



# Boosting performance in machine learning of geophysical flows via scale separation

Davide Faranda<sup>1,2,3</sup>, Mathieu Vrac<sup>1</sup>, Pascal Yiou<sup>1</sup>, Flavio Maria Emanuele Pons<sup>1</sup>, Adnane Hamid<sup>1</sup>, Giulia Carella<sup>1</sup>, Cedric Ngoungue Langué<sup>1</sup>, Soulivanh Thao<sup>1</sup>, and Valerie Gautard<sup>4</sup>

<sup>1</sup>Laboratoire des Sciences du Climat et de l'Environnement, CE Saclay l'Orme des Merisiers, UMR 8212 CEA-CNRS-UVSQ, Université Paris-Saclay & IPSL, 91191 Gif-sur-Yvette, France.

<sup>2</sup>London Mathematical Laboratory, 8 Margravine Gardens, London, W68RH, UK.

<sup>3</sup>LMD/IPSL, Ecole Normale Supérieure, PSL research University, Paris, France

<sup>4</sup>DRF/IRFU/DEDIP//LILAS Département d'Electronique des Detecteurs et d'Informatique pour la Physique, CE Saclay l'Orme des Merisiers, 91191 Gif-sur-Yvette, France.

**Correspondence:** Davide Faranda (davide.faranda@lscce.ipsl)

**Abstract.** Recent advances in statistical and machine learning have opened the possibility to forecast the behavior of chaotic systems using recurrent neural networks. In this article we investigate the applicability of such a framework to geophysical flows, known to involve multiple scales in length, time and energy and to feature intermittency. We show that both multiscale dynamics and intermittency introduce severe limitations on the applicability of recurrent neural networks, both for short-term forecasts, as well as for the reconstruction of the underlying attractor. We suggest that possible strategies to overcome such limitations should be based on separating the smooth large-scale dynamics from the intermittent/small-scale features. We test these ideas on global sea-level pressure data for the past 40 years, a proxy of the atmospheric circulation dynamics. Better short and long term forecasts of sea-level pressure data can be obtained with an optimal choice of spatial coarse grain and time filtering.

10

## 1 Introduction

The advent of high-performance computing has paved the way for advanced analyses of high-dimensional datasets (Jordan and Mitchell, 2015; LeCun et al., 2015). Those successes have naturally raised the question of whether it is possible to learn the behavior of a dynamical system without resolving or even without knowing the underlying evolution equations. Such an interest is motivated on one side by the fact that many complex systems still miss a universally accepted state equation — e.g. brain dynamics (Bassett and Sporns, 2017), macro-economical and financial systems (Quinlan et al., 2019) — and, on the other, by the need of reducing the complexity of the dynamical evolution for the systems of which the underlying equations are known — e.g. on geophysical and turbulent flows (Wang et al., 2017). Evolution equations are difficult to solve for large systems such as the geophysical flows, so that approximations and parameterizations are needed for meteorological and climatological



20 applications (Buchanan, 2019). These difficulties are enhanced by those encountered in the modelling of phase transitions  
that lead to cloud formation and convection, which are major sources of uncertainty in climate modelling (Bony et al., 2015).  
Machine Learning techniques capable of learning geophysical flows dynamics would help improve those approximations and  
avoid running costly simulations resolving explicitly all spatial/temporal scales.

Recently, several efforts have been made to apply machine learning to the prediction of geophysical data (Wu et al., 2018),  
25 to learn parameterizations of subgrid processes in climate models (Krasnopolsky and Fox-Rabinovitz, 2006; Rasp et al., 2018;  
Gentine et al., 2018), to the forecasting (Liu et al., 2015; Grover et al., 2015; Haupt et al., 2018) and nowcasting (i.e. extremely  
short-term forecasting) of weather variables (Xingjian et al., 2015; Shi et al., 2017; Sprenger et al., 2017), and to quantify the  
uncertainty of deterministic weather prediction (Scher and Messori, 2018). One of the greatest challenge is to replace equations  
of climate models with neural network capable to produce reliable long and short term forecast of meteorological variables.  
30 A first great step in this direction was the use of Echo State Networks (ESN, (Jaeger, 2001)) a particular case of Recurrent  
Neural Networks (RNN) to forecast the behavior of chaotic systems, such as the Lorenz (1963) and the Kuramoto-Sivashinsky  
dynamics (Hyman and Nicolaenko, 1986) . It was shown that RNN predictions of both systems attain performances comparable  
to those obtained with the real equations (Pathak et al., 2017, 2018). Good performance of regularized RNN in the short-term  
prediction of multidimensional chaotic time series was obtained, both from simulated and real data (Xu et al., 2018). This  
35 success motivated several follow-up studies with a focus on meteorological and climate data. These are based on the idea of  
feeding various statistical learning algorithms with data issued from dynamical systems of different complexity, in order to  
study short-term predictability and long-term capabilities of RNN in producing a surrogate dynamics of the input data. Recent  
examples include equation-informed moment-matching for the Lorenz96 model (Lorenz, 1996; Schneider et al., 2017), multi-  
layer perceptrons to reanalysis data (Scher, 2018), or convolutional neural networks to simplified climate simulation models  
40 (Dueben and Bauer, 2018; Scher and Messori, 2019). All these learning algorithms were capable to provide some short-term  
predictability, but failed in obtaining a long-term behavior coherent with the input data.

In this article we specifically focus on how to improve the performance of ESN in simulating long trajectories of large-scale  
climate fields. With respect to the results presented in Pathak et al. (2018), we aim at going beyond the predictability horizon  
and investigate the ability of machine learning algorithms in shadowing the dynamics of observed data. Such applications  
45 would avoid the use of general circulation models based on primitive equations to reproduce the evolution of a subset of  
variables and therefore obtain surrogates dynamics of existing datasets with little computational power. Previous results (Scher,  
2018; Dueben and Bauer, 2018; Scher and Messori, 2019) suggest that RNN simulations of large-scale climate fields are not  
as straightforward as those of the chaotic systems considered by (Pathak et al., 2018). We identify two main mechanisms  
responsible for these limitations: (i) the non-trivial interactions with small-scale motions carrying energy at large scale and  
50 (ii) the intermittent nature of the dynamics. Intermittency triggers large fluctuations of observables of the motion in time and  
space (Schertzer et al., 1997) and can result in non-smooth trajectories within the flow, leading to local unpredictability and  
increasing the number of degrees of freedom needed to describe the dynamics (Paladin and Vulpiani, 1987).

By applying ESN to multiscale and intermittent systems, we investigate how scale separation improves ESN predictions.  
Our goal is to reproduce a surrogate of the large-scale dynamics of global sea-level pressure fields, a proxy of the atmospheric



55 circulation. We begin by analysing three different dynamical systems: we simulate the effects of small scales by artificially  
introducing small-scale dynamics in the Lorenz 1963 equations (Lorenz, 1963) via additive noise. We investigate the Pomeau-  
Manneville equations (Manneville, 1980) stochastically perturbed with additive noise to have an example of intermittent be-  
havior. We then analyse the performance of ESN in the Lorenz 1996 system (Lorenz, 1996). This system dynamics is meant to  
mimic that of the atmospheric circulation and feature both large-scale and small-scale variables with an intermittent behavior.  
60 For all of those systems, as well as for the sea-level pressure data, we show how the performance of ESN in predicting the  
behavior of the system deteriorates rapidly when small-scale dynamics feedback to large scale is important. The idea of using  
moving average for scale separation is already established for meteorological variables (Eskridge et al., 1997). We choose the  
ESN framework following the results of Pathak et al. (2017, 2018), and an established literature about its ability to forecast  
chaotic time series and its stability to noise. For example, Shi and Han (2007); Li et al. (2012) analyse and compare the pre-  
65 dictive performance of simple and improved ESN on simulated and observed one-dimensional chaotic time series. We aim  
at understanding this sensitivity in a deeper way, while assessing the possibility to reduce its impact on prediction through  
simple noise reduction methods. The remaining of this article is organised as follows: first, we give an overview of the ESN  
method and provide the description of the systems used. Then, we show the results for the perturbed Lorenz 1963 equations,  
for the Pomeau-Manneville intermittent map, and for the Lorenz 1996 equations. We discuss the improvement in short-term  
70 prediction and the long-term attractor reconstruction obtained with the moving average filter. We conclude by testing these  
ideas on atmospheric circulation data.

## 2 Methods

Reservoir computing” is a variant of recurrent neural networks (RNN) in which the input signal is connected to a fixed and  
random dynamical system called reservoir (Hinaut, 2013). The principle of Reservoir computing consists in projecting first the  
75 input signal to a space of high dimension in order to obtain a non-linear representation of the signal; and then perform a new  
projection (linear regression or ridge regression) between the high-dimensional space and the output units. In our study ESN  
are implemented as follows. The code is given in the appendix and it shows the parameters used for the computations. Let  $u(t)$   
be the  $K$ -dimensional observable consisting of  $t = 1, 2, \dots, T$  time iterations, originating from a dynamical system and  $r(t)$  be  
the  $N$ -dimensional reservoir state, then:

$$80 \quad r(t + dt) = \tanh(Wr(t) + W_{in}u(t)) \quad (1)$$

where  $W$  is the adjacency matrix of the reservoir: its dimensions are  $N \times N$ , and  $N$  is the number of neurons of the reservoir.  
In ESN, the neuron layers of classic deep neural networks are replaced by a single layer consisting of a sparsely connected  
random network, with coefficients uniformly distributed in  $[-0.5; 0.5]$ .  $W_{in}$ , with dimensions  $N \times K$ , is the weight matrix of  
the connections between the input layer and the reservoir and the coefficients are randomly sampled, as for  $W$ . The output of  
85 the network at time step  $t + dt$  is

$$W_{out}r(t + dt) = v(t + dt) \quad (2)$$



where  $v(t + dt)$  is the ESN prediction,  $W_{out}$  with dimensions  $K \times N$ , is the weight matrix of the connections between the reservoir neurons and the output layer. We estimate  $W_{out}$  via a ridge regression (Hastie et al., 2015):

$$W_{out} = v(t + dt)r(t + dt)^T [r(t + dt)r(t + dt)^T - \lambda I]^{-1} \quad (3)$$

90 with  $\lambda = 10^{-8}$ . In the prediction phase we have a recurrent relationship:

$$r(t + dt) = \tanh(Wr(t) + W_{in}W_{out}r(t)). \quad (4)$$

## 2.1 ESN performance indicators

In this paper, we use three different indicators of performance of the ESN:

### 95 Statistical distributional test

As a first diagnostic of the performance of ESN, we aim at assessing whether the marginal distribution of the forecast values for a given dynamical system is significantly different from the invariant distribution of the system itself. To this purpose, we conduct a  $\chi^2$  test (Cochran, 1952), designed as follows. Let  $U$  be a system observable with support  $R_U$  and probability density function  $f_U(u)$ , and let  $u(t)$  be a sample trajectory from  $U$ . Let now  $\hat{f}_U(u)$  be an approximation of  $f_U(u)$ , namely the  
 100 histogram of  $u$  over  $i = 1, \dots, M$  bins. Note that, if  $u$  spans the entire phase space,  $\hat{f}_U(u)$  is the numerical approximation of the Sinai-Ruelle-Bowen measure of the dynamical system (Eckmann and Ruelle, 1985; Young, 2002). Let now  $V$  be the variable generated by the ESN forecasting, with support  $R_V = R_U$ ,  $v(t)$  the forecast sample,  $g_V(v)$  its probability density function and  $\hat{g}_V(v)$  the histogram of the forecast sample. We test the null hypothesis that the marginal distribution of the forecast sample is the same as the invariant distribution of the system, against the alternative hypothesis that the two distributions are significantly  
 105 different:

$$H_0 : f_U(u) = g_V(v) \quad \text{for every } u \in R_U$$

$$H_1 : f_U(u) \neq g_V(v) \quad \text{for any } u \in R_U$$

Under  $H_0$ ,  $\hat{f}_U(u)$  is the expected value for  $\hat{g}_V(v)$ , which implies that observed differences ( $\hat{g}_V(v) - \hat{f}_U(u)$ ) are due to random errors, and are then independent and identically distributed Gaussian random variables. Statistical theory shows that, given  $H_0$   
 110 true, the test statistics

$$\Sigma = \sum_{i=1}^M \frac{(\hat{g}_V(v) - \hat{f}_U(u))^2}{\hat{f}_U(u)} \quad (5)$$

is distributed as a chi-squared random variable with  $M$  degrees of freedom,  $\chi^2(M)$ . Then, to test the null hypothesis at the level  $\alpha$ , the observed value of the test statistics  $\Sigma$  is compared to the critical value corresponding to the  $1 - \alpha$  quantile of the chi-square distribution,  $\Sigma_c = \chi_{1-\alpha}^2(M)$ : if  $\Sigma > \Sigma_c$ , the null hypothesis must be rejected in favour of the specified alternative.



115 In our setup, we encounter two limitations in using the standard  $\chi^2$  test. First, problems may arise when  $\hat{f}_U(u)$ , i.e. if  
the sample distribution does not span the entire support of the invariant distribution of the system. We observe this in a  
relatively small number of cases; since aggregating the bins would introduce unwanted complications, we decide to discard the  
pathological cases, controlling the effect empirically as described below. Moreover, even producing relatively large samples,  
we are not able to actually observe the invariant distribution of the considered system, which would require much longer  
120 simulations. As a consequence, we would observe excessive rejection rates when testing samples generated under  $H_0$ .

We decide to control these two effects by using a Monte Carlo approach. To this purpose, we use 10000 samples using  
the system equations under the null hypothesis, and we compute the test statistic for each one according to Eq. (5). Then, we  
use the  $(1 - \alpha)$  quantile of the empirical distribution of  $\Sigma$  — instead of the theoretical  $\chi^2(M)$  — to determine the critical  
threshold  $\Sigma_c$ . As a last remark, we notice that we are making inference in repeated tests setting, as the performance of the ESN  
125 is tested 10000 times. Performing a high number of independent tests at a chosen level  $\alpha$  increases the observed rejection rate:  
in fact, even if the samples are drawn under  $H_0$ , extreme events become more likely, resulting in an increased probability to  
erroneously reject the null hypothesis. To avoid this problem, we apply the Bonferroni correction (Bonferroni, 1936), testing  
each one of the  $m = 10000$  available samples at the level  $\alpha' = \frac{\alpha}{m}$ , with  $\alpha = 0.05$ .

Averaging the test results over several sample pairs  $u(t)$ ,  $v(t)$  we obtain a rejection rate  $0 < \phi < 1$  that we use to measure  
130 the adherence of a ESN trajectory  $v(t)$  to trajectories obtained via the equations. If  $\phi = 0$ , almost all the ESN trajectories can  
shadow original trajectories, if  $\phi = 1$  none of the ESN trajectories resemble those of the systems of equations.

### Predictability Horizon

As a measure of the predictability horizon of the ESN forecast compared to the equations, we use the root mean square error  
(RMSE):

$$135 \quad RMSE(\tau) = \sqrt{\frac{1}{\tau} \sum_{t=1}^{\tau} (u(t) - v(t))^2} \quad (6)$$

and we define the predictability horizon  $\tau_s$  as the first time that RMSE exceeds a certain threshold  $s$ . We link  $s$  to the average  
separation of observations in the observable  $U$  and we fix

$$s = \frac{1}{T-1} \sum_{t=2}^{T-1} [u(t) - u(t-1)].$$

We have tested the sensitivity of results against the exact definition of  $s$ .

We interpret  $\tau_s$  as a natural measure of the Lyapunov time  $\vartheta$ , namely the time it takes for an ensemble of trajectories of a  
dynamical system to diverge (Faranda et al., 2012; Panichi and Turchetti, 2018).

### Initial Forecast Error



140 The initial error is given by  $\eta = RMSE(t = 1)$ , for the first time step after the initial condition at  $t = 0$ . We expect  $\eta$  to reduce as the training time increases. In this phase, the the smaller the initial error will be.

## 2.2 Moving average filter

Equipped with these indicators, we analyze two sets of simulations performed with and without smoothing, which was imple-  
145 mented using a moving average filter. The moving average operation is the integral of  $u(t)$  between  $t$  and  $t - w$ , where  $w$  is the window size of the moving average. The simple moving average filter can be seen as a nonparametric time series smoother (see e.g. (Brockwell and Davis, 2016), chapter 1.5). It can be applied to smooth out (relatively) high frequencies in a time series, both to de-noise the observations of a process or to estimate trend-cycle components, if present. Moving averaging consists, in practice, of replacing the observation  $u(t)$  by a value  $u^{(f)}(t)$ , obtained by averaging the previous  $w$  observations. If the time  
150 dimension is discrete (like in the Pomeau-Manneville system) it is defined as:

$$u^{(f)}(t) = \frac{1}{w} \sum_{i=0}^{w-1} u(t-i), \quad (7)$$

while for continuous time systems (like the Lorenz 1963 system), the sum is formally replaced by an integral:

$$u^{(f)}(t) = \frac{1}{w} \int_{t-w}^t u(\zeta) d\zeta. \quad (8)$$

We can define the residuals as:

$$155 \delta u(t) = u^{(f)}(t) - u(t) \quad (9)$$

In practice, the computation always refers to the discrete time case, as continuous time systems are also sampled at finite time steps. Since Echo State Networks are known to be sensitive to noise (see e.g. Shi and Han (2007)), we exploit the simple moving average filter to smooth out high-frequency noise and assess the results for different smoothing windows  $w$ . We find that the choice of the moving averaging window  $w$  must respect two conditions: it should be large enough to smooth the noise  
160 but smaller than the characteristic time  $\tau$  of the large-scale fluctuations of the system. For chaotic systems,  $\tau$  can be derived knowing the rate of exponential divergence of the trajectories, a quantity linked to the Lyapunov exponents (Wolf et al., 1985), and  $\tau$  is known as the Lyapunov time.

We also remark that we can express explicitly the original variable  $u(t)$  as a function of the filtered variable  $u^{(f)}(t)$  as:

$$165 u(t) = w(u^{(f)}(t) - u^{(f)}(t-1)) + u(t-w) \quad (10)$$

we will test this formula for stochastically perturbed systems to evaluate the error introduced by the use of residuals  $\delta u$ .



### 2.3 Testing ESN on filtered dynamics

Here we describe the algorithm used to test ESN performance on filtered dynamics:

- 170 1. Simulate the reference trajectory  $u(t)$  using the equations of the dynamical systems, where  $u(t)$  has been standardized by subtracting the mean and dividing by its standard deviation.
2. Perform the moving average filter to obtain  $u^{(f)}(t)$ .
3. Extract from  $u^{(f)}(t)$  a training set  $u_{train}^{(f)}(t)$  with  $t \in \{1, 2, \dots, T_{train}\}$ .
4. Train the ESN on  $u_{train}^{(f)}(t)$  dataset.
5. Obtain the ESN forecast  $v^{(f)}(t)$  for  $t \in \{T_{train} + 1, T_{train} + 2, \dots, T\}$ .
- 175 6. Add residuals (Eq. 9) to  $v^{(f)}(t)$  sample as  $v(t) = v^{(f)}(t) + \delta u$ , where  $\delta u$  is randomly sampled from the  $\delta u(t)$  with  $t \in \{1, 2, \dots, T_{train}\}$ .
7. Compare  $v(t)$  and  $u(t > T_{train})$  using the metrics  $\phi$ ,  $\tau$  and  $\eta$ .

As an alternative to step 6, one can also use Eq. (10) and obtain:

$$v(t) = w(v^{(f)}(t) - v^{(f)}(t-1)) + v(t-w), \quad (11)$$

180 that does not require the use of residuals  $\delta u(t)$ .

## 3 Results

The systems we analyze are the Lorenz 1963 attractor (Lorenz, 1963) with the classical parameters, discretized with a Euler scheme and a  $dt = 0.001$ , the Pomeau-Manneville intermittent map (Manneville, 1980), the Lorenz 1996 equations (Lorenz, 1996) and the NCEP sea-level pressure data (Saha et al., 2014).

185

### Lorenz 1963 equations

The Lorenz (Lorenz, 1963) system is a simplified model of Rayleigh-Benard convection, derived by E.N. Lorenz. It is an autonomous continuous dynamical system with three variables  $u \in \{x, y, z\}$  parametrizing respectively the convective motion, the horizontal temperature gradient and the vertical temperature gradient. It writes:

$$\begin{aligned} 190 \quad \frac{dx}{dt} &= \sigma(y - x) + \epsilon \xi_x(t) \\ \frac{dy}{dt} &= -xz + \rho x - y + \epsilon \xi_y(t), \\ \frac{dz}{dt} &= xy - bz + \epsilon \xi_z(t), \end{aligned} \quad (12)$$



where  $\sigma$ ,  $\varrho$  and  $b$  are three parameters,  $\sigma$  mimicking the Prandtl number and  $\varrho$  the reduced Rayleigh number. The Lorenz model is usually defined using Eq. (12), with  $\sigma = 10$ ,  $\varrho = 28$  and  $b = 8/3$ . A deterministic trajectory of the system is shown in Figure 1a). It has been obtained via integrating numerically the Lorenz equations with an Euler scheme ( $dt = 0.001$ ). The systems is perturbed via additive noise:  $\xi_x(t)$ ,  $\xi_y(t)$  and  $\xi_z(t)$  are random variable all drawn from a Gaussian distribution. The initial conditions are randomly selected within a long trajectory of  $5 \cdot 10^6$  iterations. First, we study the dependence of the ESN on the training length in the deterministic system ( $\epsilon = 0$ , Figure 1b-d). We analyse the behavior of the rejection rate  $\phi$  (panel b), the predictability horizon  $\tau_s$  (panel c) and the initial error  $\eta$  (panel d) as a function of the training sample size. Our analysis suggests that  $t \sim 10^2$  is a minimum sufficient choice for the training window. We compare this time to the typical time scales of the motion of the systems, determined via the maximum Lyapunov exponent  $\lambda$ . For the Lorenz 1963 system,  $\lambda = 0.9$ , so that the Lyapunov time  $\vartheta \approx \mathcal{O}(\frac{1}{\lambda}) \approx 1.1$ . From the previous analysis we should train the network at least for  $t > 100\vartheta$ . For the other systems analysed in this article, we take this condition as a lower boundary for the training times.

To show the effectiveness of the moving average filter in boosting the machine-learning performances we produce 10 ESN trajectories obtained without moving average (Figure 2-green) and with (Figure 2-red) a moving average window  $w = 0.01$  and compare them to the reference trajectory (blue) obtained with  $\epsilon = 0.1$ . The value of  $w = 10dt = 0.01$  respects the condition  $w \ll \vartheta$ . Indeed, the RMSE averaged over the two groups of trajectories (Figure 2-b) shows an evident gain of accuracy (a factor of  $\sim 10$ ) when the moving average procedure is applied. We now study in a more systematic way the dependence of the ESN performance on noise intensity  $\epsilon$ , network size  $N$  and for three different averaging windows  $w = 0$ ,  $w = 0.01$ ,  $w = 0.05$ . We produce, for each combination, 100 ESN forecasts. Figure 3 shows  $\phi$  (a),  $\log(\tau_{s=1})$  (b) and  $\log(\eta)$  (c) computed setting  $u \equiv x$  variable of the Lorenz 1963 system (results qualitatively do not depend on the chosen variable). In each panel from left to right the moving average window is increasing, upper sub-Panels are obtained using the exact expression in Eq. 11 and lower panels using the residuals in Eq 9. For increasing noise intensity and for small reservoirs sizes, the performances without moving average (left subpanels) rapidly get worse. The moving average smoothing with  $w = 0.01$  (central sub-panels) improves the performance for  $\log(\tau_{s=1})$  (b) and  $\log(\eta)$  (c), except when the noise is too large ( $\epsilon = 1$ ). When the moving average window is too large (right panels), the performances of  $\phi$  decrease. This failure can be attributed to the fact that residuals  $\delta u$  (Eq.9) are of the same order of magnitude of the ESN predicted fields for  $\epsilon$  large. Indeed, if we use the formula provided in Eq. 11 as an alternative to step 6, we can evaluate the error introduced in the residuals. The results shown in Figure 3 suggest that residuals can be used without problems when the noise is small compared with the dynamics. When  $\epsilon$  is close to one, the residuals overlay the deterministic dynamics and ESN forecast are poor. In this case, the exact formulation in Eq. 11 appears much better.

### Pomeau-Manneville intermittent map

Several dynamical systems, including Earth climate, display intermittency, i.e., the time series of a variable issued by the system can experience sudden chaotic fluctuations, as well as a predictable behavior where the observables have small fluctuations. In atmospheric dynamics, such a behavior is observed in the switching between zonal and meridional phases of the mid-latitude



dynamics if a time series of the wind speed at one location is observed: when a cyclonic structure passes through the area, the wind has high values and large fluctuations, when an anticyclonic structure is present the wind is low and fluctuations are smaller (Weeks et al., 1997; Faranda et al., 2016). It is then of practical interest to study the performance of ESN in Pomeau  
230 Manneville predictions as they are a first prototypical example of the intermittent behavior found in climate data.

In particular, the Pomeau-Manneville (Manneville, 1980) map is probably the simplest example of intermittent behavior, produced by a 1D (here  $u = x$ ) discrete deterministic map given by:

$$x_{t+1} = \text{mod}(x_t + x_t^{1+a}, 1) + \epsilon\xi(t), \quad (13)$$

where  $0 < a < 1$  is a parameter. We use  $a = 0.91$  in this study and a trajectory consisting of  $5 \times 10^5$  iterations. The systems  
235 is perturbed via additive noise  $\xi(t)$  drawn from a Gaussian distribution. It is well known that Pomeau-Manneville systems exhibit sub-exponential separation of nearby trajectories and then the Lyapunov exponent is  $\lambda = 0$ . However, one can define a Lyapunov exponent for the non-ergodic phase of the dynamics and extract a characteristic time scale (Korabel and Barkai, 2009). From this latter reference, we can derive a value  $\lambda \simeq 0.2$  for  $a = 0.91$ , implying  $w < \tau \simeq 5$ . We find that the best match between ESN and equations in terms of the  $\phi$  indicator are obtained for  $w = 3$ .

240

Results for the Pomeau-Manneville map are shown in Figure 4. We first observe that the ESN forecast of the intermittent dynamics of the Pomeau-Manneville map is much more challenging than for the Lorenz system as a consequence of the intermittent behavior of this system. For the simulations performed with  $w = 0$ , the ESN cannot simulate an intermittent behavior, for all noise intensities and reservoir sizes. This is reflected in the behavior of the indicators. In the deterministic limit, the  
245 ESN fails to reproduce the invariant density in 80% of the cases ( $\phi \simeq 0.8$ ). For intermediate noise intensities  $\phi > 0.9$  (Figure 4-a). The predictability horizon  $\log(\tau_{s=0.5})$  for the short term forecast is small (Figure 4d) and the initial error large (Figure 4g). The moving average procedure with  $w = 3$  partially improves the performances (Figure 4b,c,e,f,h,i) and it enables ESN to simulate an intermittent behavior (Figure 5). Performances are again better when using the exact formula (Figure 4b,e,h) than using the residuals  $\delta u$  (Figure 4c,f,i). Figure 5a) shows the intermittent behavior of the data generated with the ESN trained  
250 on moving averaged data of Pomeau-Manneville system (red) and compare to the target time series (blue). ESN simulations do not reproduce the intermittency in the average of the target signal. They only show some second order intermittency in the fluctuations. Figure 5b) displays the power spectra showing in both cases a power law decay, which are typical of turbulent phenomena. Although the intermittent behavior is captured, this realization of ESN shows that the values are concentrated around  $x = 0.5$  for the ESN prediction, whereas the non-intermittent phase peaks around  $x = 0$  for the target data.

255

### The Lorenz 1996 system

Before running the ESN algorithm on actual climate data, we test our idea in a more sophisticated, and yet still idealized, model of atmospheric dynamics, namely the Lorenz 1996 equations (Lorenz, 1996). This model explicitly separates two scales and therefore will provide a good test for our ESN algorithm. The Lorenz 1996 system consists of a lattice of large-scale resolved



260 variables  $X$ , coupled to small-scale variables  $Y$ , whose dynamics can be intermittent, so that  $u \in \{X, Y\}$ . The model is defined via two equations:

$$\begin{aligned}\frac{dX_i}{dt} &= X_{i-1}(X_{i+1} - X_{i-2}) - X_i + F - \frac{hc}{b} \sum_{j=1}^J Y_{j,i}, \\ \frac{dY_{j,i}}{dt} &= cbY_{j+1,i}(Y_{j-1,i} - Y_{j+2,i}) - cY_{j,i} + \frac{hc}{b} X_i\end{aligned}\tag{14}$$

265 where  $i = 1, \dots, I$  and  $j = 1, 2, \dots, J$  denote respectively the number of large-scale  $X$  and small-scale  $Y$  variables. Large-scale variables are meant to represent the meanders of the jet-stream driving the weather at mid-latitudes. The first term on the right-hand side represents advection, the second diffusion, while  $F$  mimics an external forcing. The system is controlled via the parameters  $b$  and  $c$  (the time scale of the the fast variables compared to the small variables) and via  $h$  (the coupling between large and small scales). From now on, we fix  $I = 30, J = 5$  and  $F = b = 10$  as these parameters are typically used to  
270 explore the behavior of the system (Frank et al., 2014). We integrate the equations with an Euler scheme ( $dt = 10^{-3}$ ) from the initial conditions  $Y_{j,i} = X_i = F$ , where only one mode is perturbed as  $X_{i=1} = F + \varepsilon$  and  $Y_{j,i=1} = F + \varepsilon^2$ . Here  $\varepsilon = 10^{-3}$ . We discard about  $2 \cdot 10^3$  iterations to reach a stationary state on the attractor, and we retain  $5 \cdot 10^4$  iterations. When  $c$  and  $h$  vary, different interactions between large and small scales can be achieved. A few examples of simulations of the first mode  $X_1$  and  $Y_1$  are given in Figure 6. Figure 6a,c show simulations obtained for  $h = 1$  by varying  $c$ : the larger  $c$  the more intermittent the  
275 behavior of the fast scales. Figure 6.b,d) show simulations obtained for different coupling  $h$  at fixed  $c = 10$ : when  $h = 0$ , there is no small-scale dynamics.

In the Lorenz 1996 model we can explore what happens to the ESN performances if we turn on and off intermittency and/or the small-to-large-scale coupling, without introducing any additional noise term. Moreover, we can also learn the Lorenz 1996  
280 dynamics on the  $X$  variables only, or learn the dynamics on both  $X$  and  $Y$  variables. The purpose of this analysis is to assess whether the ESN are capable of learning the dynamics of the large-scale variables  $X$  alone, and how this capability is influenced by the coupling and the intermittency of the small-scale variables  $Y$ . Using the same simulations presented in Figure 6, we train the ESN on the first  $2.5 \cdot 10^4$  iterations, and then perform, changing the initial conditions 100 different ESN predictions for  $2.5 \cdot 10^4$  more iterations. We apply our performance indicators not to the entire  $I$ -dimensional  $X$  variable  $(X_1, \dots, X_I)$ , as  
285 the  $\chi^2$  test becomes intractable in high dimensions, but rather to the spatial average of the large-scale variables  $X$ . The behavior of each variable  $X_i$  is similar, so the average is representative of the collective behavior. The rate of failure  $\phi$  is very high (not shown) because even when the dynamics is well captured by the ESN the variables are not scaled and centered as those of the original systems. For the following analysis, we therefore replace  $\phi$  with the  $\chi^2$  distance  $T$  (Eq. (5)). The use of  $T$  allows for better highlighting the differences in the ESN performance with respect to the chosen parameters. The same considerations  
290 also apply to the analysis of the sea-level pressure data reported in the next paragraph.



Results of the ESN simulations for the Lorenz 1996 system are reported in Figure 7. In Figure 7a,c,e) ESN predictions are obtained by varying  $c$  at fixed  $h = 1$ , while in Figure 7b,d,f) by varying  $h$  at fixed  $c = 10$ . The continuous lines refer to results obtained feeding the ESN with only the  $X$  variables, dotted lines with both  $X$  and  $Y$ . For the  $\chi^2$  distance  $T$  (Figure 7a,b),  
295 performances show a large dependence on both intermittency  $c$  and coupling  $h$ . First of all, we remark that learning both  $X$  and  $Y$  variables lead to higher distances  $T$ , except for the non intermittent case,  $c = 1$ . For  $c > 1$ , the dynamics learnt on both  $X$  and  $Y$  never settles on a stationary state resembling that of the Lorenz 1996 model. When  $c > 1$  and only the dynamics of the  $X$  variables is learnt, the dependence on  $N$  when  $h$  is varied is non monotonic and better performances are achieved for  $800 < N < 1200$ . For this range, the dynamics settles on stationary states whose spatio-temporal evolution resembles that of  
300 the Lorenz 1996 model, although the variability of time and spatial scales is different from the target. An example is provided in Figure 8, for  $N = 800$ .

Let us now analyse the two indicators of short-term forecasts. Figure 7c,d) display the predictability horizon  $\tau_s$  with  $s = 1$ . The best performances are achieved for the non-intermittent case  $c = 1$  and learning both  $X$  and  $Y$ . When only  $X$  is learnt, we  
305 again get better performances in terms of  $\tau_s$  for rather small network sizes. The performances for  $c > 1$  are better when only  $X$  variables are learnt. The good performance of ESN in learning only the large-scale variables  $X$  are even more surprising when looking at initial error  $\eta$  (Figure 7), which is one order of magnitude smaller when  $X, Y$  are learnt. Despite this advantage in the initial conditions, the ESN performances on  $(X, Y)$  are better only when the dynamics of  $Y$  is non-intermittent. We find clear indications that large intermittency ( $c = 25$ ) and strong small-to-large scale variables coupling ( $h = 1$ ) worsen the ESN  
310 performances, supporting the claims made for the Lorenz 1963 and the Pomeau-Manneville systems.

### The NCEP sea-level pressure data

We now test the effectiveness of the moving average procedure in learning the behavior of multiscale and intermittent systems on climate data issued by reanalysis projects. We use data from the National Centers for Environmental Prediction (NCEP)  
315 version 2 (Saha et al., 2014) with a horizontal resolution of  $2.5^\circ$ . We adopt the global 6 hourly sea-level pressure (SLP) field from 1979 to 31/08/2019 as the meteorological variable proxy for the atmospheric circulation. It traces cyclones (resp. anti-cyclones) with minima (resp. maxima) of the SLP fields. The major modes of variability affecting mid-latitudes weather are often defined in terms of the Empirical Orthogonal Functions (EOF) of SLP and a wealth of other atmospheric features (Hurrell, 1995; Moore et al., 2013), ranging from teleconnection patterns to storm track activity to atmospheric blocking can be  
320 diagnosed from the SLP field.

In addition to the time moving average filter, we also investigate the effect of spatial coarse-graining the SLP fields by a factor  $c$  and perform the learning on the reduced fields. We use the nearest neighbor approximation, which consist in taking from the original dataset the closest value to the coarse grid. Compared with methods based on averaging or dimension reduction techniques such as EOFs, the nearest neighbors approach has the advantage of not removing the extremes (except if the extreme  
325 is not in one of the closest gridpoint) and preserve cyclonic and anticyclonic structures. For  $c = 2$  we obtain a horizontal



resolution of  $5^\circ$  and for  $c = 4$  a resolution  $10^\circ$ . For  $c = 4$  the information on the SLP field close to the poles is lost. However, in the remaining of the geographical domain, the coarse grained fields still capture the positions of cyclonic and anticyclonic structures. Indeed, as shown in (Faranda et al., 2017), this coarse grain field still preserves the dynamical properties of the original one. There is therefore a certain amount of redundant information on the original  $2.5^\circ$  horizontal resolution SLP fields.

330 The dependence of the quality of the prediction for the sea-level pressure NCEPv2 data on the coarse graining factor  $c$  and on the moving average window size  $w$  is shown in Figure 9. We show the results obtained using the residuals (Eq. 9). Figure 9a-c) show the distance from the invariant density, using the  $\chi^2$  distance  $T$ . Here it is clear that by increasing  $w$ , we get better forecast with smaller network sizes  $N$ . A large difference for the predictability expressed as predictability horizon  $\tau_s$ ,  $s = 1.5$  hPa (Figure 9d-f) emerges when SLP fields are coarse grained. We gain up to 10h in the predictability horizon with respect  
335 to the forecasts performed on the original fields ( $c = 0$ ). This gain is also reflected by the initial error  $\eta$  (Figure 9g-i). From the combination of all the indicators, after a visual inspection, we can identify the best-set of parameters:  $w = 12$  h,  $N = 200$  and  $c = 4$ . Indeed this is the case such that, with the smallest network we get almost the minimal  $\chi^2$  distance  $T$ , the highest predictability (32 h) and one of the lowest initial errors. We also remark that, for  $c = 0$  (panels (c) and (i)), the fit always diverges for small network sizes.

340 We compare in details the results obtained for two 10-year predictions with  $w = 0$ h and  $w = 12$ h at  $N = 200$  and  $c = 4$  fixed. At the beginning of the forecast time (Supplementary Video 1), the target field (panel a) is close to both that obtained with  $w = 0$ h (panel b) and  $w = 12$ h (panel c). When looking at a very late time (Supplementary Video 2), of course we do not expect to see agreement among the three datasets. Indeed we are well beyond the predictability horizon. However, we remark that the dynamics for the run with  $w = 0$ h is steady: positions of cyclones and anticyclones barely evolve with time. Instead,  
345 the run with  $w = 12$ h shows a richer dynamical evolution with generation and annihilation of cyclones. A similar effect can be observed in the ESN prediction of the Lorenz 96 system shown in Figure 8b) where the quasi-horizontal patterns indicate less spatial mobility than the original system (Figure 8a).

In order to assess the performances of the two ESNs with and without moving average in a more quantitative way, we present the space-time distributions in Figure 10a). The distribution obtained for the moving average  $w = 12$ h has more realistic tails and matches better than the run  $w = 0$ h that of the target data. Figure 10b-d) shows the wavelet spectrograms (or scalograms) (Hudgins et al., 1993). The scalogram is the absolute value of the continuous wavelet transform of a signal, plotted  
350 as a function of time and frequency. The target data spectrogram (b) presents a rich structure at different frequencies and some interannual variability. The wavelet spectrogram of non-filtered ESN run  $w = 0$  h (c) shows no short time variability and too large interseasonal and interannual variability. The spectrogram of the target data is better matched by the run with  $w = 12$  h  
355 (d) which shows that, on time scales of days to weeks, there is a larger variability.

#### 4 Discussion

We have analysed the performance of ESN in reproducing both the short and long-term dynamics of observables of geophysical flows. The motivation for this study came from the evidence that a straightforward application of ESN to high dimensional



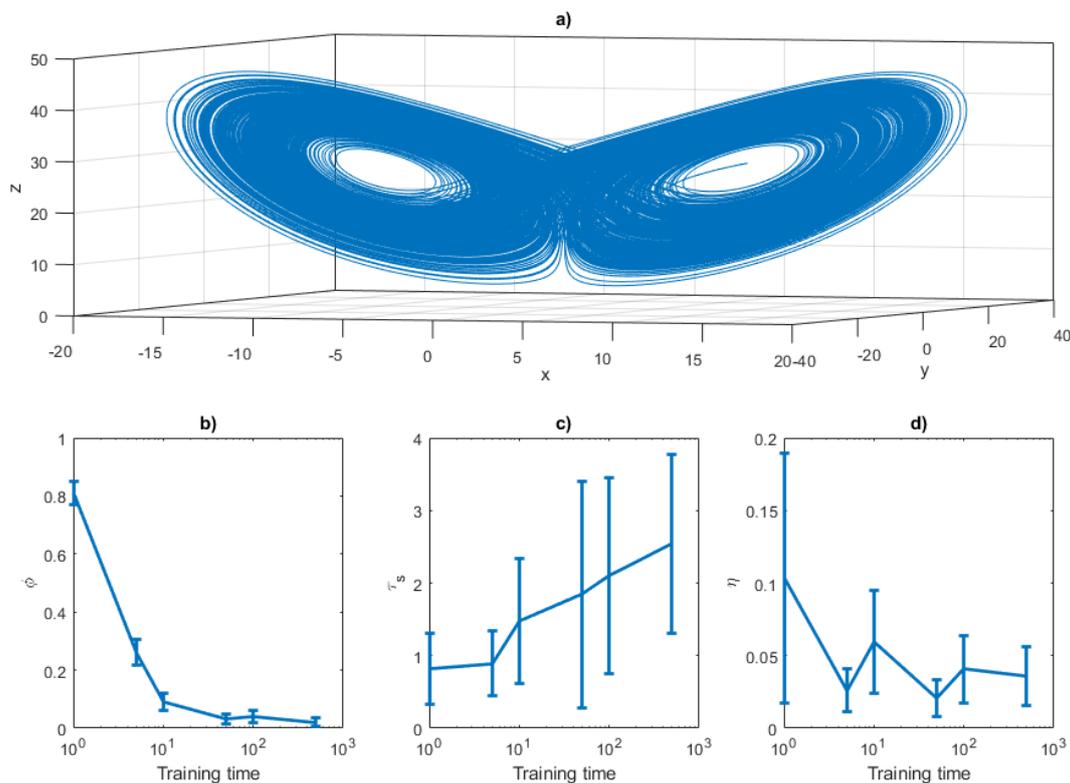
geophysical data (such as the 6 hourly global gridded sea-level pressure data) does not yield to the same results quality obtained  
360 by (Pathak et al., 2018) for the Lorenz 1963 and the Kuramoto-Sivashinsky models. Here we have investigated the causes for  
this behavior and identified two main bottlenecks: (i) intermittency and (ii) the presence of multiple dynamical scales, which  
both appear in geophysical data. In order to illustrate this effect, we have first analysed two low dimensional systems, namely  
the Lorenz (1963) and the Manneville (1980) equation. To mimic multiple dynamical scales, we have added noise terms to  
the dynamics. The performance of ESN in predicting rapidly drops when the systems are perturbed with noise. Filtering the  
365 noise allows to partially recover predictability. It also enables to simulate some qualitative intermittent behavior in the Pomeau-  
Manneville dynamics. This feature could be explored by changing the degree of intermittency in the Pomeau-Manneville map  
as well as performing parameter tuning in ESN. This is left for future work. Here we have used a simple moving-average  
filter and shown that a careful choice of the moving-average window can enhance predictability. As an intermediate step be-  
tween the low-dimensional models and the application to the sea-level pressure data, we have analysed the ESN performances  
370 on the Lorenz (1996) system. This system was introduced to mimic the behavior of the atmospheric jet at mid-latitude, and  
features a lattice of large-scale variables, each connected to small-scale variables. Both the coupling between large and small  
scales and intermittency can be tuned in the model, giving rise to a plethora of behaviors. For the Lorenz 1996 model, we did  
not have to apply a moving average filter to the data, as we can train the ESN on the large-scale variables only. Our computa-  
tions have shown that, whenever the small scales are intermittent, or the coupling is strong, learning the dynamics of the coarse  
375 grained variable is more effective, both in terms of computation time and performances. The results also apply to geophysical  
datasets: here we analysed the atmospheric circulation, represented by sea-level pressure fields. Again we have shown that both  
a spatial coarse-graining and a time moving-average filter improve the ESN performances.

Our results may appear rather counter-intuitive, as the weather and climate modelling communities are moving towards  
380 extending simulations of physical processes to small scales. As an example, we cite the use of highly-resolved convection-  
permitting simulations (Fosser et al., 2015) as well as the use of stochastic (and therefore non-smooth) parameterizations in  
weather models (Weisheimer et al., 2014). We have, however, a few heuristic arguments on why the coarse-graining and filter-  
ing operations should improve the ESN performances. First of all, the moving-average operation helps both in smoothing the  
signal and by providing the ESN with a wider temporal information. In some sense, this is reminiscent of the embedding proce-  
385 dure (Cao, 1997), where the signal behavior is reconstructed by providing not only information on the previous time step, but on  
previous times depending on the complexity. The filtering procedure can also be motivated by the fact that the active degrees of  
freedom for the sea-level pressure data are limited. This has been confirmed by Faranda et al. (2017) via coarse-graining these  
data and showing that the active degrees of freedom are independent on the resolution, in the same range explored in this study.  
In other words, including small scales in the learning of sea-level pressure data, does not provide additional information on  
390 the dynamics and push towards over-fitting and saturating the ESN with redundant information. The latter consideration poses  
also some caveats on the generality of our results: we believe that this procedure is not beneficial whenever a clear separation  
of scales is not achievable, e.g. in non-confined 3-D turbulence. Moreover, in this study, note that three sources of stochasticity  
were present: (i) in the random matrices and reservoir, (ii) in the perturbed initial conditions and (iii) in the ESN simulations

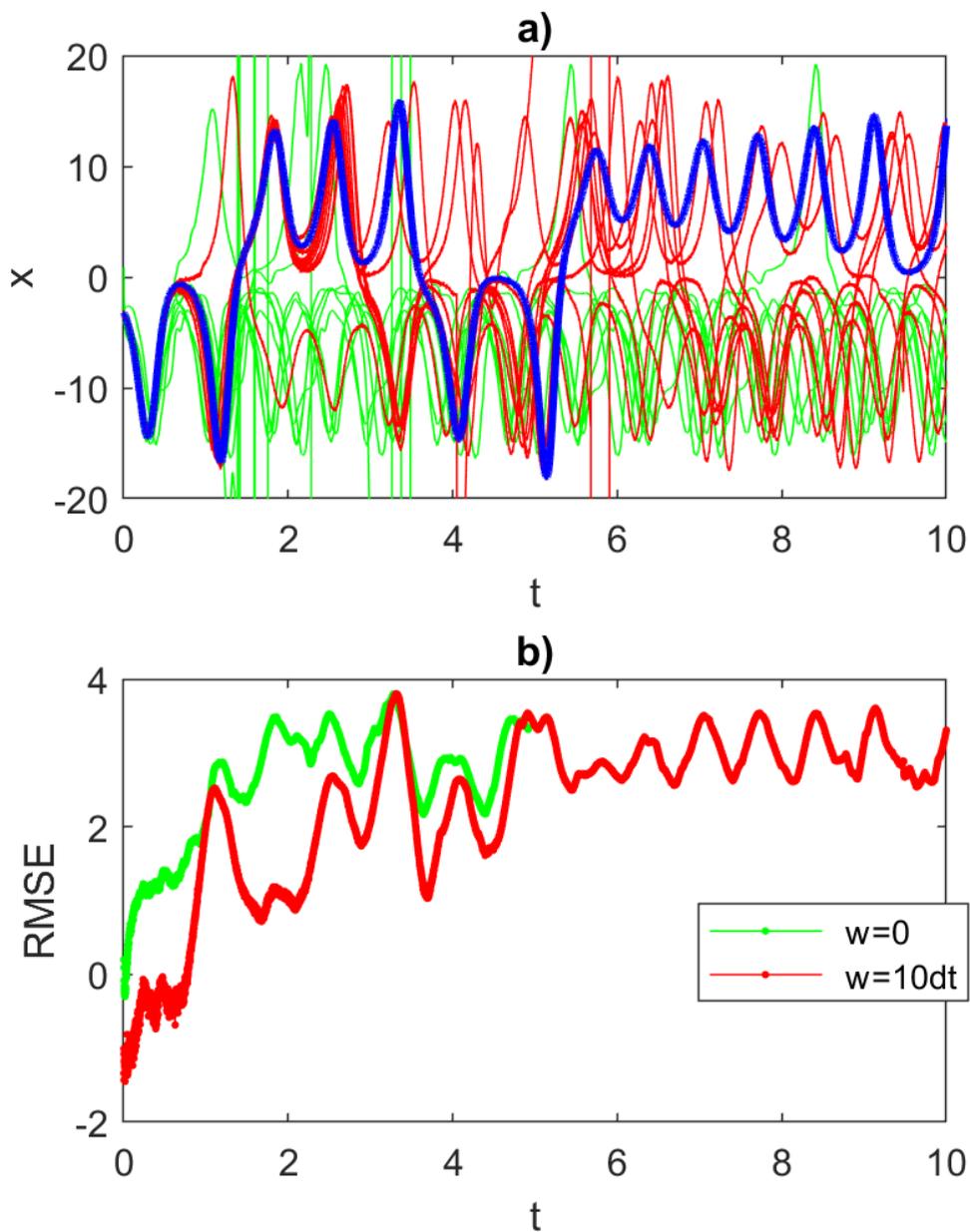


when using moving average filtered data with sampled  $\delta u$  components. The first one is inherent to the model definition. The  
395 perturbations of the starting conditions allow characterizing the sensitivity of our ESN approach to the initial conditions. The  
stochasticity induced by the additive noise  $\delta u$  provides a distributional forecast at each time  $t$ . Although this latter noise can  
be useful to simulate multiple trajectories and evaluate their long-term behaviour, in practice, i.e., in the case where a ESN  
would be used operationally to generate forecasts, one might not want to employ a stochastic formulation with an additive  
noise, but rather the explicit and deterministic formulation in Eq. 11. This exemplifies the interest of our ESN approach for  
400 possible distinction between forecasts and long-term simulations, and therefore makes it flexible to adapt to the case of interest.

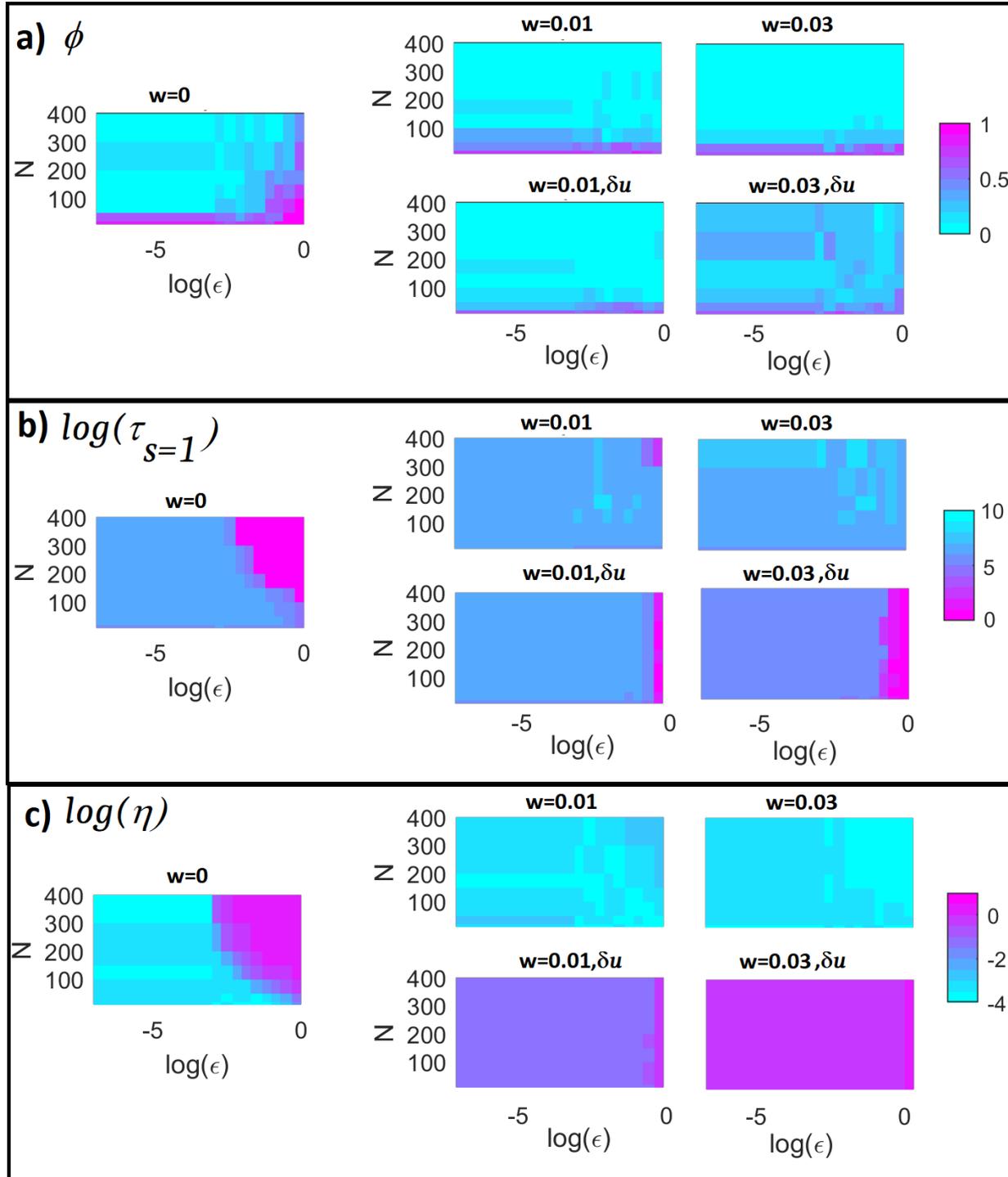
In future work, it will be interesting to use other learning architectures and other methods of separating large- from small-  
scale components (Wold et al., 1987; Froyland et al., 2014; Kwasniok, 1996). For example, our results give a more formal  
framework for applications of machine learning techniques on geophysical data. Deep-learning approaches have proven useful  
405 in performing learning at different time and spatial scales whenever each layer is specialized in learning some specific features  
of the dynamics (Bolton and Zanna, 2019; Gentine et al., 2018). Indeed, several difficulties encountered in the application  
of machine learning on climate data could be overcome if the appropriate framework is used, but this requires a critical  
understanding of the limitations of the learning techniques.



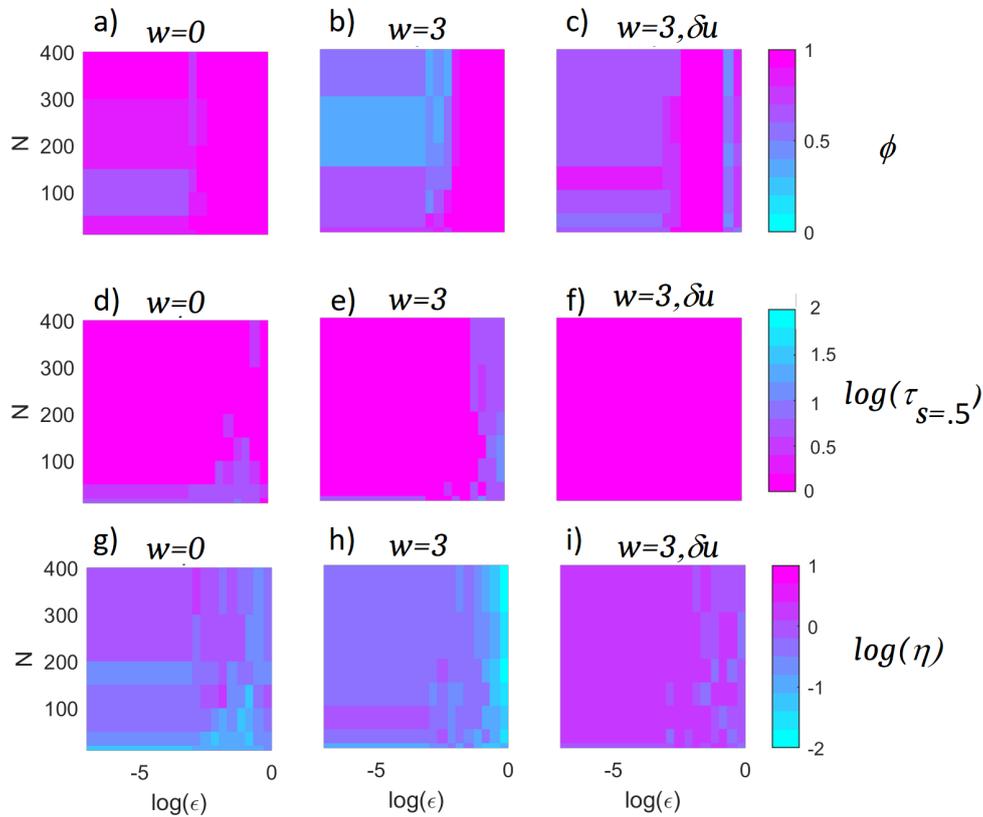
**Figure 1.** a) Lorenz 1963 attractor obtained with a Euler scheme with  $dt = 0.001$ ,  $\sigma = 10$ ,  $r = 28$  and  $b = 8/3$ . Panels b-d) show the performances indicator as a function of the training time. b) the rejection rate  $\phi$  of the invariant density test for the  $x$  variable; c) the first time  $t$  such that the  $RMSE > 1$ ; d) the initial error  $\eta$ . The error bar represents the average and the standard deviation of the mean over 100 realizations.



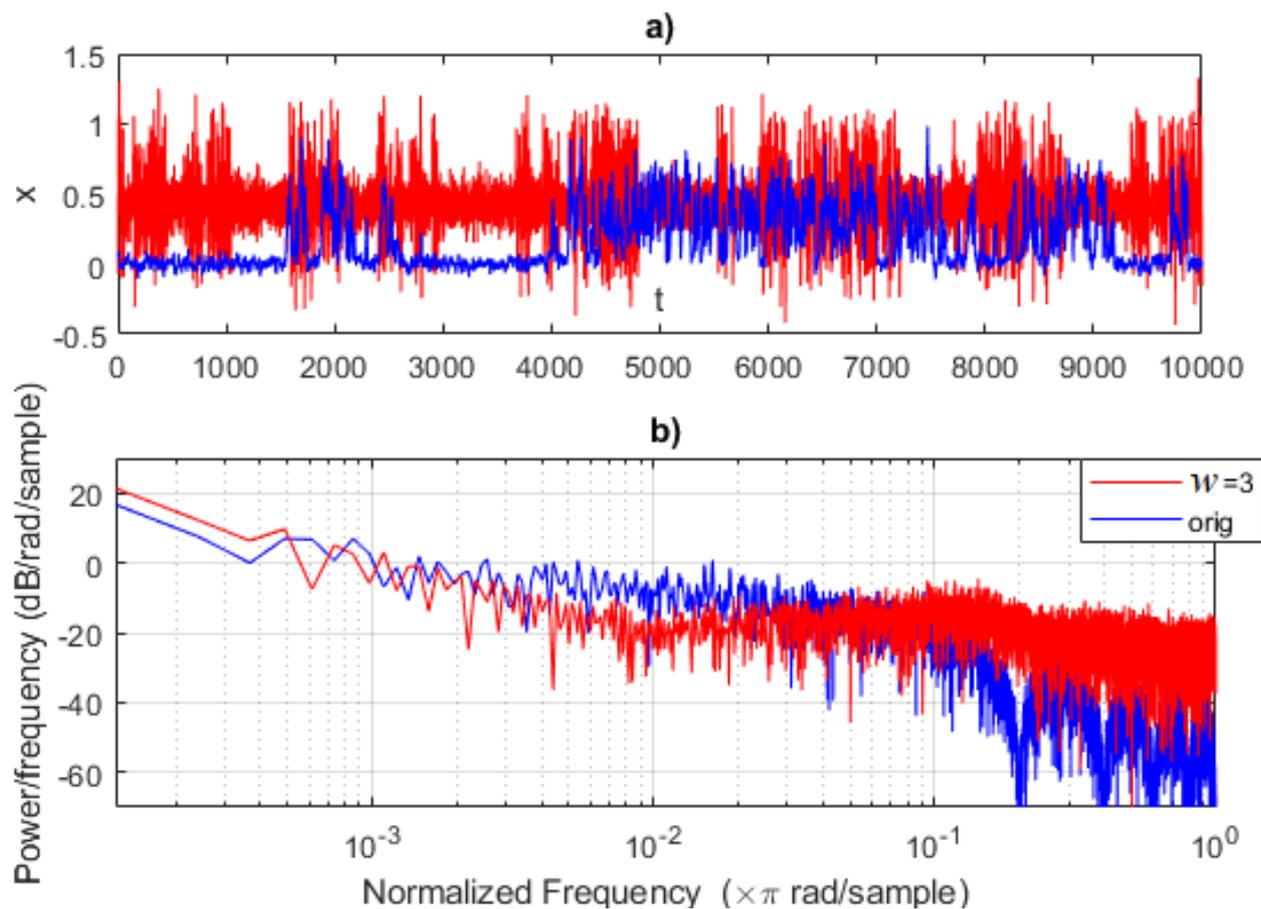
**Figure 2.** a) Trajectories predicted using ESN on the Lorenz 1963 attractor for the variable  $x$ . The attractor is perturbed with Gaussian noise with variance  $\epsilon = 0.1$ . The target trajectory is shown in blue. 10 trajectories obtained without moving average (green) show an earlier divergence compared to 10 trajectories where the moving average is performed with a window size of  $w = 10dt = 0.01$  (red). Panel (b) shows the evolution of the log(RMSE), averaged over the trajectories for the cases with  $w = 0.01$  (red) and  $w = 0$  (green). The trajectories are all obtained after training the ESN for  $10^5$  time-steps. Each trajectory consists of  $10^4$  timesteps.



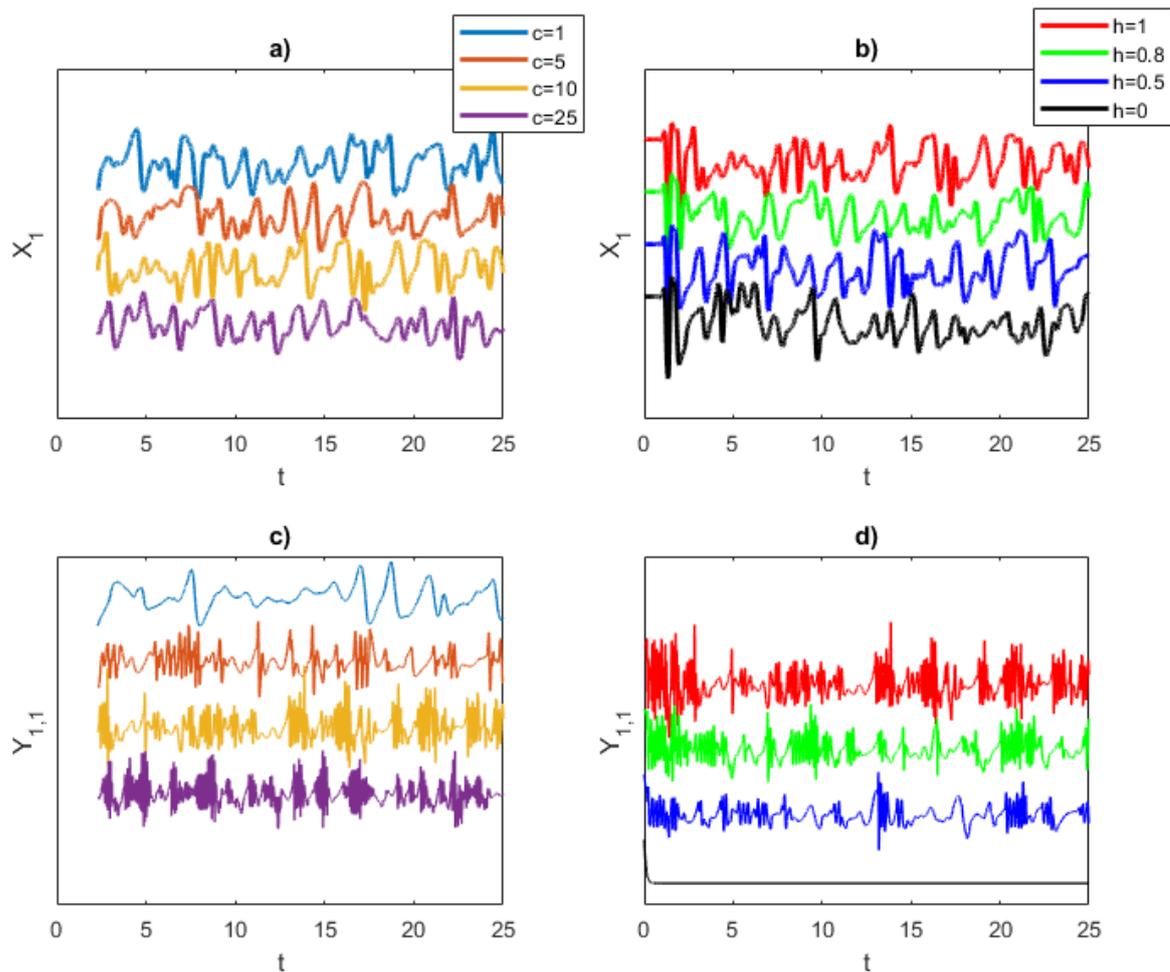
**Figure 3.** Lorenz 1963 analysis for increasing noise intensity  $\epsilon$  (x-axes), and number of neurons  $N$  (y-axes). The colorscale represents:  $\phi$  the rate of failure of the  $\chi^2$  test (size  $\alpha = 0.05$ ) (a); the logarithm of predictability horizon  $\log(\tau_{s=1})$  (b); the logarithm of initial error  $\log(\eta)$  (c). All the values are averages over 30 realizations. Left sub-panels refer to results without moving average, central sub-panels with averaging window  $w = 0.01$ , right hand-side panels with averaging window  $w = 0.03$ . Upper sub-panels are obtained using the exact expression in Eq. 11 and lower sub-panels using the residuals in Eq. 9. The trajectories are all obtained after training the ESN for  $10^5$  time-steps. Each trajectory consists of  $10^4$  timesteps.



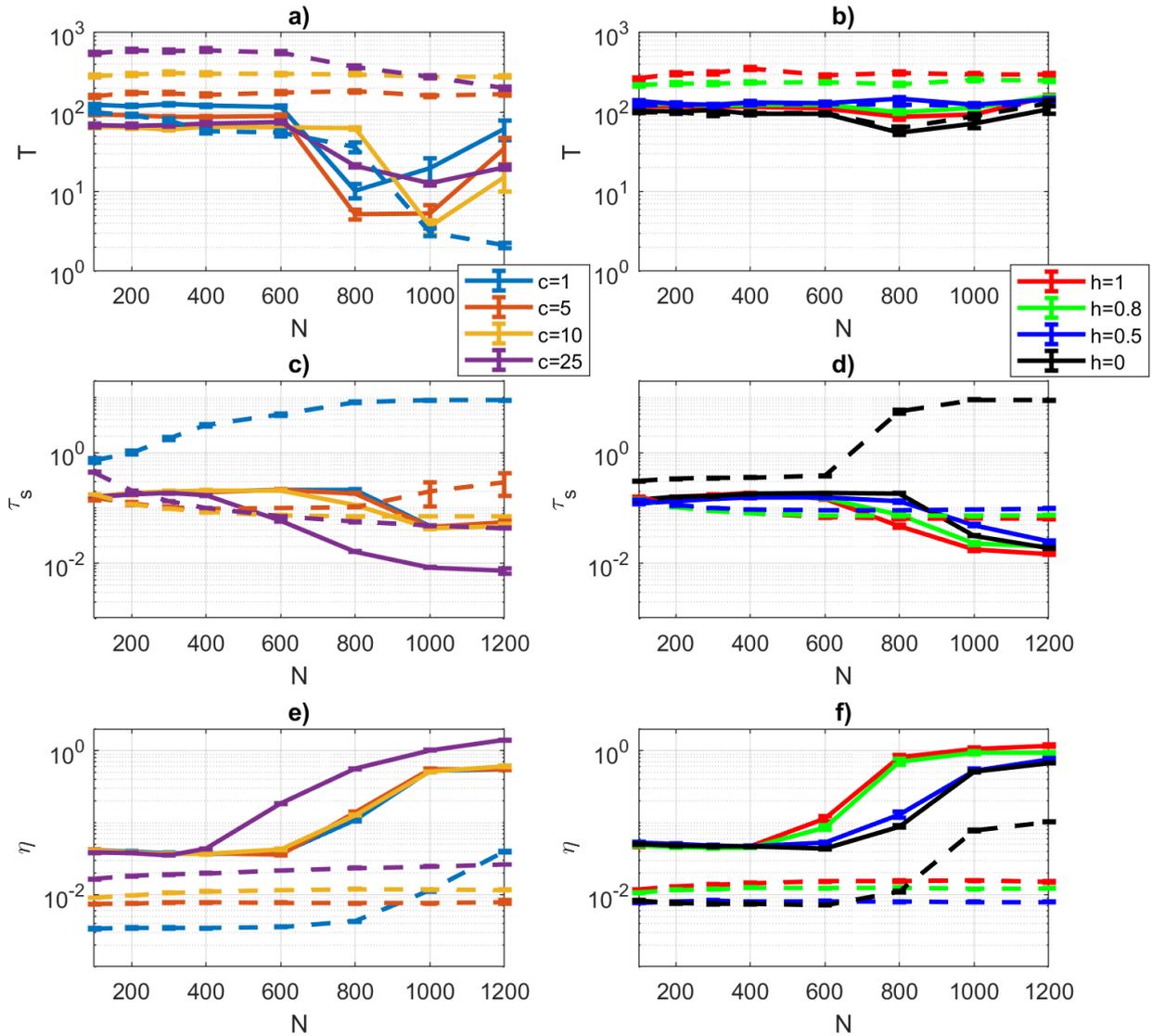
**Figure 4.** Analysis of the Pomeau-Manneville system for increasing noise intensity  $\epsilon$  (x-axes), and number of neurons  $N$  (y-axes). The colorscale represents:  $\phi$  the rate of failure of the  $\chi^2$  test (size  $\alpha = 0.05$ ) (a-c); the logarithm of predictability horizon  $\log(\tau_{s=0.5})$  (d-f); the logarithm of initial error  $\log(\eta)$  (g-i). All the values are averages over 30 realizations. Panels a,d,g) refer to results without moving average, b,c,e,f,h,i) with averaging window  $w = 3$ , c,f,i) using the residuals in Eq 9. The trajectories are all obtained after training the ESN for  $10^5$  time-steps. Each trajectory consists of  $10^4$  timesteps.



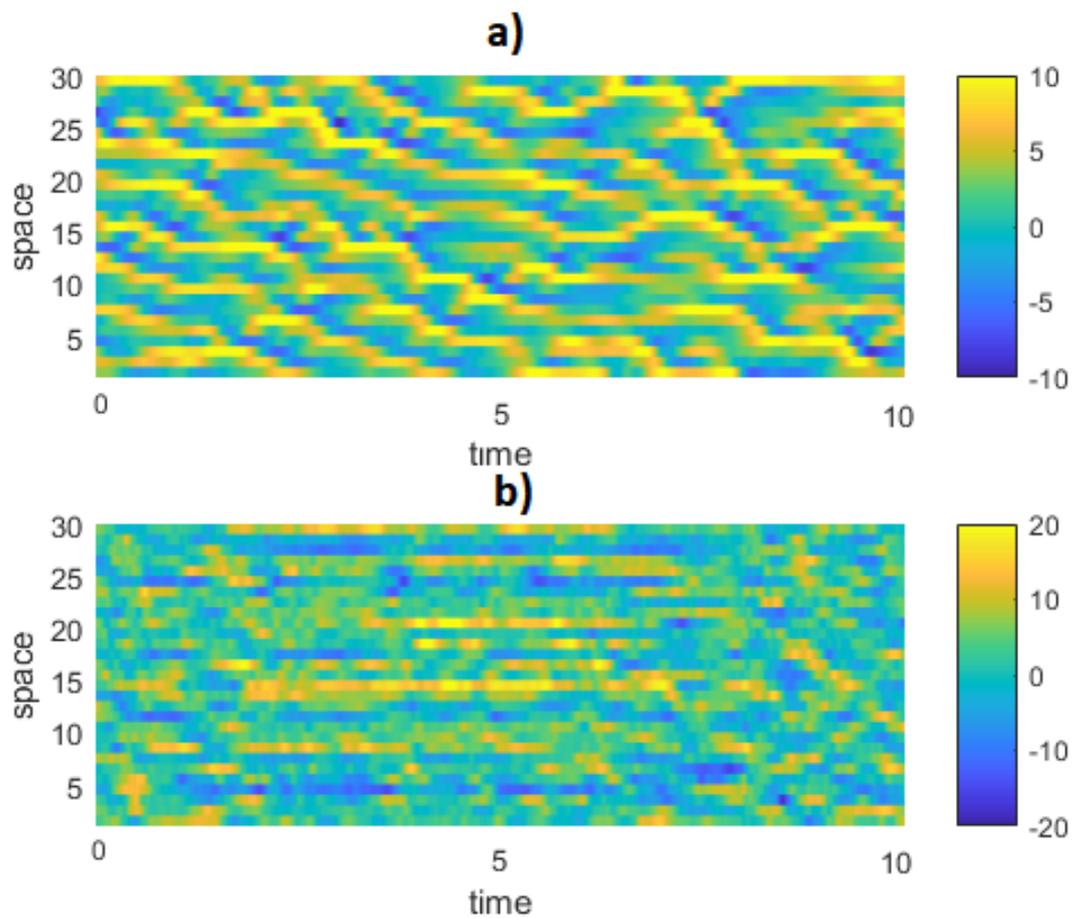
**Figure 5.** Pomeau-Manneville ESN simulation (red) showing an intermittent behavior and compared to the target trajectory (blue). The ESN trajectory is obtained after training the ESN for  $10^5$  time-steps using the moving average time series with  $w = 3$ . It consists of  $10^4$  timesteps. Cases  $w = 0$  are not shown as trajectories always diverge. Evolution of trajectories in time (a) and Fourier power spectra (b).



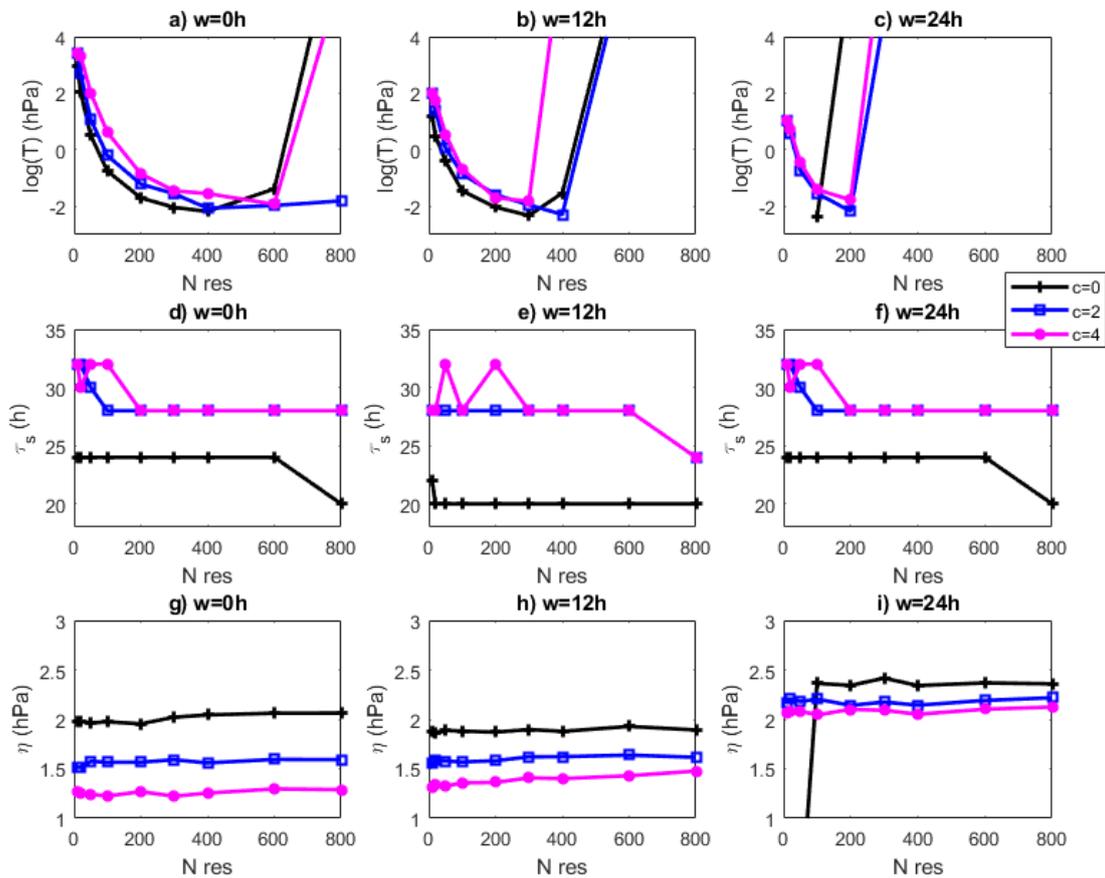
**Figure 6.** Lorenz 1996 simulations for the large-scale variable  $X_1$  (a,b) and small-scale variable  $Y_{1,1}$  (c,d). Panels (a,c) show simulations varying  $c$  at fixed  $h = 1$ . The larger  $c$ , the more intermittent the behavior of the fast scales. Panels (b,d) show simulations varying the coupling  $h$  for fixed  $c = 10$ . When  $h = 0$ , there is no small-scale dynamics.  $y$ -axes are in arbitrary units, time-series are shifted for better visibility.



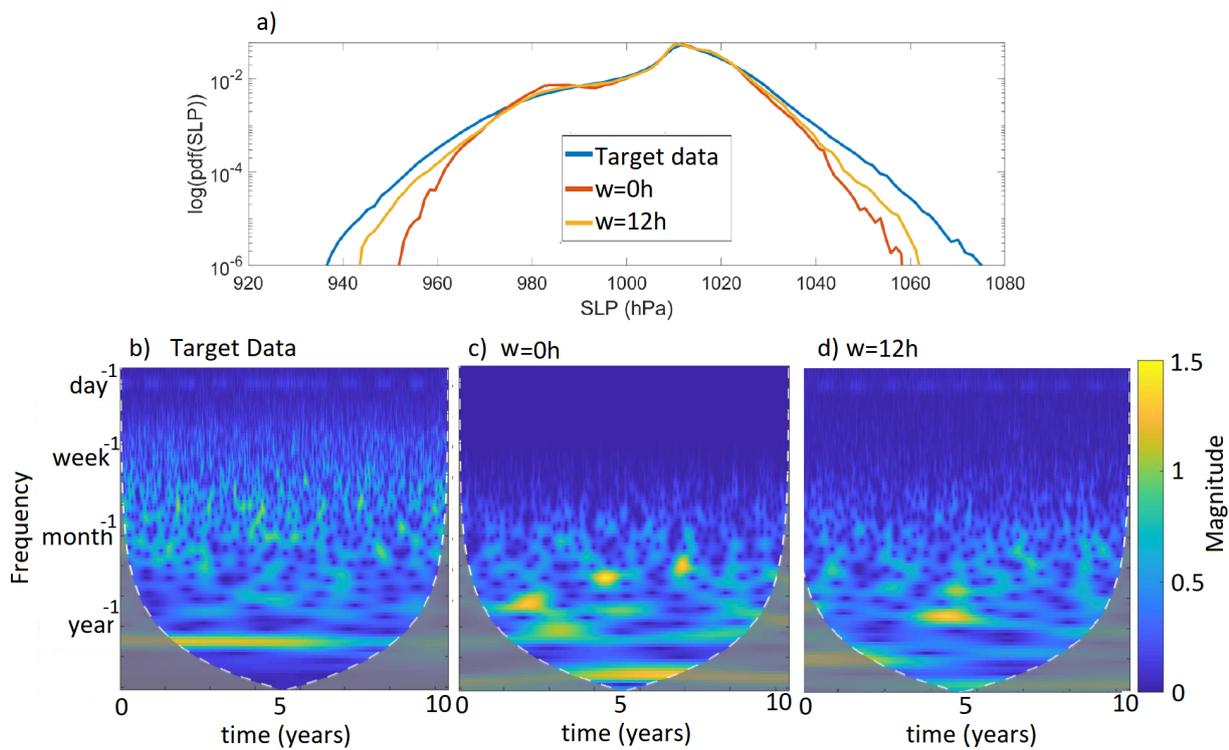
**Figure 7.** Lorenz 1996 ESN prediction performance for the large-scale variables  $X$  only. a,b)  $\chi^2$  distance  $T$ ; (c,d) the predictability horizon  $\tau_s$  with  $s = 1$ . (e,f) the initial error  $\eta$  in hPa. In (a,c,e) ESN predictions are made varying  $c$  at fixed  $h = 1$ . In (b,d,f) ESN predictions are made varying  $h$  at fixed  $c = 10$ . Continuous lines show ESN prediction performance made considering  $X$  variables only, dotted lines considering both  $X$  and  $Y$  variables.



**Figure 8.** Example of (a) target Lorenz 1996 spatio-temporal evolution of large-scale variables  $X$  for  $c = 1, h = 1$  and (b) ESN prediction realized with  $N = 800$  neurons. Note that the colors are not on the same scale for the two panels.



**Figure 9.** Dependence of the quality of the results for the prediction of the sea-level pressure NCEPv2 data on the coarse graining factor  $c$  and on the moving average window size  $w$ . a-c)  $\chi^2$  distance  $T$ ; d-f) predictability horizon (in hours)  $\tau_s$ ,  $s = 1.5$  hPa; g-i) logarithm of initial error  $\eta$ . Different coarse grain factor  $c$  are shown with different colors. a,d,g)  $w = 0$ , b,e,h)  $w = 12$  h, c,f,i)  $w = 24$  h.



**Figure 10.** a) Distributions of 10 years of 6h spatial and temporal data at all grid points obtained for the target NCEPv2 SLP data (blue), an ESN with  $c = 4$  and  $w = 0$  h (red), and an ESN with  $c = 4$  and  $w = 12$  h (orange). b-d) wavelet spectrograms for the NCEPv2 SLP target data (b), a run with  $c = 4$   $w = 0$  h (c), and with  $c = 4$  and  $w = 12$  h (d).



## Appendix A: Numerical code

410 We report here the MATLAB code used for the computation of the Echo State Network. This code is adapted from the original code available here: <https://mantas.info/code/simpleesn>

### A1 ESN Training

```
function [Win, W, Wout]=RNN_training(data,Nres)
%This function train the Echo State network using the data provided.
415 %INPUTS:
%data: a matrix of the input data to train, arranged as space X time
%Nres: the number of neurons N to be used in the training
%OUTPUTS:
%Win: the input weight matrix which consists of random weights
420 %W: the network of neurons
%Wout: the output weights, they are adjusted to match the next iterations
inSize = size(data,1);
trainLen= size(data,2);
Win = (rand(Nres,1+inSize)-0.5) .* 1;
425 W = rand(Nres,Nres)-0.5;
% normalizing and setting spectral radius
opt_disp = 0;
rhoW = abs(eigs(W,1,'LM',opt));
W = W .* ( 1.25 /rhoW);
430 % memory allocation
X = zeros(1+inSize+Nres,trainLen-1);
Yt = data(:,2:end)';
x = zeros(Nres,1);
for t = 1:trainLen-1
435 u = data(:,t);
x = tanh( Win*[1;u] + W*x );
X(:,t) = [1;u;x];
end
reg = 1e-8; % regularization coefficient
440 Wout = ((X*X' + reg*eye(1+inSize+Nres)) \ (X*Yt))';
end
```



## A2 ESN Prediction

```
function [Y_pred]=RNN_prediction(data,Win, W, Wout)
% This function returns the recurrent Echo State Network prediction
445 %INPUT:
%data: the full data matrix of the data to predict in the form (space*time)
%Win: input weights
%W: neurons matrix
%Wout: output weights
450 %OUTPUT:
%Y_pred: the RNN prediction
Y_pred = zeros(size(data,1), size(data,2) );
x = zeros(size(W,1),1);
u=data(:,1);
455 for t = 1:size(data,2)
x = tanh( Win*[1;u] + W*x );
y = Wout*[1;u;x];
Y_pred(:,t) = y;
u = y;
460 end
end
```

*Code and data availability.* The numerical code used in this article is provided in Appendix A

*Author contributions.* Davide Faranda, Mathieu Vrac, Pascal Yiou and Soulivanh Thao conceived this study. Davide Faranda, Adnane Hamid, Cedric Nguoungue Langue and Giulia Carella performed the analysis. All the authors contributed to writing and discussing the  
465 results of the paper.

*Competing interests.* The authors declare that there is no conflict of interest.



*Acknowledgements.* We acknowledge Barbara D'Alena, Julien Brajard, Venkatramani Balaji, Berengere Dubrulle, Robert Vautard, Nikki Vercauteren, Francois Daviaud, Yuzuru Sato for useful discussions. This work is supported by the CNRS INSU-LEFE-MANU grant "DINCLIC".



## 470 References

- Bassett, D. S. and Sporns, O.: Network neuroscience, *Nature neuroscience*, 20, 353, 2017.
- Bolton, T. and Zanna, L.: Applications of deep learning to ocean data inference and subgrid parameterization, *Journal of Advances in Modeling Earth Systems*, 11, 376–399, 2019.
- Bonferroni, C.: Teoria statistica delle classi e calcolo delle probabilita, *Pubblicazioni del R Istituto Superiore di Scienze Economiche e Commerciali di Firenze*, 8, 3–62, 1936.
- 475 Bony, S., Stevens, B., Frierson, D. M., Jakob, C., Kageyama, M., Pincus, R., Shepherd, T. G., Sherwood, S. C., Siebesma, A. P., Sobel, A. H., et al.: Clouds, circulation and climate sensitivity, *Nature Geoscience*, 8, 261, 2015.
- Brockwell, P. J. and Davis, R. A.: *Introduction to time series and forecasting*, Springer, 2016.
- Buchanan, M.: *The limits of machine prediction*, Ph.D. thesis, Nature Publishing Group, 2019.
- 480 Cao, L.: Practical method for determining the minimum embedding dimension of a scalar time series, *Physica D: Nonlinear Phenomena*, 110, 43–50, 1997.
- Cochran, W. G.: The  $\chi^2$  test of goodness of fit, *The Annals of Mathematical Statistics*, pp. 315–345, 1952.
- Dueben, P. D. and Bauer, P.: Challenges and design choices for global weather and climate models based on machine learning, *Geoscientific Model Development*, 11, 3999–4009, 2018.
- 485 Eckmann, J.-P. and Ruelle, D.: Ergodic theory of chaos and strange attractors, in: *The theory of chaotic attractors*, pp. 273–312, Springer, 1985.
- Eskridge, R. E., Ku, J. Y., Rao, S. T., Porter, P. S., and Zurbenko, I. G.: Separating different scales of motion in time series of meteorological variables, *Bulletin of the American Meteorological Society*, 78, 1473–1484, 1997.
- Faranda, D., Lucarini, V., Turchetti, G., and Vaienti, S.: Generalized extreme value distribution parameters as dynamical indicators of stability, *International Journal of Bifurcation and Chaos*, 22, 1250 276, 2012.
- 490 Faranda, D., Masato, G., Moloney, N., Sato, Y., Daviaud, F., Dubrulle, B., and Yiou, P.: The switching between zonal and blocked mid-latitude atmospheric circulation: a dynamical system perspective, *Climate Dynamics*, 47, 1587–1599, 2016.
- Faranda, D., Messori, G., and Yiou, P.: Dynamical proxies of North Atlantic predictability and extremes, *Scientific reports*, 7, 41 278, 2017.
- Fosser, G., Khodayar, S., and Berg, P.: Benefit of convection permitting climate model simulations in the representation of convective precipitation, *Climate Dynamics*, 44, 45–60, 2015.
- 495 Frank, M. R., Mitchell, L., Dodds, P. S., and Danforth, C. M.: Standing swells surveyed showing surprisingly stable solutions for the Lorenz'96 model, *International Journal of Bifurcation and Chaos*, 24, 1430 027, 2014.
- Froyland, G., Gottwald, G. A., and Hammerlindl, A.: A computational method to extract macroscopic variables and their dynamics in multiscale systems, *SIAM Journal on Applied Dynamical Systems*, 13, 1816–1846, 2014.
- 500 Gentine, P., Pritchard, M., Rasp, S., Reinaudi, G., and Yacalis, G.: Could machine learning break the convection parameterization deadlock?, *Geophysical Research Letters*, 45, 5742–5751, 2018.
- Grover, A., Kapoor, A., and Horvitz, E.: A deep hybrid model for weather forecasting, in: *Proceedings of the 21th ACM SIGKDD International Conference on Knowledge Discovery and Data Mining*, pp. 379–386, ACM, 2015.
- Hastie, T., Tibshirani, R., and Wainwright, M.: *Statistical learning with sparsity: the lasso and generalizations*, Chapman and Hall/CRC, 505 2015.



- Haupt, S. E., Cowie, J., Linden, S., McCandless, T., Kosovic, B., and Alessandrini, S.: Machine Learning for Applied Weather Prediction, in: 2018 IEEE 14th International Conference on e-Science (e-Science), pp. 276–277, IEEE, 2018.
- Hinaut, X.: Réseau de neurones récurrent pour le traitement de séquences abstraites et de structures grammaticales, avec une application aux interactions homme-robot, Ph.D. thesis, Thèse de doctorat, Université Claude Bernard Lyon 1, 2013.
- 510 Hudgins, L., Friehe, C. A., and Mayer, M. E.: Wavelet transforms and atmospheric turbulence, *Physical Review Letters*, 71, 3279, 1993.
- Hurrell, J. W.: Decadal trends in the North Atlantic Oscillation: regional temperatures and precipitation, *Science*, 269, 676–679, 1995.
- Hyman, J. M. and Nicolaenko, B.: The Kuramoto-Sivashinsky equation: a bridge between PDE's and dynamical systems, *Physica D: Nonlinear Phenomena*, 18, 113–126, 1986.
- Jaeger, H.: The “echo state” approach to analysing and training recurrent neural networks-with an erratum note, Bonn, Germany: German  
515 National Research Center for Information Technology GMD Technical Report, 148, 13, 2001.
- Jordan, M. I. and Mitchell, T. M.: Machine learning: Trends, perspectives, and prospects, *Science*, 349, 255–260, 2015.
- Korabel, N. and Barkai, E.: Pesin-type identity for intermittent dynamics with a zero Lyapunov exponent, *Physical review letters*, 102, 050 601, 2009.
- Krasnopolsky, V. M. and Fox-Rabinovitz, M. S.: Complex hybrid models combining deterministic and machine learning components for  
520 numerical climate modeling and weather prediction, *Neural Networks*, 19, 122–134, 2006.
- Kwasniok, F.: The reduction of complex dynamical systems using principal interaction patterns, *Physica D: Nonlinear Phenomena*, 92, 28–60, 1996.
- LeCun, Y., Bengio, Y., and Hinton, G.: Deep learning, *nature*, 521, 436, 2015.
- Li, D., Han, M., and Wang, J.: Chaotic time series prediction based on a novel robust echo state network, *IEEE Transactions on Neural  
525 Networks and Learning Systems*, 23, 787–799, 2012.
- Liu, J. N., Hu, Y., He, Y., Chan, P. W., and Lai, L.: Deep neural network modeling for big data weather forecasting, in: *Information Granularity, Big Data, and Computational Intelligence*, pp. 389–408, Springer, 2015.
- Lorenz, E. N.: Deterministic nonperiodic flow, *Journal of the atmospheric sciences*, 20, 130–141, 1963.
- Lorenz, E. N.: Predictability: A problem partly solved, in: *Proc. Seminar on predictability*, vol. 1, 1996.
- 530 Manneville, P.: Intermittency, self-similarity and  $1/f$  spectrum in dissipative dynamical systems, *Journal de Physique*, 41, 1235–1243, 1980.
- Moore, G., Renfrew, I. A., and Pickart, R. S.: Multidecadal mobility of the North Atlantic oscillation, *Journal of Climate*, 26, 2453–2466, 2013.
- Paladin, G. and Vulpiani, A.: Degrees of freedom of turbulence, *Physical Review A*, 35, 1971, 1987.
- Panichi, F. and Turchetti, G.: Lyapunov and reversibility errors for Hamiltonian flows, *Chaos, Solitons & Fractals*, 112, 83–91, 2018.
- 535 Pathak, J., Lu, Z., Hunt, B. R., Girvan, M., and Ott, E.: Using machine learning to replicate chaotic attractors and calculate Lyapunov exponents from data, *Chaos: An Interdisciplinary Journal of Nonlinear Science*, 27, 121 102, 2017.
- Pathak, J., Hunt, B., Girvan, M., Lu, Z., and Ott, E.: Model-free prediction of large spatiotemporally chaotic systems from data: A reservoir computing approach, *Physical review letters*, 120, 024 102, 2018.
- Quinlan, C., Babin, B., Carr, J., and Griffin, M.: *Business research methods*, South Western Cengage, 2019.
- 540 Rasp, S., Pritchard, M. S., and Gentine, P.: Deep learning to represent subgrid processes in climate models, *Proceedings of the National Academy of Sciences*, 115, 9684–9689, 2018.
- Saha, S., Moorthi, S., Wu, X., Wang, J., Nadiga, S., Tripp, P., Behringer, D., Hou, Y.-T., Chuang, H.-y., Iredell, M., et al.: The NCEP climate forecast system version 2, *Journal of Climate*, 27, 2185–2208, 2014.



- Scher, S.: Toward Data-Driven Weather and Climate Forecasting: Approximating a Simple General Circulation Model With Deep Learning, 545 Geophysical Research Letters, 45, 12–616, 2018.
- Scher, S. and Messori, G.: Predicting weather forecast uncertainty with machine learning, Quarterly Journal of the Royal Meteorological Society, 144, 2830–2841, 2018.
- Scher, S. and Messori, G.: Weather and climate forecasting with neural networks: using general circulation models (GCMs) with different complexity as a study ground, Geoscientific Model Development, 12, 2797–2809, 2019.
- 550 Schertzer, D., Lovejoy, S., Schmitt, F., Chigirinskaya, Y., and Marsan, D.: Multifractal cascade dynamics and turbulent intermittency, Fractals, 5, 427–471, 1997.
- Schneider, T., Lan, S., Stuart, A., and Teixeira, J.: Earth system modeling 2.0: A blueprint for models that learn from observations and targeted high-resolution simulations, Geophysical Research Letters, 44, 12–396, 2017.
- Shi, X., Gao, Z., Lausen, L., Wang, H., Yeung, D.-Y., Wong, W.-k., and Woo, W.-c.: Deep learning for precipitation nowcasting: A benchmark 555 and a new model, in: Advances in neural information processing systems, pp. 5617–5627, 2017.
- Shi, Z. and Han, M.: Support vector echo-state machine for chaotic time-series prediction, IEEE Transactions on Neural Networks, 18, 359–372, 2007.
- Sprenger, M., Schemm, S., Oechslin, R., and Jenkner, J.: Nowcasting foehn wind events using the adaboost machine learning algorithm, Weather and Forecasting, 32, 1079–1099, 2017.
- 560 Wang, J.-X., Wu, J.-L., and Xiao, H.: Physics-informed machine learning approach for reconstructing Reynolds stress modeling discrepancies based on DNS data, Physical Review Fluids, 2, 034 603, 2017.
- Weeks, E. R., Tian, Y., Urbach, J., Ide, K., Swinney, H. L., and Ghil, M.: Transitions between blocked and zonal flows in a rotating annulus with topography, Science, 278, 1598–1601, 1997.
- Weisheimer, A., Corti, S., Palmer, T., and Vitart, F.: Addressing model error through atmospheric stochastic physical parametrizations: 565 impact on the coupled ECMWF seasonal forecasting system, Philosophical Transactions of the Royal Society A: Mathematical, Physical and Engineering Sciences, 372, 20130 290, 2014.
- Wold, S., Esbensen, K., and Geladi, P.: Principal component analysis, Chemometrics and intelligent laboratory systems, 2, 37–52, 1987.
- Wolf, A., Swift, J. B., Swinney, H. L., and Vastano, J. A.: Determining Lyapunov exponents from a time series, Physica D: Nonlinear Phenomena, 16, 285–317, 1985.
- 570 Wu, J.-L., Xiao, H., and Paterson, E.: Physics-informed machine learning approach for augmenting turbulence models: A comprehensive framework, Physical Review Fluids, 3, 074 602, 2018.
- Xingjian, S., Chen, Z., Wang, H., Yeung, D.-Y., Wong, W.-K., and Woo, W.-c.: Convolutional LSTM network: A machine learning approach for precipitation nowcasting, in: Advances in neural information processing systems, pp. 802–810, 2015.
- Xu, M., Han, M., Qiu, T., and Lin, H.: Hybrid regularized echo state network for multivariate chaotic time series prediction, IEEE transactions 575 on cybernetics, 49, 2305–2315, 2018.
- Young, L.-S.: What are SRB measures, and which dynamical systems have them?, Journal of Statistical Physics, 108, 733–754, 2002.