realm of neural networks, which is the focus of this paper, the feed-forward neural network is the archetypical architecture, which may excel at classification and regression problems [14]. The feed-forward neural network, however, is not the optimal architecture for chaotic time series forecasting because it not designed to learn temporal correlations. Specifically, in time series forecasting, inputs and outputs are ordered sequentially, in other words, they are temporally correlated. To overcome the limitations of feed-forward neural networks, Recurrent Neural Networks (RNNs) [15] have been designed to learn temporal correlations. Examples of successful applications span from speech recognition [16], through language translation [17], fluids [18, 19, 20, 21, 22], to thermo-acoustic oscillations [23], among many others. RNNs take into account the sequential nature of the inputs by updating a hidden time-varying state through an internal loop. As a result of the long-lasting time dependencies of the hidden state, however, training RNNs with Back Propagation Through Time [24] is notoriously difficult. This is because the repeated backwards multiplication of intermediate gradients cause the final gradient to either vanish or become unbounded depending on the spectral radius of the gradient matrix [? 25]. This makes the training ill-posed, which may negatively affect the computational of the optimal set of weights. To overcome this problem, two main types of RNN architectures have been proposed: Gated Structures and Reservoir Computing. Gated Structures prevent gradients from vanishing or becoming unbounded by regularizing the passage of information inside the network, as accomplished in architectures such as Long Short-Term Memory (LSTM) networks [26] and Gated Recurrent Units (GRU) networks [27]. Alternatively, in Reservoir Computing (RC) [28, 29], a high-dimensional dynamical system, the reservoir, acts both as a nonlinear expansion of the inputs and as the memory of the system [30]. At each time step, the output is computed as a linear combination of the reservoir state's components, the weights of which are the only trainable parameters of the machine. Training is, therefore, reduced to a linear regression problem, which bypasses the issue of repeated gradients multiplication in RNNs.

In chaotic attractors, Reservoir Computing has been employed to achieve at least four different goals: to (i) learn ergodic properties, such as Lyapunov exponents [31, 32] and statistics [31, 23]; (ii) filter out noise to recover the deterministic dynamics [33], (iii) reconstruct unmeasured (hidden) variables [34, 35, 36] and (iv) time-accurately predict the dynamics [37, 38, 39]. In this work, we focus on the time-accurate short term prediction of chaotic attractors. A successful Reservoir Computing architecture is the Echo State Network (ESN) [28], which is a universal approximator [40, 41] suitable for the prediction of chaotic time series [37]. There are two broad categories of Echo State Networks: model-free [37, 30] and model-informed [42, 38]. On the one hand, in model-free ESNs, which are the original networks, the training is performed on data only [30]. On the other hand, in model-informed ESNs, the governing equations, or a reduced-order form of them, are embedded in the architecture, for example, in the reservoir in hybrid ESNs [42], or in the loss function in physics-informed ESNs [38]. In chaotic time series forecasting, model-informed ESNs typically outperform model-free ESNs [42, 38, 39]. Both model-free and model-informed Echo State Networks perform as well as LSTMs and GRUs, requiring less computational resources for training [43, 44]. The robustness of ESNs for chaotic time series, however, has not been fully investigated yet, which motivates the overarching objective of this study. Two key aspects may affect ESN robustness. The first aspect is random the initialization, which is required to create the reservoir [30]. Networks with different initializations may perform substantially differently, even after hyperparameter tuning [45]. For an ESN to be robust, network testing through an ensemble of network realizations is required. The second aspect is high hyperparameter sensitivity [30, 46]. The most common validation strategy to compute the hyperparameters for learning chaotic dynamics is the Single Shot Validation, which minimizes the error in an interval subsequent to the training interval. Other validation strategies have been investigated, such as the Walk Forward Validation

and the K-Fold cross Validation [47], but this study was restricted to non-chaotic systems. The computation of the optimal set of hyperparameters is typically performed by Grid Search [28, 21, 42, 37], although other optimization strategies such as Evolutionary Algorithms [48, 49], Stochastic Gradient Descent [50], Particle Swarm Optimization [51] and Bayesian Optimization [52] have been proposed. In particular, Bayesian Optimization (BO) has proved to improve the performance of reservoir-computing architectures in the prediction of chaotic time series, outperforming the commonly used Grid Search strategy [53]. Bayesian Optimization is a gradient-free search strategy, thereby, it is less sensitive to local minima with respect to gradient descent methods [52, 50]. Moreover, Bayesian Optimization is based on Gaussian Process (GP) regression [54], therefore, it naturally quantifies the uncertainty on the computation.

The objective of this paper is three-fold with a focus on learning chaotic dynamics from data. First, we investigate the robustness of the Single Shot Validation, Walk Forward Validation and the K-Fold cross Validation. Second, we propose the Recycle Validation and the chaotic version of existing validation strategies to specifically tackle the forecasting of chaotic systems. Third, we analyse Bayesian optimization for optimal hyperparameter selection. The Lorenz system [10] and the Kuznetsov oscillator [55] are considered as prototypical low-order nonlinear deterministic systems. We highlight fundamental challenges in the robustness of ESNs for chaotic solutions with a comparative investigation on quasiperiodic oscillations. Both model-free and model-informed architectures are analysed.

The paper is organized as follows. Section 2 presents the model-free and model-informed Echo State Network architectures. Section 3 describes the validation strategies. Section 4 investigates the robustness of the Single Shot Validation in forecasting chaotic time series. Section 5 analyses the new validation strategies to improve the robustness in forecasting chaotic time series. Section 6 investigates the robustness of the validation strategies in forecasting quasiperiodic time series. Finally, we summarize the results of this study and discuss future work in the conclusions (section 7).

2. Echo State Networks

As shown in Fig. 1, in the Echo State Network, at any time t_i the input vector, $\mathbf{u}_{\text{in}}(t_i) \in \mathbb{R}^{N_u}$, is mapped into the reservoir state, by the input matrix, $\mathbf{W}_{\text{in}} \in \mathbb{R}^{N_r \times N_u}$, where $N_r \gg N_u$. The reservoir state, $\mathbf{r} \in \mathbb{R}^{N_r}$, is updated at each time iteration as a function of the current input and its previous value

$$\mathbf{r}(t_{i+1}) = \tanh(\mathbf{W}_{\text{in}}\mathbf{u}_{\text{in}}(t_i) + \mathbf{W}\mathbf{r}(t_i)), \quad (1)$$

where $\mathbf{W} \in \mathbb{R}^{N_r \times N_r}$ is the state matrix. The predicted output, $\mathbf{u}_{\text{p}}(t_{i+1}) \in \mathbb{R}^{N_u}$, is obtained as

$$\mathbf{u}_{\text{p}}(t_{i+1}) = \hat{\mathbf{r}}(t_{i+1})^T \mathbf{W}_{\text{out}}, \quad \hat{\mathbf{r}}(t_{i+1}) = \mathbf{g}(\mathbf{r}(t_{i+1})); \quad (2)$$

where $\mathbf{g}(\cdot)$ is a nonlinear transformation, $\hat{\mathbf{r}} \in \mathbb{R}^{N_r}$ is the updated reservoir state, and $\mathbf{W}_{\text{out}} \in \mathbb{R}^{N_r \times N_u}$ is the output matrix. The input matrix, \mathbf{W}_{in} , and state matrix, \mathbf{W} , are (pseudo)randomly generated and fixed, while the weights of the output matrix, \mathbf{W}_{out} , are computed by training the network. In this work, the input matrix, \mathbf{W}_{in} , has only one element different from zero per row, which is sampled from a uniform distribution in $[-\sigma_{\text{in}}, \sigma_{\text{in}}]$, where σ_{in} is the input scaling. The state matrix, \mathbf{W} , is an Erdős-Renyi matrix with average sparseness s , in which each neuron (each row of \mathbf{W}) has on average only $(1-s)N_r$ connections (non-zero elements). The non-zero elements are obtained by sampling from a uniform distribution in $[-1, 1]$; the entire matrix is then rescaled by a multiplication factor to set the spectral radius, ρ . The spectral radius is key to enforcing the *echo state property*. (In a network with the echo state property, the state loses its dependence on its previous values for sufficiently large times and, therefore, it is uniquely defined by the sequence of inputs.) While the echo state property may hold for a wider range of spectral radii [56], the condition $\rho < 1$ ensures the echo state property in most situations [30].

The ESN can be run either in open-loop or closed-loop configuration. In the open-loop configuration, first, we feed data as the input at each time step to compute and store $\hat{\mathbf{r}}(t_i)$ (1-2). In the initial transient of this

process, the washout interval, we do not compute the output, $\mathbf{u}_p(t_i)$. The purpose of the washout interval is for the reservoir state to satisfy the echo state property, thereby becoming independent of the arbitrarily chosen initial reservoir state, $\mathbf{r}(t_0) = 0$. Secondly, we train the output matrix, \mathbf{W}_{out} , by minimizing the Mean Square Error (MSE) between the outputs, $\mathbf{u}_p(t_i)$, and the data, $\mathbf{u}_d(t_i)$, over a training set of N_{tr} points

$$\text{MSE} \triangleq \frac{1}{N_{\text{tr}}N_u} \sum_i^{N_{\text{tr}}} \|\mathbf{u}_p(t_i) - \mathbf{u}_d(t_i)\|^2, \quad (3)$$

where $\|\cdot\|$ is the L_2 norm. Minimizing (3) is a least-squares minimization problem, which can be solved as a linear system through ridge regression

$$(\mathbf{R}\mathbf{R}^T + \beta\mathbf{I})\mathbf{W}_{\text{out}} = \mathbf{R}\mathbf{U}_d^T, \quad (4)$$

where $\mathbf{R} \in \mathbb{R}^{N_{\hat{\mathbf{r}}} \times N_{\text{tr}}}$ and $\mathbf{U}_d \in \mathbb{R}^{N_u \times N_{\text{tr}}}$ are the horizontal concatenation of the updated reservoir states, $\hat{\mathbf{r}}$, and the data, \mathbf{u}_d , respectively; \mathbf{I} is the identity matrix and β is the user-defined Tikhonov regularization parameter [57]. We solve the linear system through the `linalg.solve` function in `numpy` [58]. In the closed-loop configuration, starting from an initial data point as an input and an initial reservoir state obtained after the washout interval, the output, \mathbf{u}_p , is fed back to the network as an input for the next time step prediction. In doing so, the network is able to autonomously evolve in the future. The closed-loop configuration is used during validation and testing.

2.1. Model-free and model-informed architectures

We consider model-free and model-informed architectures (Fig. 1). The basic model-free ESN is obtained by setting $\mathbf{g}(\mathbf{r}) = \mathbf{r}$. This architecture, however, generates symmetric solutions in the closed-loop configuration [34, 23], which can cause the predicted trajectory to stray away from the actual attractor towards a symmetric attractor, which is not a solution of the dynamical system (but it is a solution of the network). To break the symmetry, we add biases in the input and output layers

$$\mathbf{r}_{i+1} = \tanh(\mathbf{W}_{\text{in}}[\mathbf{u}_{\text{in}}; b_{\text{in}}] + \mathbf{W}\mathbf{r}_i), \quad \hat{\mathbf{r}}_i = [\mathbf{r}_{i+1}; 1], \quad \mathbf{u}_p = \hat{\mathbf{r}}_i^T \mathbf{W}_{\text{out}}. \quad (5)$$

where $[\cdot; \cdot]$ indicates vertical concatenation, b_{in} is the scalar input bias and $\mathbf{W}_{\text{in}} \in \mathbb{R}^{N_r \times (N_u + 1)}$. In the model-informed ESN, also known as hybrid as proposed by [42], information about the governing equations (model knowledge) is embedded into the model through a function of the input, $\mathcal{K}(\mathbf{u}_{\text{in}})$, which, for example, may be a reduced order model that provides information about the output at the next time step as

$$\hat{\mathbf{r}}_i = [\mathbf{r}_{i+1}; 1; \mathcal{K}(\mathbf{u}_{\text{in}})]. \quad (6)$$

In this work, we use $\mathcal{K}(\mathbf{u}_{\text{in}})$ only to update the reservoir state [39], in order to use the same input matrix, \mathbf{W}_{in} , and state matrix, \mathbf{W} , of the model-free architecture. This allows us to directly compare the performances of the model-free and model-informed architectures.

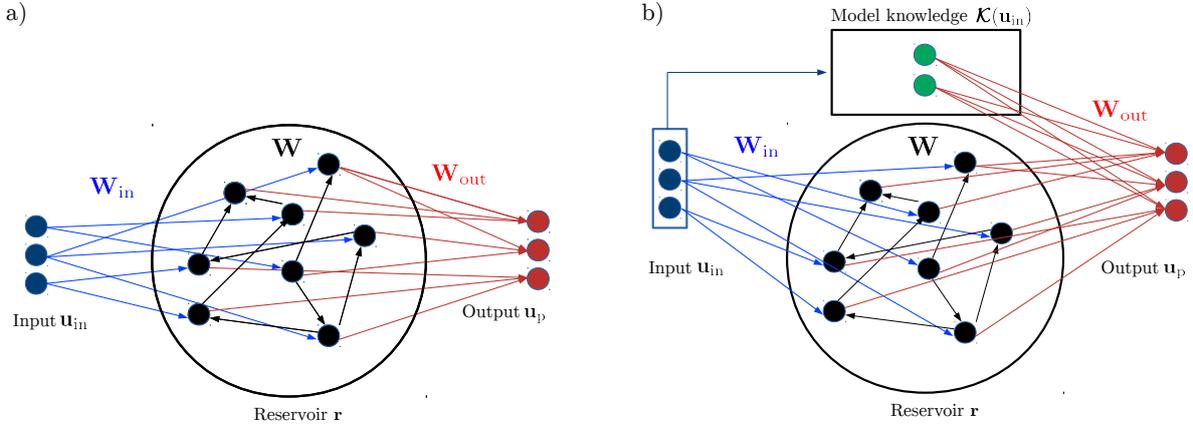


Figure 1: Schematic representation of (a) model-free and (b) model-informed Echo State Networks (ESNs).

3. Validation

The purpose of the validation is to determine the hyperparameters by minimizing an error. We make a distinction between the hyperparameters (i) that require re-initialization of \mathbf{W}_{in} and \mathbf{W} , and (ii) that do not require re-initialization. The size of the reservoir, N_r , and sparseness, s , require re-initialization, whereas the input scaling, σ_{in} , the spectral radius, ρ , the Tikhonov parameter, β , and the input bias b_{in} , do not. The fundamental difference between (i) and (ii) is that the random component of the re-initialization of \mathbf{W}_{in} and \mathbf{W} makes the objective function to be minimized random, which significantly increases the complexity of the optimization. In this study, we minimize the error with respect to the input scaling, σ_{in} , and spectral radius, ρ , which are key hyperparameters for the performance of the network [46, 30]. For convenience, we rewrite the reservoir state equation (1) as

$$\mathbf{r}_{i+1} = \tanh(\sigma_{in} \hat{\mathbf{W}}_{in}[\mathbf{u}_{in}; b_{in}] + \rho \hat{\mathbf{W}} \mathbf{r}_i), \quad (7)$$

where the non-zero elements of $\hat{\mathbf{W}}_{in}$ are sampled from the uniform distribution in $[-1, 1]$ and $\hat{\mathbf{W}}$ has been scaled to have a unitary spectral radius.

3.1. Performance metrics

We determine the hyperparameters by minimizing the Mean Squared Error (3) in the validation interval of fixed length. The networks are tested on multiple starting points along the attractor by using both the Mean Squared Error and Prediction Horizon (PH), the latter of which is defined as the time interval during which the normalized error is smaller than a user-defined threshold k [38, 42]

$$\frac{\|\mathbf{u}_p(t_i) - \mathbf{u}_d(t_i)\|}{\sqrt{\frac{1}{N_{PH}} \sum_j^{N_{PH}} \|\mathbf{u}_d(t_j)\|^2}} < k, \quad (8)$$

where N_{PH} are the number of timesteps in the Prediction Horizon. The Mean Squared Error and Prediction Horizon for the same starting point in the attractor are strictly correlated (Appendix A). We use the Mean Squared Error to partition the dataset in intervals of fixed length during validation, while we use the Prediction Horizon in the test set because it is the most physical quantity when assessing the time-accurate prediction of chaotic systems [e.g., 37, 33].

3.2. Strategies

The most common validation strategy for ESNs is the Single Shot Validation (SSV), which splits the available data in a training set and a single subsequent validation set (Fig. 2a). The time interval of the validation set, during which the hyperparameters are tuned, is small and represents only a fraction of the attractor. In nonlinear time series prediction, the choice of the validation strategy has to take into account (i) the intervals we are interested in predicting and (ii) the nature of the signal we are trying to reproduce. Here, we are interested in predicting multiple intervals as the trajectory spans the attractor, rather than a specific interval starting from a specific initial condition. Moreover, an ergodic trajectory of the attractor has no underlying time-varying statistics, e.g there is no time-dependency of the mean of the signal, hence trajectories return indefinitely in nearby regions of the attractor. This means that the intervals we are interested in predicting are potentially similar to any interval of the trajectory that constitutes our dataset, regardless of the interval position in time within the dataset. For this reason, as shown in section 4, the Single Shot Validation strategies should not be employed in chaotic time series.

These observations lead us to use validation strategies based on multiple validation intervals, which may precede the training set, such as the Walk Forward Validation (WFFV) and the K-Fold cross Validation (KFV). We also propose an ad-hoc robust validation strategy—the Recycle Validation (RV). The objective of these strategies is to tune the hyperparameters over an effectively larger portion of the trajectory, by minimizing the average of the objective function (error) over multiple validation intervals. The regular version of these strategies consists of creating subsequent folds by moving forward in time the validation set by its own length. Additionally, we propose the chaotic version, in which we move the fold forward in time by one Lyapunov Time (LT) (Fig. 2). The Lyapunov Time is a key time scale in chaotic dynamical systems, which is defined as the inverse of the leading Lyapunov exponent Λ of the system, which, in turn, is the exponential rate at which infinitesimally close trajectories, $\delta\mathbf{q}(0)$ diverge [e.g., 2]

$$\|\delta\mathbf{q}(t)\| \sim \|\delta\mathbf{q}(0)\| \exp(\Lambda t) \quad t \rightarrow \infty, \quad \|\delta\mathbf{q}(0)\| \rightarrow 0. \quad (9)$$

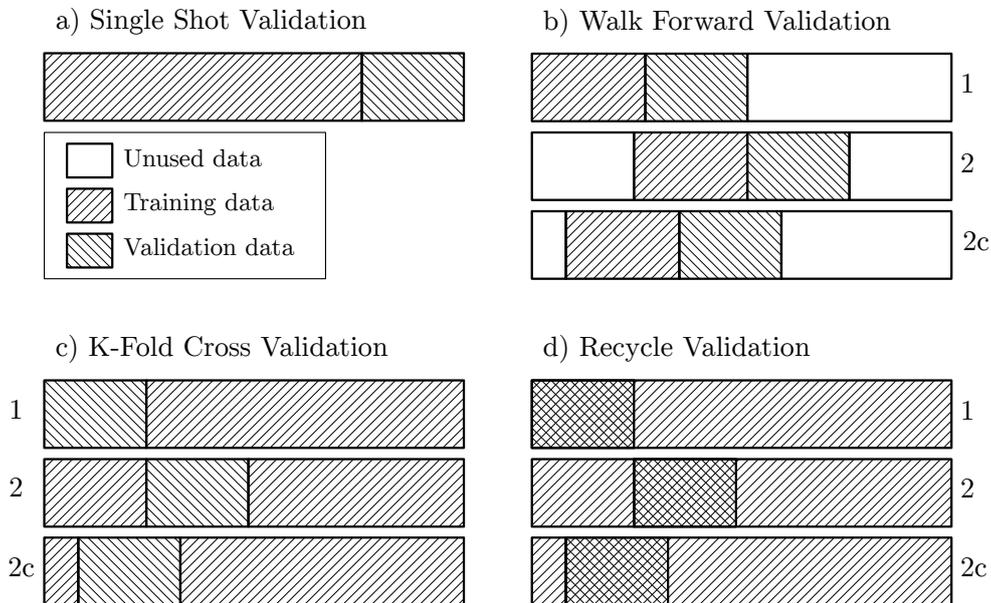


Figure 2: Partition of the data in the different validation strategies. In (b-d), bar 1 shows the first fold, bar 2 shows the second fold, and bar 2c shows the second fold in the chaotic version (shifted by one Lyapunov time).

Walk Forward Validation. In the Walk Forward Validation (WFV) (Fig. 2b), we partition the available data in multiple splits, while maintaining sequentiality of the data. From a starting dataset of length n , the first m points ($m < n$) are taken as the first fold, with N_t points for training and v points for validation ($v + N_t = m$). These quantities must respect $(n - m) = (k_1 - 1)v$; $k_1 \in \mathbb{N}$. The remaining $(k_1 - 1)$ folds are generated by moving the training plus validation set forward in time by a number of points v . This way, the original dataset is partitioned in k_1 folds and the hyperparameters are selected to minimize the average MSE over the folds.

K-Fold cross Validation. Although the K-Fold cross Validation (KFV) (Fig. 2c) is a common strategy in regression and classification, it is not commonly used in time series prediction because the validation and training intervals are not sequential to each other. This strategy partitions the available data in k_2 splits. Over the entire dataset of length n , after an initial bv points, with $0 \leq b < 1$, needed to have an integer number of splits, the remaining $n - bv$ points are used as k_2 validation intervals, each of length v . For each validation interval we define a different fold, in which we use all the remaining data points for training. We determine the hyperparameters by minimizing the average of the MSE between the folds.

Recycle Validation. We propose a the Recycle Validation (RV) (Fig. 2d), which exploits the information obtained by both open-loop and closed-loop configurations. Because the network works in two different configurations, it can obtain additional information when validating on data already used in training. To do so, first, we train the output weights using the entire dataset of n points. Second, we validate the network on k_2 splits of length v from data that has already been used to train the output weights. Each split is imposed by moving forward in time the previous validation interval by v points. After an initial bv points, with $0 \leq b < 1$, needed to have an integer number of splits, the remaining $n - bv$ points are used as k_2 validation intervals. We determine the hyperparameters by computing the average of the MSE between the splits. This strategy has four main advantages. First, it can be used in small datasets, where the partition of the dataset in separate training and validation sets may cause the other strategies to perform poorly. In small datasets, the validation intervals represent a larger percentage of the dataset since each validation interval needs to be multiple Lyapunov Times to capture the divergence of chaotic trajectories. Therefore, the training set becomes substantially smaller than the dataset and the output matrix used during validation differs substantially from the output matrix of the whole dataset. This results in a poor selection of hyperparameters. Second, for a given dataset, we maximize the number of validation splits, using the same validation intervals of the K-Fold cross Validation. Third, we tune the hyperparameters using the same output matrix, \mathbf{W}_{out} , that we use in the test set. Fourth, it has lower computational cost than the K-Fold cross Validation because it does not require retraining the output matrix for the different folds. which makes it computationally cheaper (Appendix B).

Chaotic version. The chaotic version consists of shifting the validation intervals forward in time, not by their own length, but by one Lyapunov Time when constructing the next fold. In doing so, different splits will overlap, but, since the trajectory related to the split that started one Lyapunov Time (LT) earlier has strayed away from the attractor on average by $e^{\Lambda \times 1LT} = e$, the two intervals contain different information. The purpose of this version is to further increase the number of intervals on which the network is validated. The regular and chaotic versions for each validation strategy are shown in frames (b-d) in Fig. 2 in bars 2 and 2c, respectively. The *chaotic versions* of the Walk Forward Validation, the K-fold cross Validation and the Recycle Validation are denoted by the subscript c .

3.3. Grid search and Bayesian optimization

To find the minimum of the Mean Squared Error (3) of the validation set in the hyperparameter space, we use Bayesian Optimization (BO), which is compared to Grid Search (GS). Bayesian Optimization has been shown to outperform other state-of-the-art optimization methods when the number of evaluations of an expensive objective function is limited [59, 60]. It is a global search method, which is able to incorporate prior knowledge about the objective function and to use information from the entire search space. It treats the objective function as a black box, therefore, it does not require gradient information. Starting from

an initial N_{st} evaluations of the objective function, BO performs a Gaussian Process (GP) regression [54] to reconstruct the function in the search space, using function evaluations as data. Once the GP fitting is available, we select the new point at which to evaluate the objective function so that the new point maximizes the acquisition function. The acquisition function is evaluated on the mean and standard deviation of the GP reconstruction. After the objective function is evaluated at a new point, the enlarged data set, comprising of the new point, is used to perform another GP regression, select a new point and so on and forth. In this work, we use the gp-hedge Bayesian Optimization algorithm implemented in `scikit-optimize` library in Python [61, 62]. The details of the formulation are explained in Appendix C and the Supplementary Material (S.2).

4. Robustness of the Single Shot Validation

The first testcase we investigate is the Lorenz system [10], which is a reduced-order model of Rayleigh–Bénard convection

$$\begin{aligned}\dot{x} &= \sigma_L(y - x) \\ \dot{y} &= x(\rho_L - z) - y \\ \dot{z} &= xy - \beta_L z,\end{aligned}\tag{10}$$

where $[\sigma_L, \beta_L, \rho_L] = [10, 8/3, 28]$ is selected to generate chaotic solutions [e.g., 10]. The system is integrated with a forward Euler scheme with step $dt = 0.009$ LT. The Lyapunov Time is $LT = \Lambda^{-1} \approx 1.1$ [63].

We analyse the performance of the used Single Shot Validation (SSV), which is employed for training (1LT to 9LTs), validation (9LTs to 12 LTs), and testing (12 LTs to 15LTs), as shown in Fig. 3. The input, \mathbf{u}_{in} , is normalized by its maximum variation. (This is done because we are using a single scalar quantity σ_{in} to scale all the components of the input.) The network has a fixed number of neurons, $N_r = 100$, sparseness, $s = 97\%$, Tikhonov parameter, $\beta_t = 10^{-11}$ and input bias, $b_{in} = 1$ [30]. The bias, b_{in} , is set for it to have the same order of magnitude of the normalized input. The input scaling, σ_{in} , and spectral radius, ρ , are tuned during validation in the range $[0.5, 5] \times [0.1, 1]$ to minimize the $\log_{10}(\text{MSE})$. The range of the spectral radius, ρ , is selected for the network to respect the echo state property, whereas the range of the input scaling, σ_{in} , is selected to normalize the inputs. The optimization is performed with (i) Grid Search (GS) consisting of 7×7 points, and (ii) Bayesian Optimization (BO) consisting of 5×5 starting points and 24 points acquired by the gp-hedge algorithm (Appendix C). The two optimization schemes are applied to an ensemble of $N_{ens} = 50$ networks, which differ by the random initialization of the input matrix, \mathbf{W}_{in} , and the state matrix, \mathbf{W} . N_{ens} is selected after a test on statistical convergence (Appendix D).

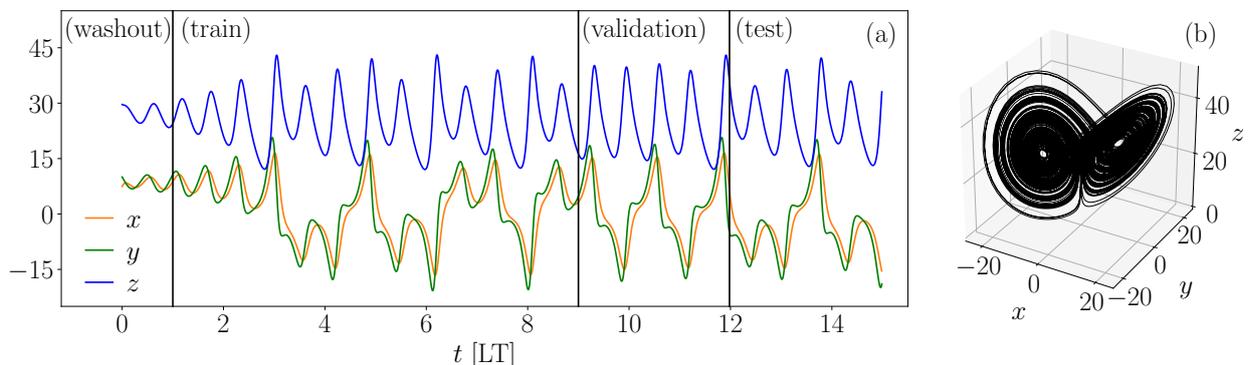


Figure 3: Solution of the Lorenz system. (a) Time series, and (b) phase plot for a longer time window. Time is expressed in Lyapunov time (LT) units.

Figure 4 shows the performance of the optimal hyperparameters computed by Grid Search and Bayesian Optimization for the ensemble members. First, we analyse the performance in validation (panel (a)). As shown by the medians reported in the caption, Bayesian optimization markedly outperforms Grid Search. Second, we analyse the performance in the test set (panel (b)). The performance of each network is assessed by computing the MSE in the test set for the hyperparameters found in the validation set. For this, the output matrix, \mathbf{W}_{out} , of the test set is obtained by retraining over both the training and validation sets. As shown by the medians, the overall performance of the networks and the benefit of using Bayesian Optimization are markedly reduced. This is a signature of chaos, whose unpredictability results in a weak correlation between validation and test sets. This is further verified by computing the mean of the Gaussian process reconstruction from a 30×30 grid of $\log_{10}(\text{MSE})$ for a representative network of the ensemble (Fig. 5). The performance of the optimal hyperparameters of the validation set can deteriorate by four, or more, orders of magnitude in the test set (panel (c)).

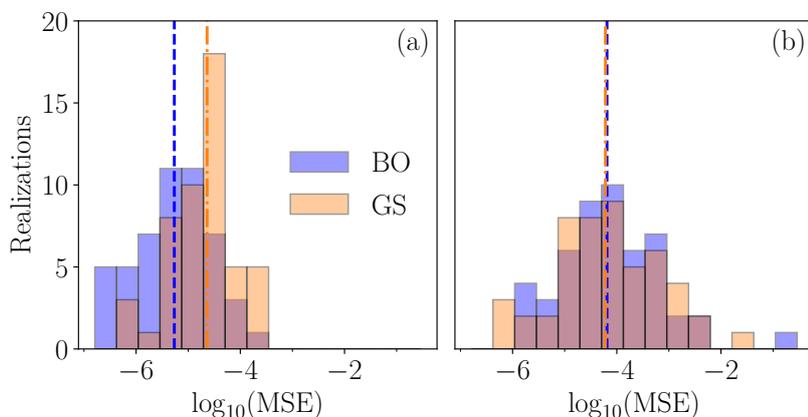


Figure 4: Performance of the optimal hyperparameters computed by Grid Search (GS) and Bayesian Optimization (BO) in (a) validation and (b) test sets. Vertical lines indicate the median of Grid Search (dash-dotted) and Bayesian Optimization (dashed). The medians are $[5.4, 23.0] \times 10^{-6}$ in the validation set and $[64.8, 60.5] \times 10^{-6}$ in the test set for BO and GS, respectively.

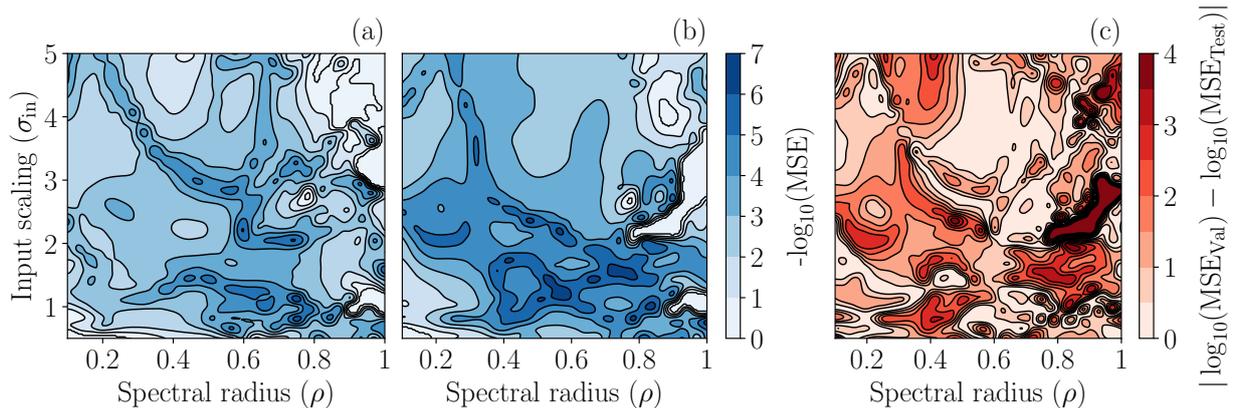


Figure 5: Mean of the Gaussian Process reconstruction in the (a) validation and (b) test sets for a representative network in the ensemble. Frame (c) shows the difference between the two sets. The MSE is saturated to be ≤ 1 in (a,b), whereas the error is saturated to be $\leq 10^4$ in (c). The reconstruction is performed on a grid of 30×30 evaluations of $\log_{10}(\text{MSE})$. For the same hyperparameters, the MSE can differ by orders of magnitude between the validation and test sets.

To assess quantitatively the correlation of the optimal hyperparameters' performance between the vali-

ation and test sets, we use the Spearman coefficient [64]

$$\tilde{r}_S(\mathbf{x}, \mathbf{y}) = \frac{\sum_i (z(x)_i - N_{\text{ens}})(z(y)_i - N_{\text{ens}})}{\sqrt{\sum_i (z(x)_i - N_{\text{ens}})^2} \sqrt{\sum_i (z(y)_i - N_{\text{ens}})^2}},$$

$$\mathbf{x} = \begin{bmatrix} \mathbf{m}_{\text{Val}}^{(\text{BO})} \\ \mathbf{m}_{\text{Val}}^{(\text{GS})} \end{bmatrix}, \quad \mathbf{y} = \begin{bmatrix} \mathbf{m}_{\text{Test}}^{(\text{BO})} \\ \mathbf{m}_{\text{Test}}^{(\text{GS})} \end{bmatrix}, \quad (11)$$

where $\mathbf{z}(\mathbf{x})$ is the ranking function; $\mathbf{m} \in \mathbb{R}^{N_{\text{ens}}}$ contains the MSE for the optimal hyperparameters in validation (subscript Val), or test (subscript Test) obtained by Bayesian Optimization (superscript BO), or Grid Search (superscript GS). \tilde{r}_S quantifies the correlation between the MSE of the optimal hyperparameters obtained during validation and the MSE for the same hyperparameters in the test set over the ensemble. The values $\tilde{r}_S = \{-1, 0, 1\}$ indicate anticorrelation, no correlation and correlation, respectively.

Figure 6 shows the correlation analysis. The scatter plot for \mathbf{x} and \mathbf{y} (panel (a)) shows that the MSE of the optimal hyperparameters in the validation and test sets are weakly correlated with $\tilde{r}_S = 0.32$, independently on whether they are computed with Bayesian Optimization or Grid Search. Panels (b,c) show the values of the optimal hyperparameters, which vary substantially from one network realization to another.

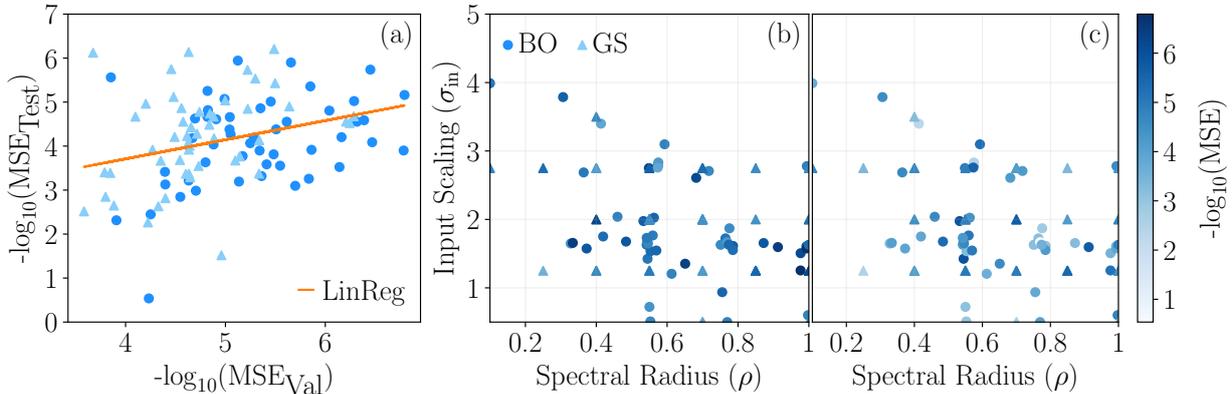


Figure 6: (a) Linear regression (LinReg) and scatter plot of the MSE of the optimal hyperparameters obtained from Bayesian Optimization (BO) and Grid Search (GS) for each network. Optimal hyperparameters for each network and corresponding MSE in (b) validation and (c) test sets. For different networks the optimal hyperparameters, and their performance, vary significantly.

4.1. Remarks

First, because the MSE and optimal hyperparameters vary significantly in different network realizations, we advise performing optimization separately for each network to increase the accuracy (as further verified in Appendix E). Second, hyperparameters that are optimal in the validation set may have a poor performance in the test set, which may greatly reduce the benefit of using Bayesian hyperparameter optimization. This highlights a fundamental challenge in learning chaotic solutions, in which validation and test sets may be topologically different portions of the attractor. We, thus, advise that the Single Shot Validation not be used in the validation of Echo State Networks in chaotic attractors. Robust validation strategies (section 3.2) are next analysed.

5. Validation for chaotic solutions

5.1. Hyperparameter optimization

We compare different validation strategies on the ensemble of $N_{\text{ens}} = 50$ networks in a “short” dataset (12 LTs) and a “long” dataset (24 LTs). The long dataset is obtained by the integration of the time series

in Fig. 3. In addition to the short dataset, we analyse the long dataset for two reasons. First, we wish to test validation strategies that require larger datasets to fully perform, such as the Walk Forward Validation. Second, we wish to investigate how the robustness is affected by the size of the dataset. We use the Single Shot Validation (SSV), Walk Forward Validation (WFV), K-Fold Validation (KFV), Recycle Validation (RV), and corresponding chaotic versions (subscript c). The long dataset allows us to define an additional chaotic Walk Forward Validation (WFV $_c$) denoted by the superscript $*$ as detailed in the Supplementary Material (S.1).

The test set has $N_t = 100$ starting points on the attractor to sample different regions of the solutions (more details in Appendix D). The Prediction Horizon is globally quantified as an arithmetic mean, $\overline{\text{PH}}_{\text{test}}$, with threshold $k = 0.2$; whereas the Mean Squared Error is globally quantified as a geometric mean, $\overline{\text{MSE}}_{\text{Test}}$, in intervals of 3LTs .

5.1.1. Model-free ESN

Figure 7 shows the mean of the Gaussian Process reconstruction of $\log_{10}(\text{MSE})$ in the hyperparameter space for a representative network of the ensemble. Panels (a,b,c) show the performance of three validation strategies in the validation set, whereas panel (d) shows the performance of the network in the test set. Because the error in (b,c) is similar to the error in (d), and the error in (a) differs from (d), we conclude that in the test set the hyperparameters computed through KFV $_c$ and RV $_c$ perform well, but the hyperparameters computed through SSV perform poorly.

A correlation analysis is shown in Tab. 1 with the Spearman correlation coefficients, \tilde{r}_s (11) (short and long datasets); and Fig. 8 with scatter plots of the optimal hyperparameters' performance (long dataset, for brevity). The Single Shot Validation has the lowest correlation among all the validation strategies in both datasets. The chaotic versions of the validation strategies correlate better than the corresponding regular versions. In particular, the chaotic K-Fold Validation and the chaotic Recycle Validation have the highest correlations. In general, increasing the size of the dataset increases the correlation, but the Single Shot Validation in the long dataset has a lower correlation than the K-Fold Validations and the Recycle Validations in the short dataset. This further demonstrates the poor robustness of the Single Shot Validation. Last, but not least, the Recycle Validation is computationally cheaper than the K-Fold Validation because the output matrix is the same for the different folds (more analysis on the computational time can be found in Appendix B).

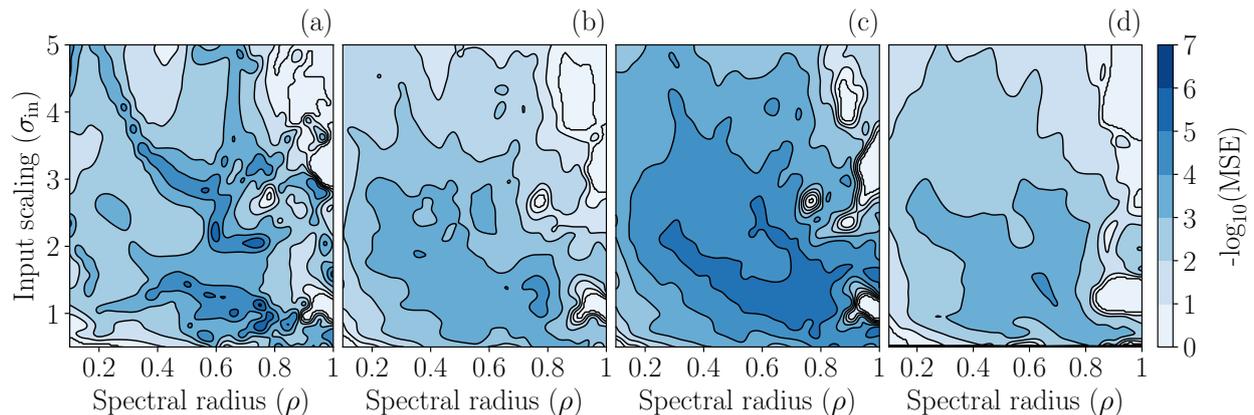


Figure 7: Mean of the Gaussian Process reconstruction for the short dataset for a representative network of the ensemble. Validation set for (a) Single Shot Validation (SSV), (b) chaotic K-Fold Validation (KFV $_c$), and (c) chaotic Recycle Validation (RV $_c$); and test set (d). The MSE is saturated to be ≤ 1 . The Gaussian Process is based on a grid of 30×30 data points.

Table 1: Spearman coefficients between validation and test sets. Bold text indicates the highest correlation in the dataset.

\tilde{r}_S	SSV	WFV	WFV _c	WFV _c *	KFV	KFV _c	RV	RV _c
Short dataset (12LTs)	0.31	0.31	0.50	-	0.60	0.65	0.59	0.62
Long dataset (24LTs)	0.49	0.51	0.61	0.70	0.70	0.85	0.67	0.81

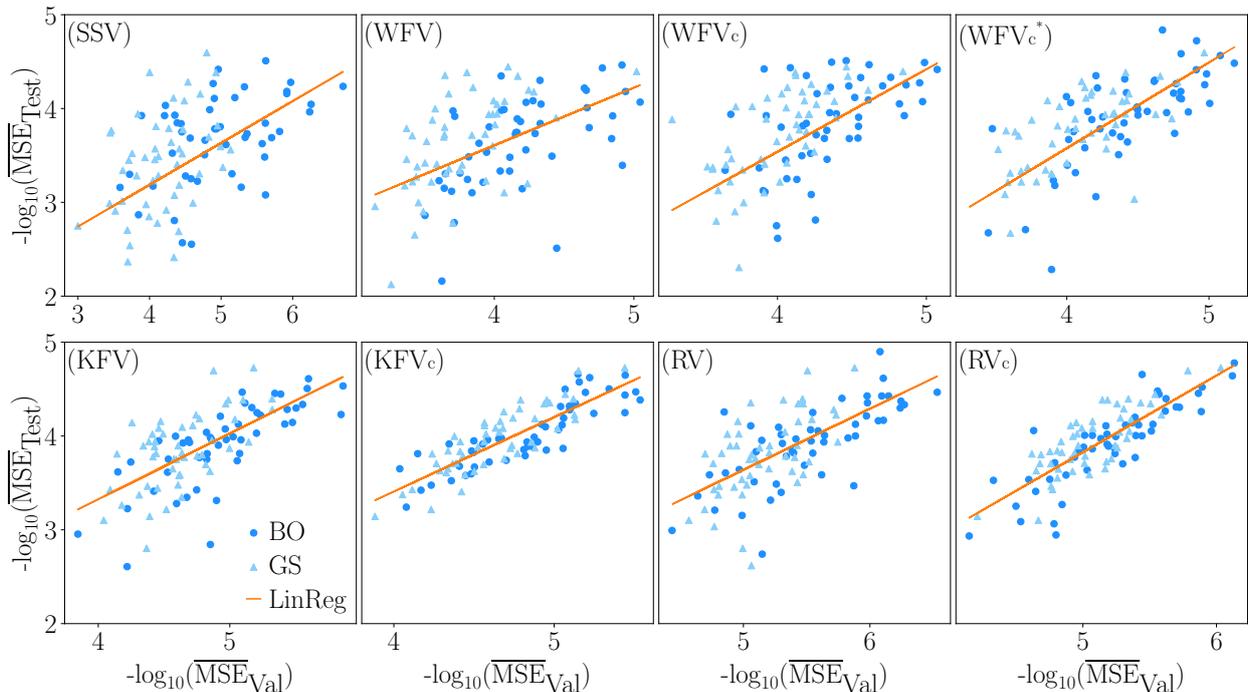


Figure 8: Linear regression (LinReg) and scatter plot of the MSE of the optimal hyperparameters obtained from Bayesian Optimization (BO) and Grid Search (GS) for each network. Single Shot Validation (SSV), Walk Forward Validation (WFV), K-Fold Validation (KFV), Recycle Validation (RV), and their chaotic versions (subscript c). Long dataset.

A comparison between Bayesian Optimization (BO) and Grid Search (GS) is shown in Fig. 9. Panels (a,b) show the ratio of the MSE between the optimal hyperparameters obtained by Bayesian Optimization and Grid Search in the validation and test sets. In both datasets, Bayesian Optimization outperforms Grid Search in the validation set in $\sim 75\%$ of the networks (except for one outlier). However, BO and GS perform similarly in the test set, especially in the short dataset (a). In the long dataset (b), Bayesian Optimization on average outperforms Grid Search, although there is a decrease in performance with respect to the validation set. Panels (c,d) show the Prediction Horizon (PH) in the test set. The chaotic K-Fold Validation and the chaotic Recycle Validation increase the Prediction Horizon by 0.5 LTs on average with respect to the Single Shot Validation. The Prediction Horizon of the long datasets (d) is $\gtrsim 0.5$ LTs larger than that of the short dataset (c). This results in the performance of the KFV_c and RV_c in the short dataset being closer to the performance of the SSV in the long dataset. Because Bayesian Optimization does not produce a substantial increase in the Prediction Horizon with respect to Grid Search, we conclude that the performance of the networks is more sensitive to the validation strategy rather than the optimization scheme.

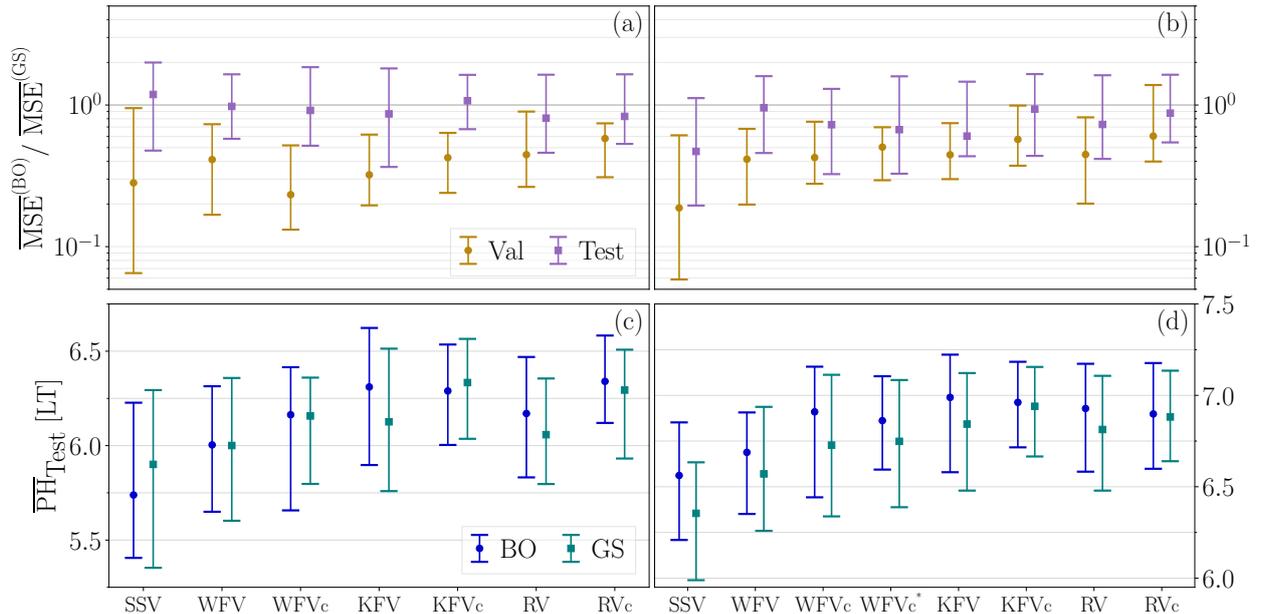


Figure 9: Comparison between hyperparameter optimization by Bayesian Optimization (BO) and Grid Search (GS). The performance metrics are the Mean Square error (MSE) and Predictability Horizon (PH). 25th (lower bar), 50th (marker) and 75th (upper bar) percentiles. (a,c) short dataset, (b,d) long dataset.

5.1.2. Model-informed ESN

We leverage knowledge about the governing equations through $\mathcal{K}(\mathbf{u}_{\text{in}})$ in the model in Eq. (6). In this testcase, we use a reduced-order model obtained through Proper Orthogonal Decomposition (POD) [65, 66] to define a POD-informed ESN. POD provides a fixed rank subspace \mathcal{E} of the state space, in which the projection of the original state vector optimally preserves its energy. The POD modes / energies are the eigenvectors / eigenvalues, of the data covariance matrix $\mathbf{C} = \frac{1}{m-1} \mathbf{U}^T \mathbf{U}$. The $M \times N_u$ matrix \mathbf{U} is the vertical concatenation of the M snapshots of the N_u -dimensional timeseries used for washout, training and validation of the network, from which its mean, $\mathbf{d} \in \mathbb{R}^{N_u}$, is subtracted columns-wise. We create an N_{POD} -dimensional reduced-order model by taking the modes ϕ_i associated with the N_{POD} largest eigenvalues of \mathbf{C} . Because \mathbf{C} is a symmetric matrix, its eigenvectors form an orthonormal basis, which is stored in the orthogonal matrix $\Phi = [\phi_1; \dots; \phi_{N_{\text{POD}}}]$. The state vector \mathbf{q} is expressed as a function of its components ξ in the subspace \mathcal{E} spanned by Φ , and its components η in the orthogonal complement of \mathcal{E} spanned by the basis Ψ

$$\mathbf{q} = \Phi \xi + \Psi \eta + \mathbf{d}. \quad (12)$$

The evolution equations are then obtained by using a *flat Galerkin approximation* [67], which neglects the contribution of the orthogonal complement: $\Psi \dot{\eta} \simeq 0$. The nonlinear dynamical system $\dot{\mathbf{q}} = \mathbf{f}(\mathbf{q})$ is projected onto \mathcal{E} through $\dot{\xi} = \Phi^T (\mathbf{f}(\Phi \xi + \mathbf{d}) - \mathbf{d})$ as

$$\dot{\xi} = \Phi^T \mathbf{f}(\Phi \xi + \mathbf{d}) \quad (13)$$

In the POD-informed ESN model ($\mathbf{q} \equiv \mathbf{u}_{\text{in}}$), we use $N_{\text{POD}} = 2$ to generate the reduced-order model, which accounts for 96% of the energy of the original signal. We use the evolution of the trajectory on the POD subspace, \mathcal{E} , to inform the ESN through $\mathcal{K}(\mathbf{u}_{\text{in}}(t_i)) = \xi(t_{i+1})$. We solve the ODE system in Eq. (13) using at each time step forward Euler with initial condition $\xi(t_i) = \Phi^T (\mathbf{u}_{\text{in}}(t_i) - \mathbf{d})$, which is the projection of the input to the network onto \mathcal{E} . The projection of the trajectory and the autonomous evolution of the flat Galerkin approximation are shown in Fig. 10. The reduced order model dynamics, ξ , differ significantly

from the dynamics of the entire state, $\Phi^T(\mathbf{q}-\mathbf{b})$. However, we show in the next paragraph that embedding model knowledge, yet imperfect, can improve the performance of the networks.

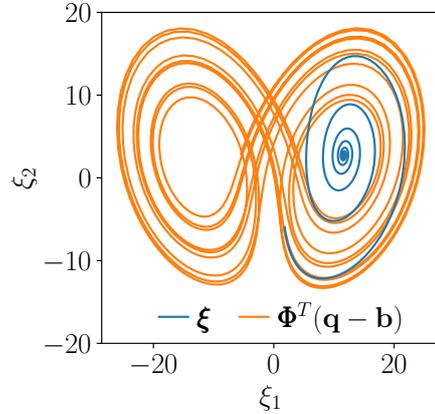


Figure 10: Projection of the trajectory onto the 2-dimensional model ($\Phi^T(\mathbf{q}-\mathbf{b})$) and autonomous evolution of the flat Galerkin approximation (ξ) in the POD-informed Echo State Network.

As compared to the model-free ESN, there is a decrease in correlation between the validation and test sets for almost all the validation strategies (Tab. 2). The Single Shot Validation is still outperformed by the other strategies, while the chaotic K-Fold Validation and chaotic Recycle Validation have the highest correlation. Panels (a,b) of Fig. 11 show similar results to the model-free case: Bayesian Optimization outperforms Grid Search in the validation set, but the two schemes perform similarly in the test set. The only exception are the chaotic K-Fold Validation and chaotic Recycle Validation in the long dataset (b), in which Bayesian Optimization outperforms Grid Search for up to 75% of the networks in the test set. Panels (c,d) show that embedding knowledge of the governing equation produces an increase of 1LT in the Prediction Horizon with respect to the model-free case (see Fig. 9). The qualitative behaviour of the validation strategies remains similar to the model-free case. To conclude, although the POD-informed architecture does increase the performance, it does not increase the robustness of the networks with respect to the model-free ESN.

Table 2: Spearman coefficients between validation and test sets. Bold text indicates the highest correlation in the dataset.

\tilde{r}_S	SSV	WFV	WFV _c	WFV _c *	KFV	KFV _c	RV	RV _c
Short dataset (12LTs)	0.15	0.39	0.34	-	0.32	0.73	0.41	0.56
Long dataset (24LTs)	0.19	0.42	0.36	0.51	0.59	0.80	0.55	0.80

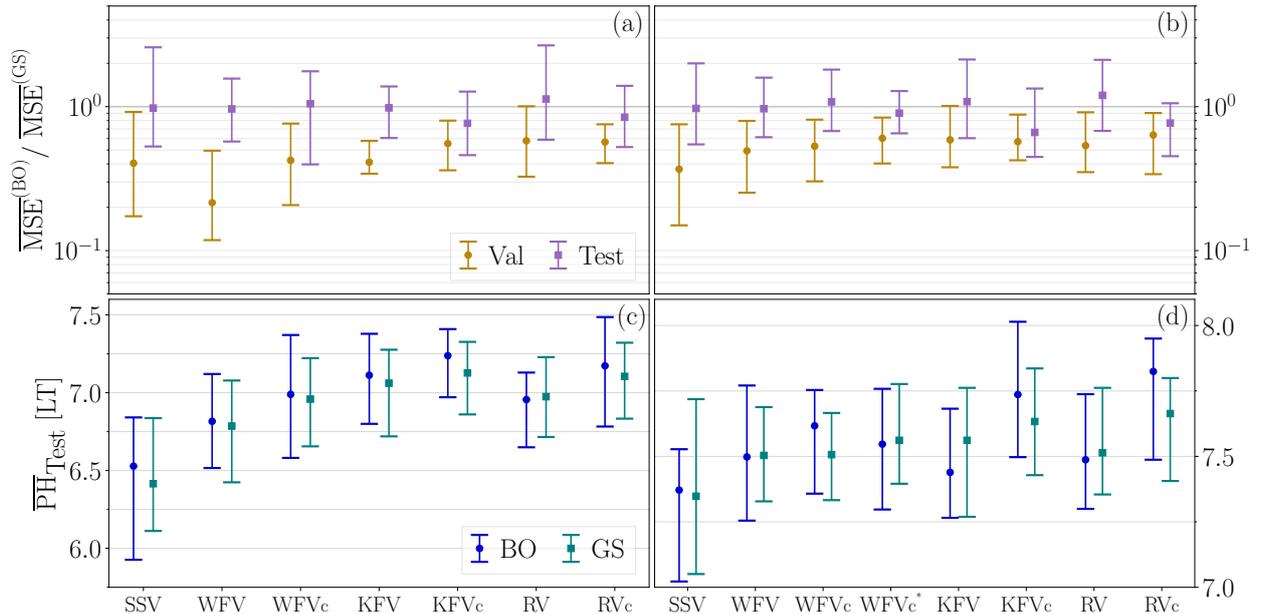


Figure 11: Same as Fig. 9 for the POD-informed ESN.

6. Validation for quasiperiodic solutions

We analyse the nonlinear oscillator proposed by Kuznetsov et al. [55], which physically represents a self-oscillatory discharge in an electric circuit. The oscillator is a three-dimensional system, which can display periodic, quasiperiodic and chaotic behaviours as a function of the parameters $[\lambda, \omega_0, \mu]$

$$\begin{aligned}
 \dot{x} &= y, \\
 \dot{y} &= y(\lambda + z + x^2 - \frac{1}{2}x^4) - \omega_0^2 x, \\
 \dot{z} &= \mu - x^2.
 \end{aligned} \tag{14}$$

The primary purpose of this testcase is to compare the robustness of Echo State Networks in forecasting quasiperiodic solutions versus chaotic solutions. This enables us to determine whether the challenges encountered in the Lorenz system are specific to learning chaotic time series. We obtain quasiperiodic and chaotic solutions by setting $\lambda = 0$, $\omega_0 = 2.7$ and $\mu_{Qp} = 0.9$ and $\mu_{Ch} = 0.5$, as shown in Figs. 12(b,d), respectively. (For completeness, in this section, we report the Kuznetsov chaotic solution as well.) The datasets of 7.5 LTs that we use for washout, training and validation are shown in panels (a,c).

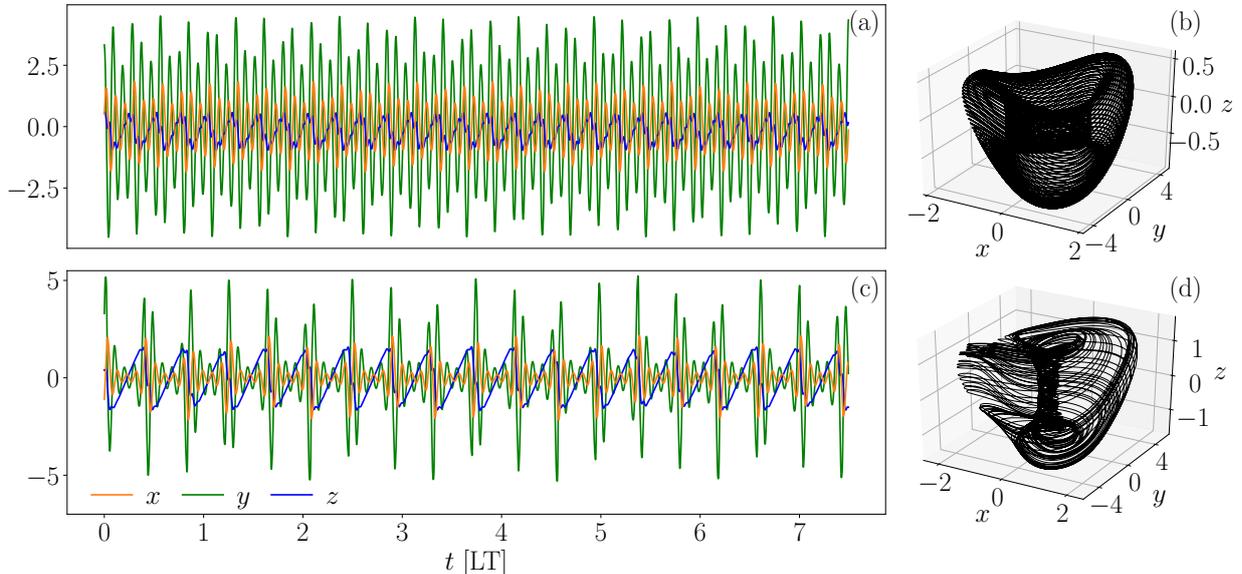


Figure 12: Kuznetsov oscillator. (a) Quasiperiodic and (c) chaotic time series; (b) quasiperiodic and (d) chaotic phase plots for a longer time window. Time is expressed in the Lyapunov time (LT) of the chaotic case ($LT \approx 25$ [55]).

6.1. Hyperparameter optimization

The network parameters, the size of the ensemble, and the optimization strategies are the same as those of section 4. We modify the input scaling, $b_{in} = 0.1$, for it to have the same order of magnitude of the input, which is obtained by normalizing the signal by its maximum variation component-wise. We study the enlarged interval $\rho = [0.01, 1]$ because we observed empirically that the optimal hyperparameters often lie in the range $\rho \leq 0.1$. Given the multiple orders of magnitude of the spectral radius, the hyperparameter space is analysed in a logarithmic scale. We use the same architecture and validation strategies for the quasiperiodic and chaotic case (as detailed in the Supplementary Material, S.1). The different strategies are tested by computing the arithmetic mean \overline{PH}_{test} of the Prediction Horizon on N_t starting points for the chaotic case, and by computing the geometric mean \overline{MSE}_{test} of the Mean Squared Error in 2 LTs intervals starting from the same points. In the chaotic case, we select $N_t = 75$, whereas in the quasiperiodic we select $N_t = 50$ through the procedure described in Appendix D. The performance in the quasiperiodic dataset is assessed only through the Mean Squared Error because the Prediction Horizon is infinite, i.e., a quasiperiodic solution has zero dominant Lyapunov exponents [68].

6.1.1. Model-free ESN

Figure 13 shows the MSE in the hyperparameter space for the quasiperiodic case. The plots for three validation strategies, (a-c), are very close to the MSE in the test set, (d), which means that hyperparameters that perform well in the validation set, perform as well in the test set. This is in contrast with the behaviour in chaotic solutions (see Fig. 7). The Spearman coefficients (Tab. 3) confirm that the correlation between validation and test sets is higher in the quasiperiodic dataset than the chaotic dataset. Notably, the peak $\tilde{r}_s = 0.97$ obtained in the Recycle Validations indicates almost complete correlation. As before, the Single Shot Validation is outperformed by the K-Fold Validation and Recycle Validation, but its correlation in the quasiperiodic dataset is higher than that of chaotic cases. The high correlation in the quasiperiodic dataset is identified as the dense clustering around the linear regression of Fig. 14. Two remarks can be made. On the one hand, the high correlation in the quasiperiodic dataset implies that the challenges in producing robust results in Echo State Networks in chaotic attractors are due to the complexity of the chaotic signal, rather than the properties of the networks. On the other hand, the marked difference in performance between

different networks is still present in the quasiperiodic dataset, which means that ESNs are sensitive to the realizations (further analysis is reported in Appendix E). Practically, we advise that different networks be optimized independently in the quasiperiodic case as well.

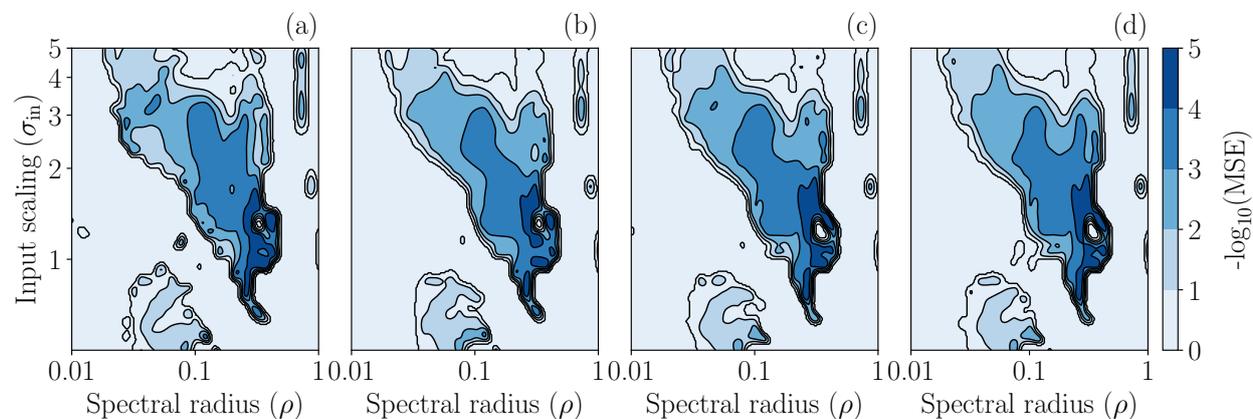


Figure 13: Mean of the Gaussian Process reconstruction for the quasiperiodic dataset for a representative network of the ensemble. Validation set for (a) Single Shot Validation, (b) chaotic K-Fold Validation, and (c) chaotic Recycle Validation; and (d) test set. The MSE is saturated to be ≤ 1 . The Gaussian Process is based on a grid of 30×30 data points.

Table 3: Spearman coefficients between validation and test sets. Bold text indicates the highest correlation in the dataset.

\tilde{r}_S	SSV	WFV	WFV _c	KFV	KFV _c	RV	RV _c
Quasiperiodic dataset	0.80	0.75	0.71	0.93	0.92	0.97	0.97
Chaotic dataset	0.49	0.48	0.58	0.70	0.76	0.66	0.81

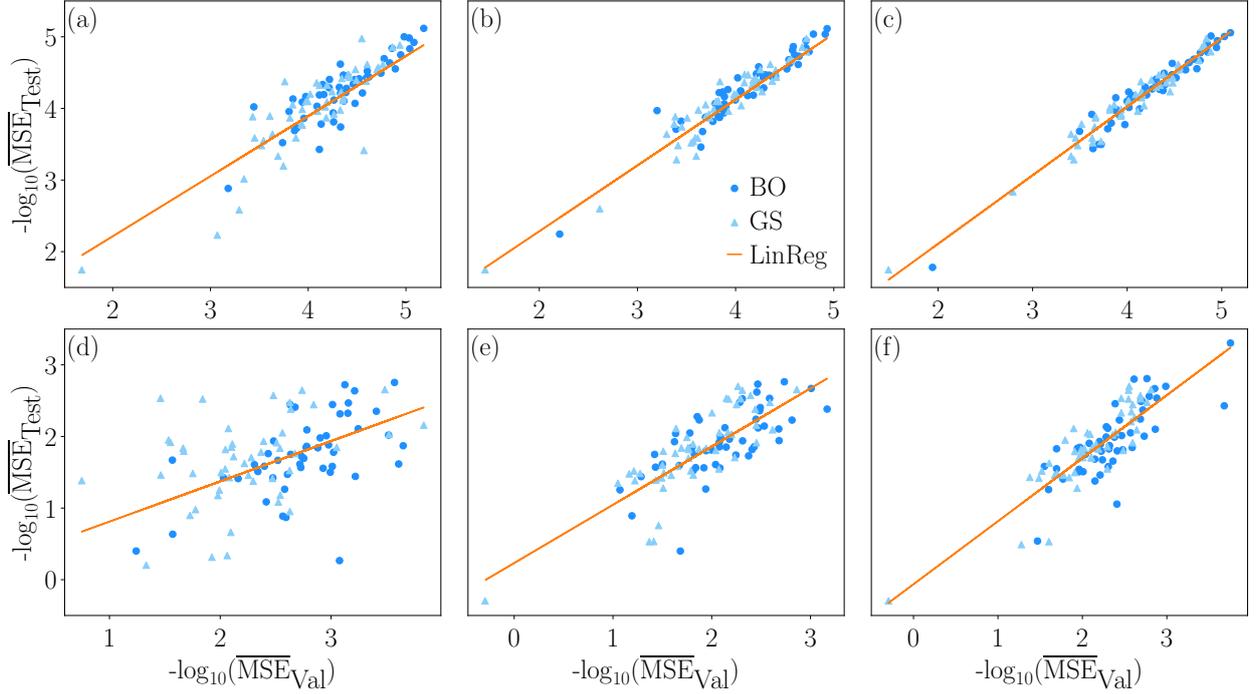


Figure 14: Linear regression (LinReg) and scatter plot of the MSE of the optimal hyperparameters obtained from Bayesian Optimization (BO) and Grid Search (GS) for each network. (a,d) Single Shot Validation, (b,e) chaotic K-Fold Validation and (c,f) chaotic Recycle Validation. (a-c) quasiperiodic and (d-f) chaotic datasets.

Panels (a,b) of Fig. 15 show the ratio of the MSE between the optimal hyperparameters obtained by Bayesian Optimization (BO) and the optimal hyperparameters from Grid Search (GS) in the validation and test sets. On the one hand, in the quasiperiodic case (a) the performance in the validation set is similar to the test set. On the other hand, in the chaotic case (b) BO outperforms GS in the validation set, although the two schemes perform similarly in the test set. In panels (c,d), we show the performance of the networks in the test set using the MSE for the quasiperiodic dataset (c) and the Prediction Horizon in the chaotic dataset (d). In the quasiperiodic dataset, Bayesian Optimization outperforms Grid Search, and the K-Fold Validations and Recycle Validations outperform the other validation strategies. In the chaotic dataset, as seen in the Lorenz system, Bayesian Optimization only slightly outperforms Grid Search, while the K-fold Validations and Recycle Validations still outperform the other validation strategies.

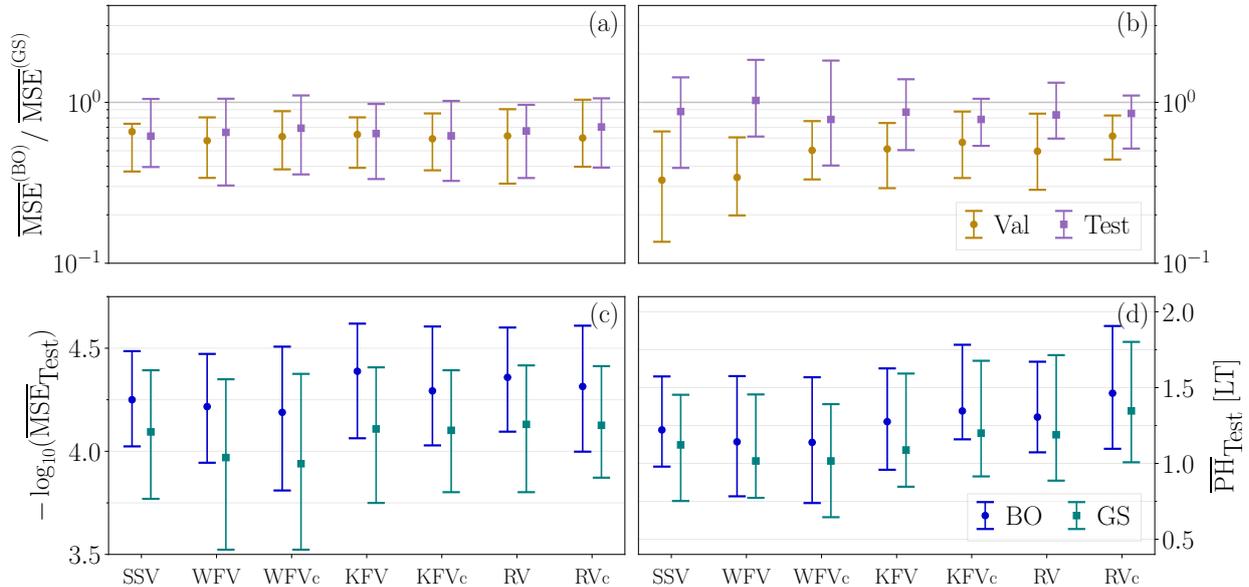


Figure 15: Comparison between hyperparameter optimization by Bayesian Optimization (BO) and Grid Search (GS) for the two performance metrics (MSE, PH). 25th (lower bar), 50th (marker) and 75th (upper bar) percentiles. (a,c) quasiperiodic dataset, (b,d.) chaotic dataset.

6.1.2. Model-informed ESN

We design a Forward Euler (FE) informed ESN (6) by integrating in time with forward Euler the y equation (14) only

$$\mathcal{K}(\mathbf{u}_{\text{in}}) = y + dt(y(\lambda + z + x^2 - \frac{1}{2}x^4) - \omega_0^2 x). \quad (15)$$

Tab. 4 shows the Spearman coefficients for the FE-informed model. In the quasiperiodic dataset, the correlation decreases for all the validation strategies with respect to the model-free case (see Tab. 3) except for the Recycle Validations, which have the highest correlation. However, in the chaotic dataset, the correlation increases for all the validation strategies. Here, the chaotic K-Fold Validation and chaotic Recycle Validation are the strategies with the highest correlation. Fig.16(a) shows that the decrease in correlation in the quasiperiodic dataset causes Bayesian Optimization to generate larger MSE than Grid Search with respect to the model-free case. In panel (b), we observe that there is still a marked discrepancy between the performance of the optimization schemes in the validation and test sets. Panels (c,d) show the performance of the FE-informed ESN in the test set. In both datasets, the performance improves when leveraging knowledge about the governing equations: the MSE decreases by about two orders of magnitude, and the Prediction Horizon improves by $\gtrsim 2$ Lyapunov Times with respect to the model-free case (see Fig. 15). The improvement in performance, however, does not correspond to a consistent increase in correlation between validation and test sets. The performance of Bayesian Optimization with respect to Grid Search in the test set does not necessarily improve. In the same fashion as the Lorenz system, the FE-informed architecture *per se* does enhance the performance, but it does not enhance the robustness of Echo State Networks.

Table 4: Spearman coefficients between validation and test sets. Bold text indicates the highest correlation in the dataset.

\tilde{r}_S	SSV	WFV	WFV _c	KfV	KfV _c	RV	RV _c
Quasiperiodic dataset	0.78	0.65	0.67	0.71	0.80	0.98	0.98
Chaotic dataset	0.57	0.63	0.63	0.75	0.79	0.71	0.85

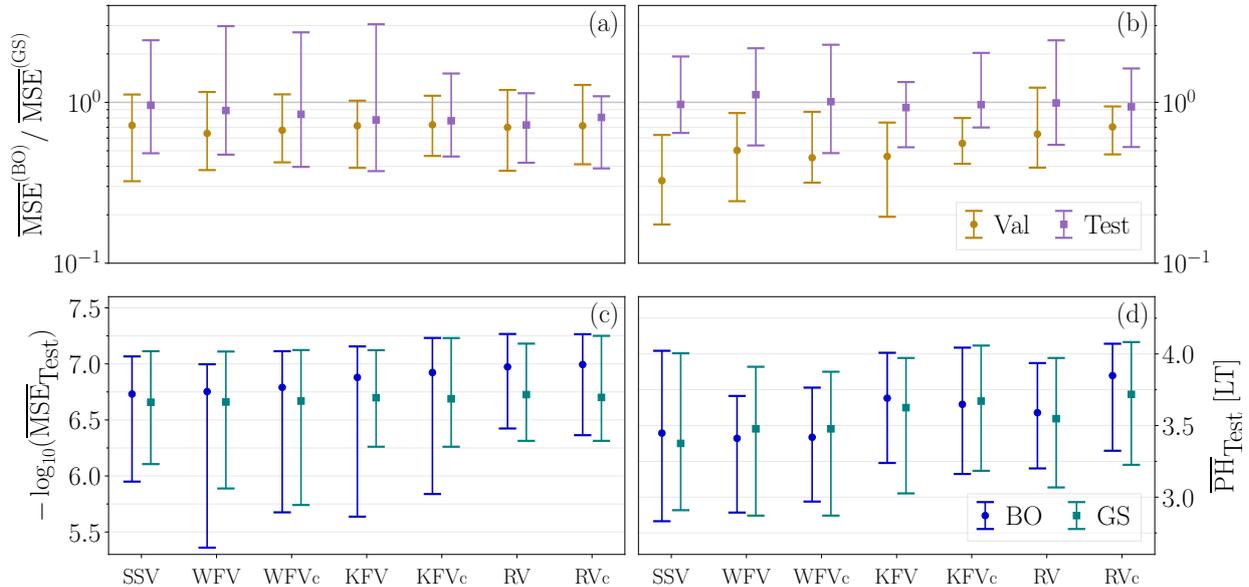


Figure 16: Same as Fig. 15 for the FE-informed ESN.

7. Conclusions

The Echo State Network (ESN) is a reservoir computing architecture that is able to learn accurately the nonlinear dynamics of systems from data. The overarching objective of this paper is to investigate and improve the robustness of ESNs, with a focus on the forecasting of chaotic systems. First, we analyse the Single Shot Validation, which is the commonly used strategy to select the hyperparameters. We show that the Single Shot Validation is the least performing strategy to fine-tune the hyperparameters. Second, we validate the ESNs on multiple points of the chaotic attractor, for which the validation set is not necessarily subsequent in time to the training set. We propose the Recycle Validation and the chaotic version of existing validation strategies based on multiple folds, such as the Walk Forward Validation and the K-Fold Cross Validation. The K-Fold Validation and Recycle Validation offer the greatest robustness and performance, with their chaotic versions outperforming the corresponding regular versions. Importantly, the Recycle Validation is computationally cheaper than the K-Fold Cross Validation. Third, we compare Bayesian Optimization with Grid Search to compute the optimal hyperparameters. We find that Bayesian Optimization is an optimization scheme that consistently finds a set of hyperparameters that perform significantly better than the Grid Search in the validation set. On the one hand, in learning quasiperiodic solutions, hyperparameters that work optimally in the validation set continue to work optimally in the test set. This is because quasiperiodic solutions are predictable (i.e., they do not have positive Lyapunov exponents). This finding is, thus, expected to generalize to other predictable solutions, such as frequency-locked solutions and limit cycles. On the other hand, in learning chaotic solutions, hyperparameters that work optimally in the validation set do not necessarily work optimally in the test set. We argue that this occurs because of the chaotic nature of the attractor, in which the nonlinear dynamics, although deterministic, manifest themselves as unpredictable variations. Fourth, we analyse the model-free ESN, which is fully data-driven, and the model-informed ESN, which leverages knowledge of the governing equations. We find that the model-informed architecture markedly improves the network’s prediction capabilities, but it does not improve the robustness. Finally, we find that the optimal hyperparameters are significantly sensitive to the random initialization of the ESN. Practically, when working with an ensemble of ESNs, we recommend computing the optimal hyperparameters for each network. In the test performed in the paper, this can increase up to six Lyapunov Times the network’s Prediction Horizon as compared to using the same set of

hyperparameters for all realizations.

This work opens up new possibilities for using Echo State Networks and, in general, recurrent neural networks, for robust learning of chaotic dynamics.

Acknowledgements

A. Racca is supported by the EPSRC-DTP and the Cambridge Commonwealth, European & International Trust under a Cambridge European Scholarship. L. Magri is supported by the Royal Academy of Engineering Research Fellowship scheme and the visiting fellowship at the Technical University of Munich – Institute for Advanced Study, funded by the German Excellence Initiative and the European Union Seventh Framework Programme under grant agreement n. 291763. The authors would like to thank Dr. N. A. K. Doan and F. Huhn for insightful discussions.

References

- [1] R. G. Deissler, Is Navier-Stokes turbulence chaotic?, *Phys. Fluids* 29 (5 , May 1986) (1986) 1453–1457. doi:10.1063/1.865663. URL <https://aip.scitation.org/doi/10.1063/1.865663>
- [2] G. Boffetta, M. Cencini, M. Falcioni, A. Vulpiani, Predictability: A way to characterize complexity, *Physics Reports* 356 (6) (2002) 367–474. arXiv:0101029, doi:10.1016/S0370-1573(01)00025-4.
- [3] J. Bec, L. Biferale, G. Boffetta, M. Cencini, S. Musacchio, F. Toschi, Lyapunov exponents of heavy particles in turbulence, *Physics of Fluids* 18 (9) (2006) 1–5. arXiv:0606024, doi:10.1063/1.2349587.
- [4] F. C. Moon, S. W. Shaw, Chaotic vibrations of a beam with non-linear boundary conditions, *International Journal of non-linear Mechanics* 18 (6) (1983) 465–477.
- [5] M. Kennedy, R. Rovatti, G. Setti, *Chaotic electronics in telecommunications*, CRC press, 2000.
- [6] H.-J. Stöckmann, *Quantum chaos: an introduction* (2000).
- [7] G. Nastac, J. W. Labahn, L. Magri, M. Ihme, Lyapunov exponent as a metric for assessing the dynamic content and predictability of large-eddy simulations, *Physical Review Fluids* 2 (9) (2017) 094606. doi:10.1103/PhysRevFluids.2.094606.
- [8] M. Hassanaly, V. Raman, Ensemble-LES analysis of perturbation response of turbulent partially-premixed flames, *Proceedings of the Combustion Institute* 37 (2) (2019) 2249–2257. doi:10.1016/j.proci.2018.06.209. URL <https://doi.org/10.1016/j.proci.2018.06.209>
- [9] B. M. Bolker, B. T. Grenfell, Chaos and biological complexity in measles dynamics, *Proceedings of the Royal Society of London. Series B: Biological Sciences* 251 (1330) (1993) 75–81.
- [10] E. N. Lorenz, Deterministic nonperiodic flow, *J. Atmos. Sci.* 20 (2) (1963) 130–141.
- [11] F. Takens, Detecting strange attractors in turbulence, in: *Dynamical systems and turbulence*, Warwick 1980, Springer, 1981, pp. 366–381.
- [12] J. Guckenheimer, P. Holmes, *Nonlinear oscillations, dynamical systems, and bifurcations of vector fields*, Vol. 42, Springer Science & Business Media, 2013.
- [13] N. Baker, F. Alexander, T. Bremer, A. Hagberg, Y. Kevrekidis, H. Najm, M. Parashar, A. Patra, J. Sethian, S. Wild, et al., Workshop report on basic research needs for scientific machine learning: Core technologies for artificial intelligence, Tech. rep., USDOE Office of Science (SC), Washington, DC (United States) (2019).
- [14] I. Goodfellow, Y. Bengio, A. Courville, *Deep learning*, MIT press, 2016.
- [15] D. E. Rumelhart, G. E. Hinton, R. J. Williams, Learning representations by back-propagating errors, *nature* 323 (6088) (1986) 533–536.
- [16] H. Sak, A. W. Senior, F. Beaufays, Long short-term memory recurrent neural network architectures for large scale acoustic modeling (2014).
- [17] I. Sutskever, O. Vinyals, Q. V. Le, Sequence to sequence learning with neural networks, arXiv preprint arXiv:1409.3215 (2014).
- [18] S. L. Brunton, B. R. Noack, P. Koumoutsakos, Machine learning for fluid mechanics, *Annual Review of Fluid Mechanics* 52 (2020) 477–508.
- [19] K. Nakai, Y. Saiki, Machine-learning inference of fluid variables from data using reservoir computing, *Physical Review E* 98 (2) (2018) 023111.
- [20] Z. Y. Wan, T. P. Sapsis, Machine learning the kinematics of spherical particles in fluid flows, *Journal of Fluid Mechanics* 857 (2018).
- [21] N. A. K. Doan, W. Polifke, L. Magri, A physics-aware machine to predict extreme events in turbulence, arXiv preprint arXiv:1912.10994 (2019).
- [22] P. R. Vlachas, W. Byeon, Z. Y. Wan, T. P. Sapsis, P. Koumoutsakos, 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 (2213) (2018) 20170844.

- [23] F. Huhn, L. Magri, Learning ergodic averages in chaotic systems, arXiv preprint arXiv:2001.04027 (2020).
- [24] P. J. Werbos, Backpropagation through time: what it does and how to do it, *Proceedings of the IEEE* 78 (10) (1990) 1550–1560.
- [25] Y. Bengio, P. Simard, P. Frasconi, Learning long-term dependencies with gradient descent is difficult, *IEEE transactions on neural networks* 5 (2) (1994) 157–166.
- [26] S. Hochreiter, J. Schmidhuber, Long short-term memory, *Neural computation* 9 (8) (1997) 1735–1780.
- [27] K. Cho, B. Van Merriënboer, C. Gulcehre, D. Bahdanau, F. Bougares, H. Schwenk, Y. Bengio, Learning phrase representations using rnn encoder-decoder for statistical machine translation, arXiv preprint arXiv:1406.1078 (2014).
- [28] H. Jaeger, H. Haas, Harnessing nonlinearity: Predicting chaotic systems and saving energy in wireless communication, *Science* 304 (5667) (2004) 78–80.
- [29] W. Maass, T. Natschläger, H. Markram, Real-time computing without stable states: A new framework for neural computation based on perturbations, *Neural computation* 14 (11) (2002) 2531–2560.
- [30] M. Lukoševičius, A practical guide to applying echo state networks, in: *Neural networks: Tricks of the trade*, Springer, 2012, pp. 659–686.
- [31] Z. Lu, B. R. Hunt, E. Ott, Attractor reconstruction by machine learning, *Chaos: An Interdisciplinary Journal of Nonlinear Science* 28 (6) (2018) 061104.
- [32] J. Pathak, Z. Lu, B. R. Hunt, M. Girvan, E. Ott, Using machine learning to replicate chaotic attractors and calculate lyapunov exponents from data, *Chaos: An Interdisciplinary Journal of Nonlinear Science* 27 (12) (2017) 121102.
- [33] N. Doan, W. Polifke, L. Magri, Physics-informed echo state networks, *Journal of Computational Science* 47 (2020) 101237. doi:<https://doi.org/10.1016/j.jocs.2020.101237>. URL <http://www.sciencedirect.com/science/article/pii/S1877750320305408>
- [34] Z. Lu, J. Pathak, B. Hunt, M. Girvan, R. Brockett, E. Ott, Reservoir observers: Model-free inference of unmeasured variables in chaotic systems, *Chaos: An Interdisciplinary Journal of Nonlinear Science* 27 (4) (2017) 041102.
- [35] N. A. K. Doan, W. Polifke, L. Magri, Learning hidden states in a chaotic system: A physics-informed echo state network approach, in: *ICCS*, Springer, 2020, pp. 117–123.
- [36] A. Racca, L. Magri, Automatic-differentiated physics-informed echo state network (api-esn), arXiv preprint arXiv:2101.00002 (2020).
- [37] J. Pathak, B. Hunt, M. Girvan, Z. Lu, E. Ott, Model-free prediction of large spatiotemporally chaotic systems from data: A reservoir computing approach, *Physical review letters* 120 (2) (2018) 024102.
- [38] N. A. K. Doan, W. Polifke, L. Magri, Physics-informed echo state networks for chaotic systems forecasting, in: *ICCS*, Springer, 2019, pp. 192–198.
- [39] A. Wikner, J. Pathak, B. Hunt, M. Girvan, T. Arcomano, I. Szunyogh, A. Pomerance, E. Ott, Combining machine learning with knowledge-based modeling for scalable forecasting and subgrid-scale closure of large, complex, spatiotemporal systems, *Chaos: An Interdisciplinary Journal of Nonlinear Science* 30 (5) (2020) 053111.
- [40] L. Grigoryeva, J.-P. Ortega, Echo state networks are universal, *Neural Networks* 108 (2018) 495–508. doi:<https://doi.org/10.1016/j.neunet.2018.08.025>. URL <https://www.sciencedirect.com/science/article/pii/S089360801830251X>
- [41] L. Gonon, J.-P. Ortega, Fading memory echo state networks are universal, *Neural Networks* (2021). doi:<https://doi.org/10.1016/j.neunet.2021.01.025>. URL <https://www.sciencedirect.com/science/article/pii/S0893608021000332>
- [42] J. Pathak, A. Wikner, R. Fussell, S. Chandra, B. R. Hunt, M. Girvan, E. Ott, Hybrid forecasting of chaotic processes: Using machine learning in conjunction with a knowledge-based model, *Chaos: An Interdisciplinary Journal of Nonlinear Science* 28 (4) (2018) 041101.
- [43] P. R. Vlachas, J. Pathak, B. R. Hunt, T. P. Sapsis, M. Girvan, E. Ott, P. Koumoutsakos, Backpropagation algorithms and reservoir computing in recurrent neural networks for the forecasting of complex spatiotemporal dynamics, *Neural Networks* (2020).
- [44] A. Chattopadhyay, P. Hassanzadeh, K. Palem, D. Subramanian, Data-driven prediction of a multi-scale lorenz 96 chaotic system using a hierarchy of deep learning methods: Reservoir computing, ann, and rnn-lstm, arXiv preprint arXiv:1906.08829 (2019).
- [45] A. Haluszczynski, C. R ath, Good and bad predictions: Assessing and improving the replication of chaotic attractors by means of reservoir computing, *Chaos: An Interdisciplinary Journal of Nonlinear Science* 29 (10) (2019) 103143.
- [46] J. Jiang, Y.-C. Lai, Model-free prediction of spatiotemporal dynamical systems with recurrent neural networks: Role of network spectral radius, *Physical Review Research* 1 (3) (2019) 033056.
- [47] M. Lukoševičius, A. Uselis, Efficient cross-validation of echo state networks, in: *International Conference on Artificial Neural Networks*, Springer, 2019, pp. 121–133.
- [48] K. Ishu, T. van der Zant, V. Becanovic, P. Ploger, Identification of motion with echo state network, in: *Oceans '04 MTS/IEEE Techno-Ocean '04* (IEEE Cat. No.04CH37600), Vol. 3, 2004, pp. 1205–1210 Vol.3.
- [49] A. A. Ferreira, T. B. Ludermitz, R. R. De Aquino, An approach to reservoir computing design and training, *Expert systems with applications* 40 (10) (2013) 4172–4182.
- [50] L. A. Thiede, U. Parlitz, Gradient based hyperparameter optimization in echo state networks, *Neural Networks* 115 (2019) 23–29.
- [51] H. Wang, X. Yan, Optimizing the echo state network with a binary particle swarm optimization algorithm, *Knowledge-Based Systems* 86 (2015) 182–193.
- [52] J. Yperman, T. Becker, Bayesian optimization of hyper-parameters in reservoir computing, arXiv preprint arXiv:1611.05193 (2016).

- [53] A. Griffith, A. Pomerance, D. J. Gauthier, Forecasting chaotic systems with very low connectivity reservoir computers, *Chaos: An Interdisciplinary Journal of Nonlinear Science* 29 (12) (2019) 123108.
- [54] C. E. Rasmussen, Gaussian processes in machine learning, in: *Summer School on Machine Learning*, Springer, 2003, pp. 63–71.
- [55] A. Kuznetsov, S. Kuznetsov, N. Stankevich, A simple autonomous quasiperiodic self-oscillator, *Communications in Nonlinear Science and Numerical Simulation* 15 (6) (2010) 1676–1681.
- [56] I. B. Yildiz, H. Jaeger, S. J. Kiebel, Re-visiting the echo state property, *Neural Networks* 35 (2012) 1–9. doi:<https://doi.org/10.1016/j.neunet.2012.07.005>. URL <https://www.sciencedirect.com/science/article/pii/S0893608012001852>
- [57] A. N. Tikhonov, A. Goncharsky, V. Stepanov, A. G. Yagola, *Numerical methods for the solution of ill-posed problems*, Vol. 328, Springer Science & Business Media, 2013.
- [58] C. R. Harris, K. J. Millman, S. J. van der Walt, R. Gommers, P. Virtanen, D. Cournapeau, E. Wieser, J. Taylor, S. Berg, N. J. Smith, R. Kern, M. Picus, S. Hoyer, M. H. van Kerkwijk, M. Brett, A. Haldane, J. F. del R'io, M. Wiebe, P. Peterson, P. G'erard-Marchant, K. Sheppard, T. Reddy, W. Weckesser, H. Abbasi, C. Gohlke, T. E. Oliphant, *Array programming with NumPy*, *Nature* 585 (7825) (2020) 357–362. doi:[10.1038/s41586-020-2649-2](https://doi.org/10.1038/s41586-020-2649-2). URL <https://doi.org/10.1038/s41586-020-2649-2>
- [59] E. Brochu, V. M. Cora, N. De Freitas, A tutorial on bayesian optimization of expensive cost functions, with application to active user modeling and hierarchical reinforcement learning, *arXiv preprint arXiv:1012.2599* (2010).
- [60] J. Snoek, H. Larochelle, R. P. Adams, Practical bayesian optimization of machine learning algorithms, in: *Advances in neural information processing systems*, 2012, pp. 2951–2959.
- [61] M. D. Hoffman, E. Brochu, N. de Freitas, Portfolio allocation for bayesian optimization., in: *UAI*, Citeseer, 2011, pp. 327–336.
- [62] P. Virtanen, et al., *SciPy 1.0: Fundamental Algorithms for Scientific Computing in Python*, *Nature Methods* 17 (2020) 261–272.
- [63] D. Viswanath, Lyapunov exponents from random fibonacci sequences to the lorenzequations, *Tech. rep.*, Cornell University (1998).
- [64] C. Spearman, The proof and measurement of association between two things, *The American Journal of Psychology* 15 (1) (1904) 72–101. URL <http://www.jstor.org/stable/1412159>
- [65] J. L. Lumley, *The structure of inhomogeneous turbulent flows*, *Atmospheric turbulence and radio wave propagation* (1967).
- [66] J. Weiss, A tutorial on the proper orthogonal decomposition, in: *AIAA Aviation 2019 Forum*, 2019, p. 3333.
- [67] H. G. Matthies, M. Meyer, Nonlinear galerkin methods for the model reduction of nonlinear dynamical systems, *Computers & Structures* 81 (12) (2003) 1277–1286.
- [68] H. Kantz, T. Schreiber, *Nonlinear time series analysis*, Vol. 7, Cambridge university press, 2004.

Appendix A. Correlation between the mean-squared error and predictability horizon

We show the high correlation between the Mean Squared Error and the Predictability Horizon given the same starting point for prediction. Figure A.17 shows the Gaussian Process reconstruction from 900 (30×30) grid points in the hyperparameter space in the $N_t = 100$ test set (Appendix D) of the Lorenz system for a representative network realization. The two quantities show almost identical behaviour. Figure A.18 shows the scatter plots for the Prediction Horizon and the Mean Square Error in the test set for the optimal hyperparameters for the ensemble in the Lorenz system. The two quantities are highly correlated, with a Spearman coefficient, $r_s \geq 0.95$, for all the validation strategies.

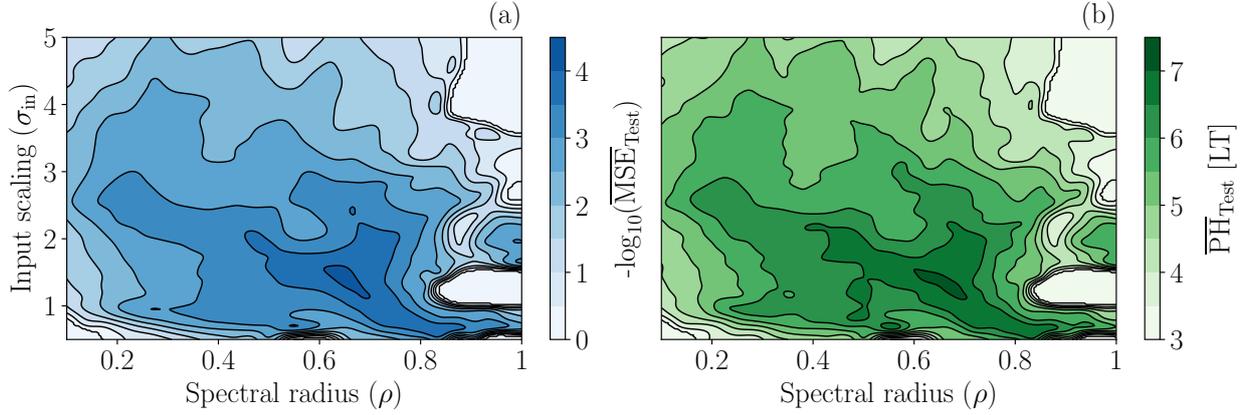


Figure A.17: Mean of the Gaussian Process reconstruction from a 900 points grid in the test set for (a) $\log_{10}(\text{MSE})$ and (b) Prediction Horizon (PH) for a representative network in the short dataset. For visualization purposes we saturate the MSE to be ≤ 1 , and the PH to be ≥ 3 . The MSE and PH closely resemble one another.

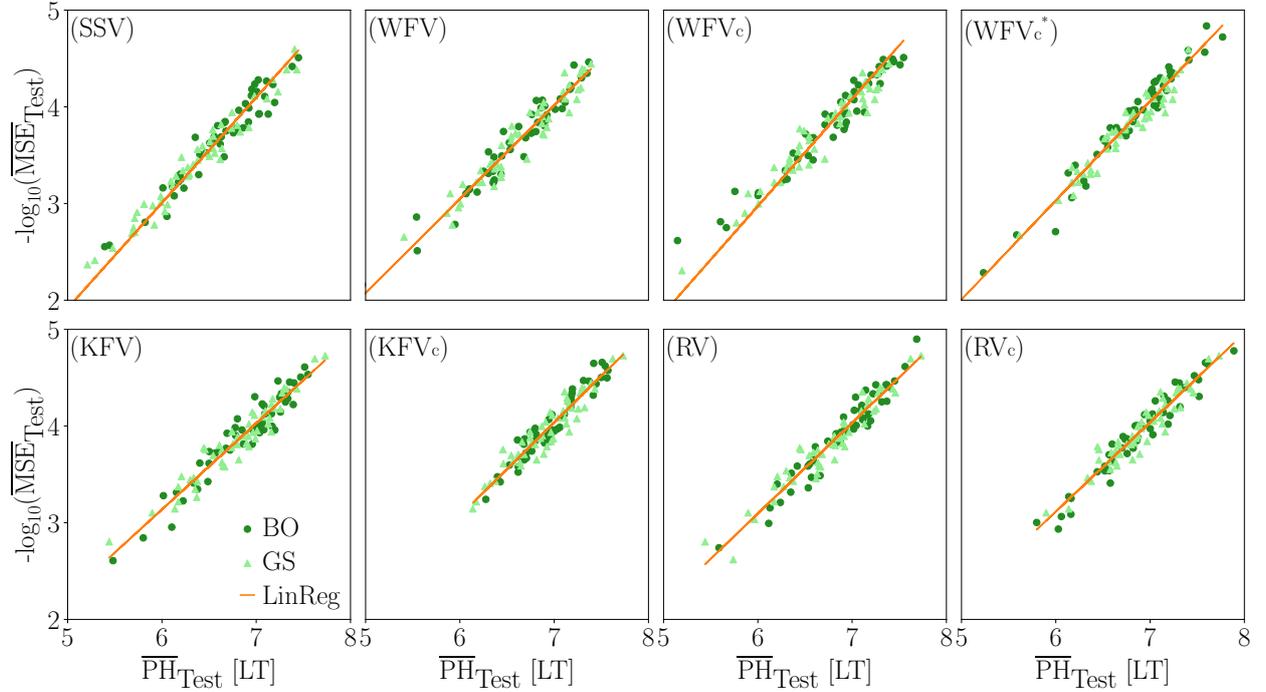


Figure A.18: Linear regression (LinReg) and scatter plot for the test set MSE and Prediction Horizon in the long dataset of the optimal hyperparameters from Bayesian Optimization (BO) and Grid Search (GS) for the Single Shot Validation (SSV), Walk Forward Validation (WFV), K-Fold Validation (KFV), Recycle Validation (RV), and their chaotic versions, with subscript c . The trends are highly correlated.

Appendix B. Computational time

In Fig. B.19, we show the CPU time required by the validation strategies to perform a Grid Search in hyperparameters space for a single network. The computational advantage of the Recycle Validation increases with the size of the dataset and the size of the reservoir. We expect the improvement in computational time to be more significant in RNN architectures whose training is more expensive, such as LSTMs and GRUs.

The Bayesian Optimization described in section 4 costs approximately 6 seconds more per network in all the cases shown. This because the additional cost of the Bayesian Optimization is independent of the cost of the evaluation function.

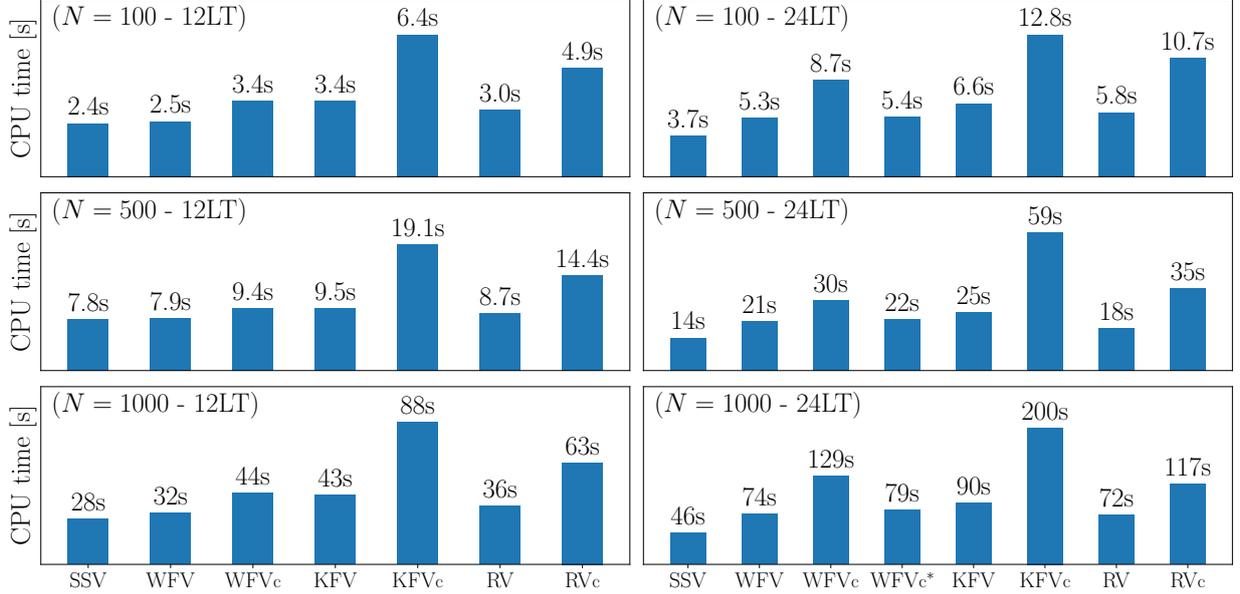


Figure B.19: CPU time required for a single network of size N to perform a 7×7 Grid Search in hyperparameters space in the 12LT and 24LT datasets of the Lorenz system. The validation strategies are the Single Shot Validation (SSV), the Walk Forward Validation (WFV), K-Fold Validation (KFV), Recycle Validation (RV), and respective chaotic versions (subscript c). The runs are on a single Intel i7-8750H processor.

Appendix C. Bayesian Optimization for hyperparameters

After we evaluate the objective function at N_{st} starting points, the objective function is reconstructed in the hyperparameter search space using the function evaluations as data points for noise-free Gaussian Process Regression. The computational cost of the regression is proportional to N_d^3 , where N_d is the number of data points, because of the inversion of the covariance matrix. The inversion is performed by Cholesky factorization regularized by the addition of $\alpha = 10^{-10}$ on the diagonal elements.

Once the Gaussian Process is performed, the next point at which to evaluate the objective function is selected in the hyperparameter space to maximize the acquisition function. The acquisition function evaluates a potential point usefulness in finding the global minimum, so that points with a high value of the acquisition function are selected during the search. A new point can be chosen for one of two reasons: (i) to try to find a new minimum by using current knowledge of the search space and (ii) to increase the knowledge of the space by exploring new regions. This trade-off is called balance between exploitation and exploration. Practically, the most used acquisition functions in the literature are the Probability of Improvement (PI), the Expected Improvement (EI) and the Lower Confidence Bound (LCB) [59]. On a given testcase, it is difficult to determine a priori which acquisition function will perform better. For this reason we use the gp-hedge algorithm [61], which improves the performance with respect to the single acquisition functions. In the algorithm, when deciding the next point of the search, the three acquisition functions are evaluated over the search space. Each acquisition function provides its own optimal point as a candidate. The next point at which the function is going to be evaluated is selected among the three candidates with probability given by the softmax function. The softmax function is evaluated on the cumulative reward from previous candidate points proposed by the acquisition functions, so that the strategy leans towards exploitation as the search progress. Once the point is selected, the Gaussian Process Regression is performed again using

the updated set of data points, until the prescribed maximum number of function evaluations is reached. More details are reported in the Supplementary Material, S.2.

Appendix D. Ensemble size and number of starting points in the test set

First, we select the number of networks in the ensemble through the convergence of the low-order moments of the statistics of the ensemble in the validation set. In Fig. D.20, we show the convergence of the Mean Squared Error (MSE) in the validation set for the chaotic Recycle Validation and chaotic K-fold Validation for the Lorenz system. For $N_{\text{ens}} = 50$ networks, indicated by the vertical line, the 25th, 50th and 75th percentiles have approximately converged to their asymptotic values. Second, we select the number of starting points in the test set, N_t , through the convergence of the statistical properties of the ensemble in the test set. The starting points are equally spaced by 3 LTs, and start from 24 LTs in the time series of Fig. 3. In Fig. D.21, we show the convergence of the Prediction Horizon. For $N_t = 100$ starting points, indicated by the vertical line, the 25th, 50th and 75th percentiles have approximately converged to their asymptotic values. We repeat the procedure to decide the number of starting points for the chaotic, $N_t = 75$, and quasiperiodic, $N_t = 50$, datasets in the Kutznetsov oscillator (results not shown). The starting points are equally spaced by 2 LTs, and start from 7.5 LTs in the time series of Fig. 12.

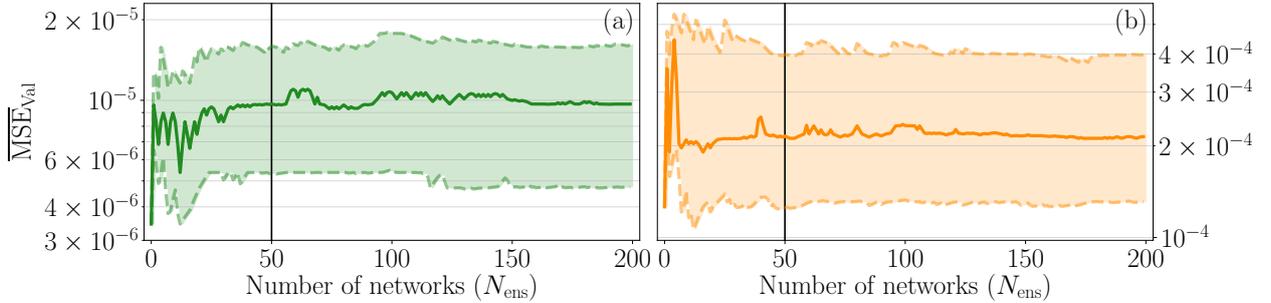


Figure D.20: 50th (continuous line) and 25th and 75th percentiles (dashed lines) for the Mean Squared Error in the validation set as a function of the number of networks in the ensemble in the short dataset of the Lorenz system. The hyperparameters are obtained through Bayesian Optimization in (a) chaotic Recycle Validation and (b) chaotic K-Fold Validation.

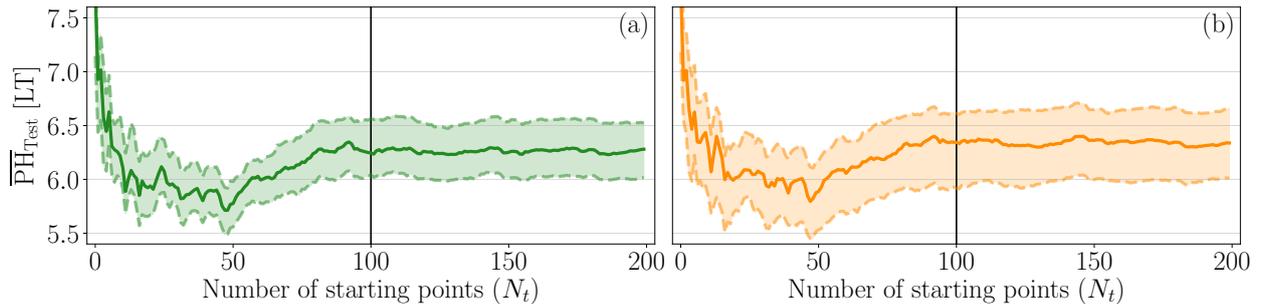


Figure D.21: 50th (continuous line) and 25th and 75th (dashed lines) percentiles for the Prediction Horizon in the test set for the ensemble as a function of the number of starting points in the test set in the short dataset of the Lorenz system. The hyperparameters are obtained through Bayesian Optimization in (a) chaotic Recycle Validation and (b) chaotic K-Fold Validation.

Appendix E. Hyperparameter variations for different realizations

As shown in Fig. 6, different network realizations have different optimal hyperparameters, which vary significantly from one network realization to another other. This suggests that different networks need to be

trained independently. If we select a fixed set of hyperparameters, some networks will perform poorly [45]. In this section, we quantify the difference in performance between optimizing the network independently and using a fixed set of hyperparameters for the entire ensemble. Figure E.22 shows the mean of the Gaussian Process reconstruction of the $\log_{10}(\text{MSE})$ in the test set. In panels (a,b), we show the MSE in the test set for two representative networks from the ensemble, while in panel (c), we show the error between the two networks. The two networks differ substantially. The same hyperparameters may result in MSEs that differ by more than four orders of magnitude.

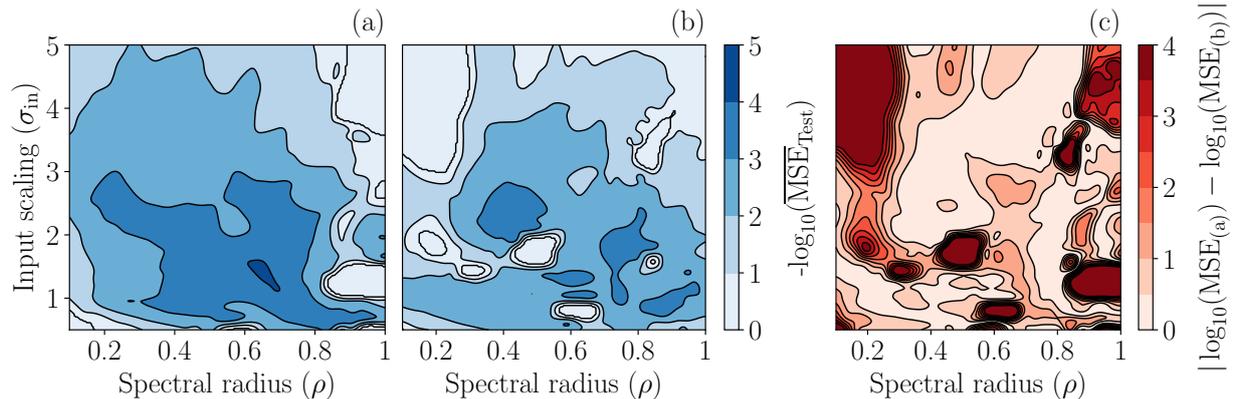


Figure E.22: Mean of the Gaussian Process reconstruction of the MSE in the test set for (a,b) two representative networks in the short dataset, and (c) difference between the two networks. For visualization purposes we saturate the MSE to be ≤ 1 and the error to be $\leq 10^4$. The Gaussian Process is based on a grid of 30×30 data points. For the same hyperparameters, the MSE can differ by orders of magnitude between the two networks.

To quantitatively evaluate the performance of the networks, we assess two possible choices of fixed hyperparameters: (i) we search the optimal fixed hyperparameters by minimizing the geometric mean over the 50 networks of the MSE in the validation set; (ii) we use the hyperparameters obtained by performing the search on a representative network from the ensemble and use that hyperparameters for all the networks. In both (i) and (ii), we perform the search using Bayesian Optimization in the chaotic K-Fold Cross Validation (KFV_c) and chaotic Recycle Validation (RV_c). Figure E.23 shows the violin plots and 25th, 50th and 75th percentiles for the Prediction Horizon in the test set for the Lorenz system. Using fixed hyperparameters yields a decrease in performance in the percentiles of around 0.5 LTs when using (i), and of more than 1 LTs when using (ii). In addition, the tail of the distribution prolongates to values of the Prediction Horizon below 1 LT, which means that the fixed hyperparameters perform poorly in a fraction of the networks. Finally, we note that the decrease in the Prediction Horizon percentiles for (ii) is larger than the improvement that we obtain when using the new validation strategies, the increased size of the dataset or the model-informed architecture. This means that optimizing the network independently, and therefore not using hyperparameters obtained from validating a network in another network, is key in Echo State Networks. Similar conclusions can be drawn for the quasiperiodic and chaotic datasets in the Kuznetsov oscillator (Fig. E.24,E.25).

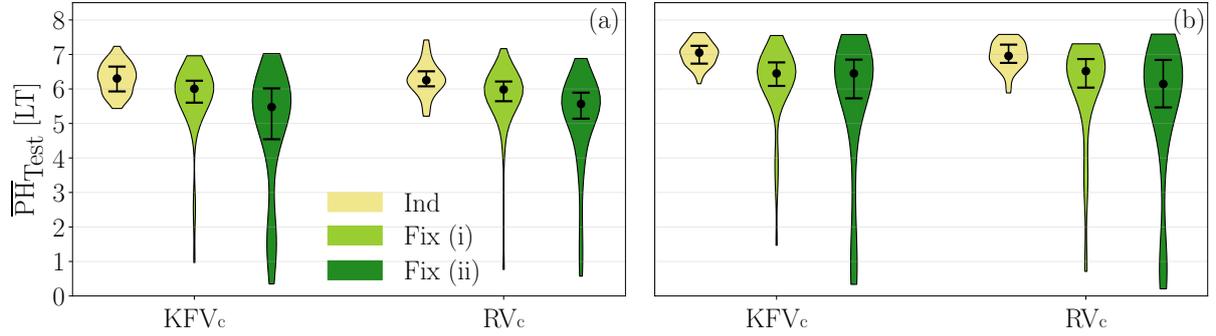


Figure E.23: Violin plots and 25th (lower bar), 50th (marker) and 75th (upper bar) percentiles of the Prediction Horizon in the test set for the 50 networks ensemble in the (a) short (b) and long datasets in the Lorenz system. Independent optimization (Ind) of each network, optimal set of fixed hyperparameters (Fix (i)), and optimal hyperparameters of a single network (Fix (ii)). We use Bayesian Optimization in the chaotic K-Fold Validation (KFV_c) and chaotic Recycle Validation (RV_c).

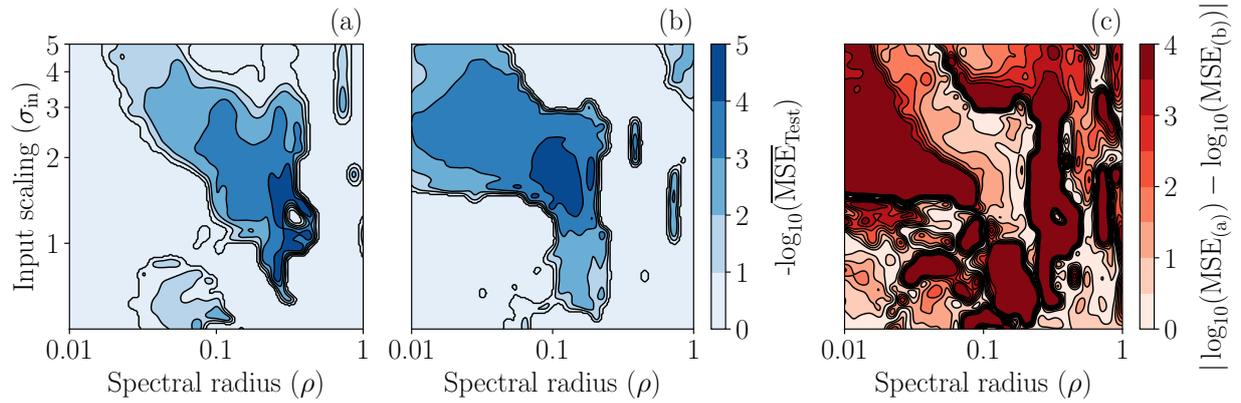


Figure E.24: Mean of the Gaussian Process reconstruction of the MSE in the test set for (a,b) two representative networks in the quasiperiodic dataset, and (c) difference between the two networks. For visualization purposes we saturate the MSE to be ≤ 1 and the error to be $\leq 10^4$. The Gaussian Process is based on a grid of 30×30 data points.

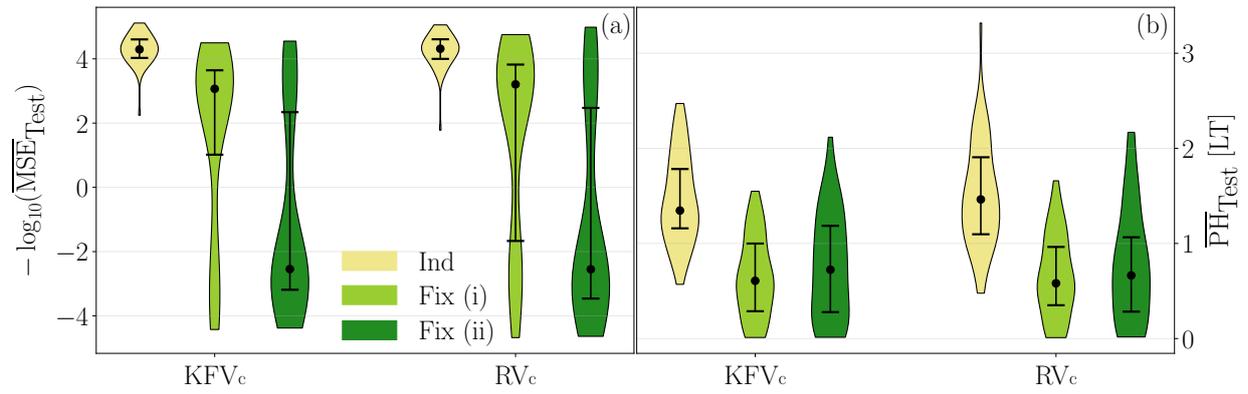


Figure E.25: Violin plots and 25th (lower bar), 50th (marker) and 75th (upper bar) percentiles for the 50 networks ensemble of the MSE in the quasiperiodic dataset, (a), and the Prediction Horizon for the chaotic dataset, (b), in the test set in the Kuznetsov Oscillator. Independent optimization (Ind) of each network, optimal set of fixed hyperparameters (Fix (i)), and optimal hyperparameters of a single network (Fix (ii)). We use Bayesian Optimization in the chaotic K-Fold Validation (KFV_c) and chaotic Recycle Validation (RV_c).