

# Analyzing Echo-state Networks Using Fractal Dimension

Norbert Michael Mayer

*Department of Electrical Engineering and  
Advanced Institute of Manufacturing with High-tech Innovations  
National Chung Cheng University  
Taiwan  
nmmayer@gmail.com*

Oliver Obst

*Centre for Research in Mathematics and Data Science  
Western Sydney University  
New South Wales  
Australia  
O.Obst@westernsydney.edu.au*

**Abstract**—This work joins aspects of reservoir optimization, information-theoretic optimal encoding, and at its center fractal analysis. We build on the observation that, due to the recursive nature of recurrent neural networks, input sequences appear as fractal patterns in their hidden state representation. These patterns have a fractal dimension that is lower than the number of units in the reservoir. We show potential usage of this fractal dimension with regard to optimization of recurrent neural network initialization. We connect the idea of “ideal” reservoirs to lossless optimal encoding using arithmetic encoders. Our investigation suggests that the fractal dimension of the mapping from input to hidden state shall be close to the number of units in the network. This connection between fractal dimension and network connectivity is an interesting new direction for recurrent neural network initialization and reservoir computing.

**Index Terms**—Reservoir Computing, Echo-state Networks, Recurrent Neural Networks, Fractals, Arithmetic Encoding.

## I. INTRODUCTION

ONE of the key problems for Recurrent Neural Network (RNN) initialization, and in particular for Reservoir Computing (RC) methods like Echo State Networks (ESNs)[1], [2], [3] is that it is still unclear what kind of connectivity results in the best performance. For recurrent neural networks, some heuristics have shown much better performance than others: Dependent on the task, ESNs with orthogonal recurrent connectivity matrices that scale with the size of the network have shown better performance than other strategies [4], [5], [6]. An initialization method that can create a good recurrent weight matrix for a given task would be a compromise between an arbitrarily chosen connectivity matrix, and a connectivity matrix found by more expensive optimization, e.g., by backpropagation through time. Insights that lead to better connectivity matrices are useful for “standard” RNNs, where they are further trained, as well as for RC approaches, where they are used without further training.

With this in mind, the purpose of our work is to investigate how the input statistics is mapped to the ESNs reservoir. It appears self-evident that the performance of the network

relates to an (ideally) optimal usage of the memory capacity of the recurrent layer neurons. Of the two related questions about the memory capacity (1) “How much memory does a specific task require”, and (2) “How much memory does a network of a given size provide”, we are addressing the latter. It should be clear that some tasks require more memory capacity than others, and that, beyond a certain point, simply increasing memory will not contribute to a better performance. We are interested in making the best possible use of a given network, but the question of how large that network should be, for a given task, is outside the scope of this paper.

In order to approach this optimal capacity of networks with a given size, we follow the concepts of lossless memory compression. Our results reveal that the recursive nature of the update function in RNNs, and in ESNs in particular, give rise to the occurrence of fractals with regard to certain features of the reservoir in a natural way. Suggesting the fractal dimension as a measure for reservoir quality, we discuss new directions for how to improve initialization and performance of ESNs.

To provide better insights into properties of optimal coding recurrent neural networks, we restrict our investigation to use a discrete set of two input values that drive a network, which incidentally also allows a very natural application of (Shannon) information theory.

## II. MATHEMATICAL BACKGROUND AND DEFINITIONS

In this section, we briefly revisit the ESNs formalism, and some of the background for fractal dimensions and lossless compression.

### A. Echo State Networks

An ESN is described by the following equation:

$$x_{t+1} = \sigma(\alpha W x_t + \beta w_{in} u_t), \quad (1)$$

where  $\sigma$  is a (non-linear) transfer function for which Lipschitz continuity is required,  $u_t$  is an external input to the network, i.e., the stimulus. The  $m$  recurrent neurons are represented by  $x_t$ , an  $m$  dimensional vector. We set appropriate norms on both matrices  $\|W\|_2 = \|w_{in}\|_2 = 1$ . The scalar values of  $\alpha$  and  $\beta$  can be used for tuning the network and their quantitative analysis.

This work was financially supported by the Advanced Institute of Manufacturing with High-tech Innovations (AIM-HI) from The Featured Areas Research Center Program within the framework of the Higher Education Sprout Project by the Ministry of Education (MOE) in Taiwan.

To prevent divergence, ESNs should satisfy the so-called echo state property (ESP) [7], [8]. Heuristics show that, for some tasks, performance of ESNs are best near the limit of what is permissible according to the ESP. Likewise, several theoretical works point in this direction [9], [6], [10].

Input values  $u_t$  form an input stream (for example from sensory input, where the ESN forms or is embedded into a larger system). These values usually follow certain restrictions and form a statistics, which can be described in a likelihood relation,

$$p(u_t) \equiv p(u_t | u_{t-1}, u_{t-2} \dots).$$

One possible way to describe such statistics is to consider the set  $\mathbb{U}$ , where each possible time series  $u_{(-\infty, t]}$  is in  $\mathbb{U}$  which may be combined with the relative probability of its occurrence  $p(u_{(-\infty, t]})$ , in order to make the description of the statistics mathematically accurate.

Information propagation in RNNs can be viewed as a recursive process. We can describe an internal state of the network as  $x_t = f(x_0, u_{[0, t]})$ . Due to the ESP, it is also that the state  $x_t = f(u_{(-\infty, t]})$ , i.e., the current internal state is independent of the initial state, if the network has received sufficiently long input sequences. Viewed differently, the state vector  $x_t$  is set into a state space  $\mathbb{X}$ . Dimensionality of this space is equal to the number of hidden layer neurons. With a bounded transition function like  $\sigma = \tanh$ , output values of recurrent neurons are restricted to  $[-1, 1]$ , and thus  $\mathbb{X}$  forms a hypercube.

We restrict the input  $u$  to the RNN so that  $u_t \in \{-1, 1\}$ . In this case, the numerical results can be used to visualize interesting characteristics. In the following, we focus on the fact that the representations of the entire set  $\mathbb{U}$  – within one experimental setting (i.e., with certain statistics) onto the internal state-set  $\mathbb{X}$  – reveals fractals as overall phenomenon. In general, ESNs will usually have 100 or more units. For visualization, we confined the number of neurons to two neurons. The representation of  $\mathbb{U}$  onto the internal state neurons reveals the fractal nature of the resulting point clouds for the naked eye (cf. Fig. 2).

### B. Fractal dimension and lossless compression

Mathematician Benoit Mandelbrot in 1975 firstly used the term *fractal* (“broken”, or “fractured”) to apply the theoretical concept of fractal dimensions (FDs) to geometric patterns in nature. Generally speaking, fractals are wrinkled objects that, to some extent, escape more conventional measures such as length and area, and are better distinguished by their FD. The geometry of fractals comes up in many mathematical models for different objects in nature – such as mountains, coastal lines and clouds [11], [12]. One way to think about FD is based on the concept of self-similarity, i.e., the relationship between the number of identically shaped smaller objects the original object can be divided into, and their masses. For perfect self-similarity, this relationships can be described as

$$d_f = \lim_{r \rightarrow 0} \frac{\log(N_r)}{\log \frac{1}{r}}$$

where  $N_r$  is the least number of distinct copies of the original object  $A$  at scale  $r$ . The union of  $N_r$  distinct copies must cover  $A$  completely. In general, fractals are not necessarily (perfectly) self-similar, but the FD is a useful attribute of fractals in nature and some fields of natural sciences, in addition to their properties described by Euclidean geometry.

Different methods have been proposed to calculate and estimate the FD. Work in [13] classifies the methods of evaluating FDs into three main groups: (1) box-counting methods, (2) variance methods, and (3) spectral methods. In particular box-counting methods became very popular for their simplicity and auto-computability [12], and are applied in various fields. A number of different box-counting methods has been brought forward [14], [15], [16], [17].

### III. RNNs AND THEIR FRACTAL DIMENSION

As we investigate the mapping of the complete input statistics onto the reservoir manifold, one can see that representations of input histories emerge as fractal sets for many parameter settings of the connectivity matrix. Essentially, these fractals are an indicator of the ‘coarseness’ of the occupied states in the underlying reservoir. Here, one can understand coarseness as a scale-invariant feature. In combination with the total number of states that the reservoir can assume, one can calculate the memory capacity of the reservoir for a given input statistics. The resulting fractal dimension are key components to the better information-theoretic understanding of reservoir computing and RNN initialization. We discuss connections between optimal encoding in ESNs and arithmetic encoders, Shannon information in reservoirs, and implications for other types of networks.

Though popular due to their efficient training, ESNs also have some basic limitations [18]. A first limitation are potentially ill-conditioned solutions, resulting in models that might diverge from the real system that is supposed to be represented. This, and other limitations have also been observed in other works, e.g., [19], [20], [21], [22], and improvements have been suggested in [23], [24], [25]. A better use of the available memory would alleviate some of these limitations.

One basic insight at the start of our research was that an optimal representation of states in a reservoir has an analogue in the representation pattern of arithmetic encoders [26]. Arithmetic encoders have capabilities comparable to Huffman codes: they encode a sequence of symbols in to a real valued number, that can then be transmitted with a given accuracy. We use them as a guideline that can lead us to an optimal representation of the input for a given statistics on the hyperspace of the recurrent layer neurons’ activities. It has been outlined in [27] that arithmetic encoders can trivially be transformed into a recurrent filter, where the output of the modified arithmetic encoder fulfills the ESP, i.e., in this case the arithmetic encoder can be interpreted as a reservoir of exactly one unit in the sense of an ESN.

The fractal dimension of a network, as calculated by, e.g., box-counting, is smaller or equal to the number of neurons in the recurrent layer. Before we illustrate the effect of measuring

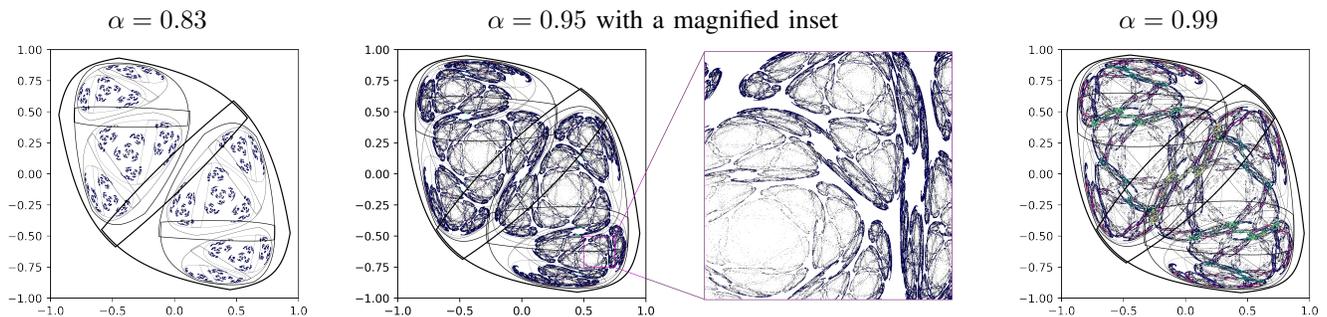


Figure 1: State space representation of an ESN’s recurrent layer with 2 neurons. The closed black curves depict the transformed state space limits.

FD, we demonstrate the reasons why the FD can be a relevant measure for ESNs:

The graphs in Figs. 1 and 2 show random state space representations of a recurrent layer with two neurons. The tanh transfer function restricts the possible values of states to the interval between  $-1$  and  $1$ . For two neurons this defines a square within which all permissible states are located. One can analyze how this square of permissible values is transferred into and transformed for the next iteration. We get two different transformed curves, for the two possible inputs of  $1$  and  $-1$ , at each iteration, cf. Fig. 2a. We also depict representations of different samples of  $\mathbb{U}$ , where the color of the points depends on whether the last value that the respective time series has received, was,  $1$  or  $-1$ . Fig. 2b shows the same results of the same network after five iterations.

Although the transformed input spaces overlap in the depicted example we do not see areas where the point clouds of different colors inundate each other.

The fractal patterns of Fig. 1 are a direct result of the recursive definition of RNNs. Other known examples of recipes for fractal sets, for example the well known Barnsley fern [28], which is constructed by a recursive rule. In (deep) feedforward networks, FractalNet [29] creates the architecture of the network based on ideas of self-similarity, but with a different motivation of influencing path lengths for subsequent training, without further investigation of the resulting coding.

Returning to RNNs and reservoirs, it appears obvious that with a broader coverage of  $\mathbb{U}$  within a recurrent layer, information processing capabilities will improve. It might be also worthwhile to consider the maximal information content that is available within the network:

$$I_{\max} = m \cdot m_b,$$

where  $m$  is the number of neurons and  $m_b$  is the logarithm of the accuracy where each of the neurons are represented. In other words,  $I_{\max}$  is the total number of bits available within the reservoir layer, or in  $\mathbb{X}$ , equivalently.

For example, in the case when values are represented by one double-precision IEEE 754 standard floating-point number, i.e., the 64 bits reserved for each number are distributed in the following manner: 1 bit is for the sign, 11 bits for the exponent,

and 52 bits for the significant bits (mantissa). In this example,  $m_b = 53$ : the number of bits of the mantissa plus the sign bit. For sake of simplicity, we assume here and for the following considerations that the smallest difference (“accuracy”) with which each state can be represented is  $2^{-m_b}$ . Thus, we neglect the fact when using floating point representations a number that is very near to zero is significantly more accurately represented than numbers with a larger absolute values<sup>1</sup>.

With this as a starting point, we can now imagine the reservoir as a  $m$ -dimensional hypercube where each side is again subdivided in  $2^{m_b}$  segments, forming a total  $2^{I_{\max}}$  sub-hyper-cubes. The ideal usage of the reservoir is achieved if all of those sub-hyper-cubes are occupied with a non-empty subset of  $\mathbb{U}$  and are not mapped to the same hypercube for as long as possible lasting input histories and the probability for the occurrence of each of the subsets is equal for all hyper-cubes.

Usually, and different from the ideal case, less than the maximal number of those hyper-sub-cubes are actually occupied. With an estimation of FD, the number of occupied boxes is

$$2^{d_f \cdot m_b}.$$

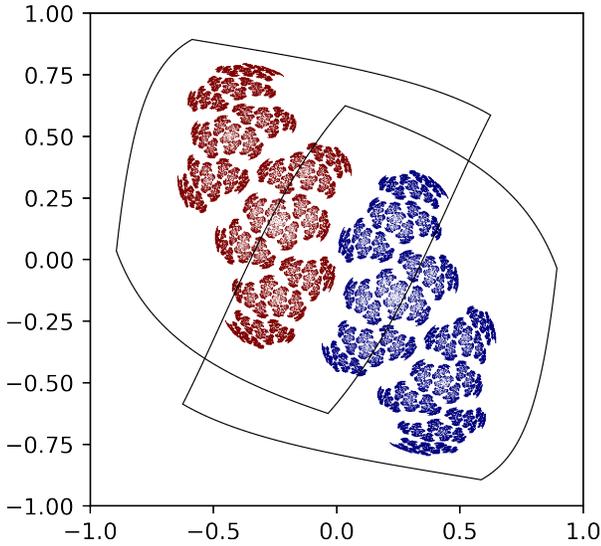
A plausible estimate of the mutual information between the input histories and the internal states thus is

$$I(\mathbb{U}; \mathbb{X}) \leq d_f \times m_b,$$

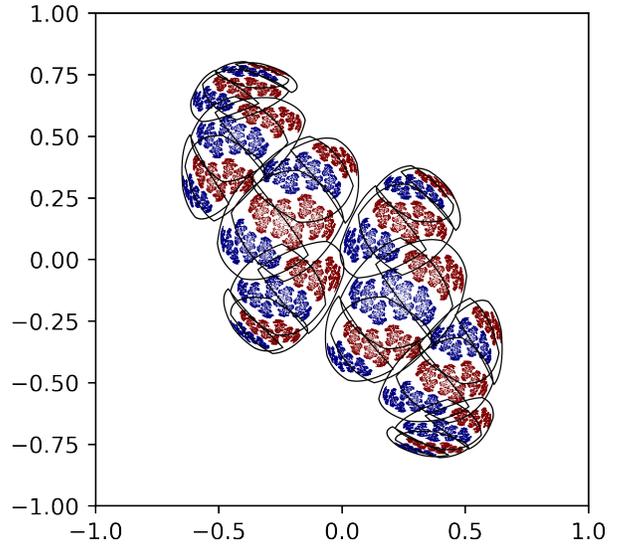
where  $d_f$  is the fractal dimension.

This suggests we can improve the performance of ESNs by maximizing the FD, by varying the parameters  $\alpha$  and  $\beta$  of the recurrent connectivity. Practically, we need to lift the fractal dimension of the data representation close to the number of neurons such that possible values of the recurrent layer in the ESN represent at least one input history. This achieves that the data of the recurrent layer represents a compressed image of an input sequence, for as much as it is possible for the history of  $u_t$ . Consequently, the representation of the recurrent layer would become something like a Huffman coded representation

<sup>1</sup>The concept may be extended to other forms of reservoirs in the sense of RC, in physical systems one may consider the accuracy of reliable readout down to thermal fluctuations, in case of super-cooling systems down to quantum effects, etc.



(a) After 1 iteration



(b) After 5 iterations

Figure 2: State space representation of an ESN's recurrent layer with 2 neurons. The closed black curves depict the transformed state space limits.

of the input history, as much as it is possible according to the entropy resulting from the statistics of  $\mathbb{U}$ .

The basic idea of box counting methods to estimate FD of an object is to count the number of boxes  $N(\epsilon_i)$  occupied by that object, with boxes at scales  $\epsilon_i$ . Using a number of different scales  $\epsilon_i$ , and plotted on a log-log scale, the slope of that resulting line yields the FD ("Richardson effect").

$$\log N(\epsilon_i) = d_f \cdot \log \epsilon + b.$$

For perfect fractals, this slope is estimated in the limit, but for practical estimation, objects are assumed to be fractal when the slope remains approximately constant over a range of scales. We use a standard box-counting method [12] directly for points which we have obtained from the hidden state output.

1) *Arithmetic Encoding*: Arithmetic encoding describes the idea to convert a sequence of symbols into a real-valued number in an iterative procedure. For the encoding, each symbol is assigned an interval with a size proportional to the probability of the occurrence of that particular symbol. In this way, recursively, a sequence of symbols can be encoded until the encoder hits the limit of accuracy of the underlying numerical representation. Here, arithmetic encoding is to map a finite sequence of a discrete set of symbols into a real value  $x_t$  in the range  $[0, 1]$ . This real value that represents the sequence can then be transmitted to a receiver. The number of digits transmitted has to be large enough to make it possible to identify the initial sequence.

Originally, arithmetic encoders have been developed at IBM [26] for data compression, competing with Huffman coding. Interesting applications of arithmetic encoders can be

found in, e.g., [30], [31], [32], [33]. Recent investigations have revealed an analogy between arithmetic encoders and optimal ESNs [6], [34], [9], [35], [5], [36]. Making use of this idea, the connectivity in an RNN would be arranged such that for the given input statistics in  $\mathbb{U}$ , the network performs a data compression so that as many as possible past input values are represented in the reservoir.

We rephrase the definition of arithmetic encoders as a recursive filter, and as a result its conformity with the ESP becomes obvious<sup>2</sup>. As a pre-requisite, it is necessary to estimate or quantify the probability  $p(\kappa_i)$  of a discrete, finite set of symbols  $\kappa_i$  with  $i \leq I$ .

Now we can define a function  $0 \leq g(i) < 1$ :

$$g(\kappa_i) = \sum_{j < i} p(\kappa_j).$$

$g(i)$  is a step-wise monotonously increasing function which can be reverted to

$$\kappa_i = \tilde{g}(x).$$

To encode a finite sequence  $u_t$  from  $0 \leq t \leq T$  of those symbols, we choose an initial value  $x_0 = a$ , where  $a$  can be an arbitrary real number in the range  $[0, 1]$ .

We then can calculate recursively for all  $0 < i \leq T$

$$x_{i+1} = p(u_i) \cdot x_i + g(u_i).$$

$x_{T+1}$  can then be transmitted. However, this value can only be transmitted using a finite accuracy. The required accuracy

<sup>2</sup>For this purpose the impact of the time axis has to be inverted from the original approach [26].

can be derived from considering the initializing setting of  $x_0 = a$ . One can then define

$$\begin{aligned} x^- &= x_{T+1}(u_T, \dots, u_0, x_0 = 0) \\ x^+ &= x_{T+1}(u_T, \dots, u_0, x_0 = 1) \end{aligned}$$

The necessary number of digits that have to be transmitted to receive the message unambiguously has to be chosen such that the receiver can know  $x \in [x^-, x^+]$  but also  $x \notin [0, x^-]$  and  $x \notin [x^+, 1]$ . The necessary digits then form an optimally short code, equivalent to a Huffman code.

At the same time it is possible to estimate an upper limit to

$$x^+ - x^- \geq \left( \max_i p(\kappa_i) \right)^T.$$

If we assume our time series to have non-zero entropy, we know  $\max_i p(\kappa_i) < 1$ . Thus, in the limit for  $T \rightarrow \infty$ , this value is zero, i.e., the arithmetic encoder is uniformly state contracting and, as a result, conform with the ESP. Arithmetic encoding is a one unit reservoir ESN.

Using the arithmetic encoder, with correctly chosen distribution, and if  $p(\kappa_i)$  is i.i.d. with regard to the position of the symbol within the sequence, the representation of possible symbol sequences is dense on  $x_T$ , for sufficient large  $T$ . Thus, the fractal dimension of the mapping of  $\mathbb{U}$  on  $\mathbb{X}$  is 1.0.

The fractal dimension will be lower than 1.0 if we were using imperfect estimates for  $p(\kappa_i)$ . As an example, if we were to assume a code statistics of three identical and independent distributed symbols A, B, C with equal probability  $\frac{1}{3}$  but in reality the sequences contained only the symbols A and C with equal probability  $\frac{1}{2}$ , one can organize the structure of the resulting gaps in the corresponding representation as a Cantor set. In this case, we achieve a fractal dimension of  $\frac{\log 2}{\log 3} \approx 0.63$ .

#### IV. RESULTS

In this section we show results of determining the FD and subsequent classification using a support vector machine (SVM). The SVM is used to show whether the input sequence mapping to the hidden states of the network is separable.

The graphs in Fig. 3 are result of using the FD obtained by applying the box-counting method, where  $x$ -axis and  $y$ -axis represents the first neuron and second neuron, respectively. Fig. 3a is the state space visualization generated for  $\alpha = 1$  and  $\beta = 0.45$ , a promising result. Fig. 3b is an example with poor results for  $\alpha = 1$  and  $\beta = 0.8$ . For low  $\beta$  values, there will not be any gaps between the points, and the goal is to avoid an overlap between points. In this example,  $\beta < 0.45$  results in overlapping and  $\beta > 0.45$  results in large gaps between the points as shown in Fig. 3b. Using Richardson effect's box-counting supports this result, see Fig. 3c. For  $\beta = 0.45$ , we achieve an almost optimal dimensionality of close to 2, equal to the number of neurons. Fig. 3d shows the unsatisfactory result for  $\beta = 0.8$ . with a dimensionality of 1.2. In Figs. 3c and 3d, the  $x$ - and  $y$ -axes represent the  $\log \epsilon$  and  $\log N_s$  values, respectively.

Fig. 4 shows a 3d-visualization of FDs on multiple scales of  $\alpha$  and  $\beta$ . The decrease in  $\beta$  is clearly visible as the dimension value increases (the values range from 2.0 to 0.2 from left to right). As  $\alpha$  increases, the fractal dimension increases (or vice-versa). That is to say, the measured dimension strongly depends on tuning the  $\alpha$  and  $\beta$  values.

We also observe a trend in the number of support vectors while alpha and beta changes: The complexity of the hyperplane as generated by the SVM model is approximately proportional to the number of support vectors, because each support vector will define one or more twists or turns.

Figure 5 shows results of a classification with a SVM. We plot the number of support vectors for different  $\alpha$ -values ( $x$ -axis), with different beta values (indicated by color). When  $\alpha$  increases, the number of support vectors has a tendency to increase.

Along with the result of Fig. 4, we aim to use the relationship between box-counting and the SVM to analyze ESN. When  $\alpha = 1$  and  $\beta = 0.2$  or  $\beta = 0.4$ , we achieve to lift the FD close to the number of neurons. The same parameters have been used for experiments in both figures.

#### A. Simulation details

For simulations in Fig. 2 we used  $W = \mathcal{O}(2.0)$ , where  $\mathcal{O}$  is a  $2 \times 2$  rotation matrix around the angle 0.5 rad;  $w_{in} = (\sqrt{0.5}, -\sqrt{0.5})$ ,  $\alpha = 0.8$  and  $\beta = 0.5$ . The network has been stimulated by all possible combinations of 1 and  $-1$  up to a limit length emulating iid. random sequences.

For Fig. 3 and 4, we used 100 input sequences of 1 and -1, each sequence with about 1 million temporal steps. Hence, for each pair of  $\alpha$  and  $\beta$ , 100 million points were generated in total, using the network with weight matrix

$$W^{in} = \begin{bmatrix} 0.8436 & 0.7381 \end{bmatrix}, W = \begin{bmatrix} 0.0169 & 0.5711 \\ 1.2895 & 0.2509 \end{bmatrix},$$

Moreover, the interval of  $\alpha$  and  $\beta$  were set to

$$\begin{aligned} \alpha &: [0.4, 0.45, 0.5 \dots, 1] = \{X \in \mathbb{R} | x_n - x_{n-1} = 0.5\}; \\ \beta &: [0.15, 0.2, 0.25 \dots, 2] = \{X \in \mathbb{R} | x_n - x_{n-1} = 0.5\} \end{aligned}$$

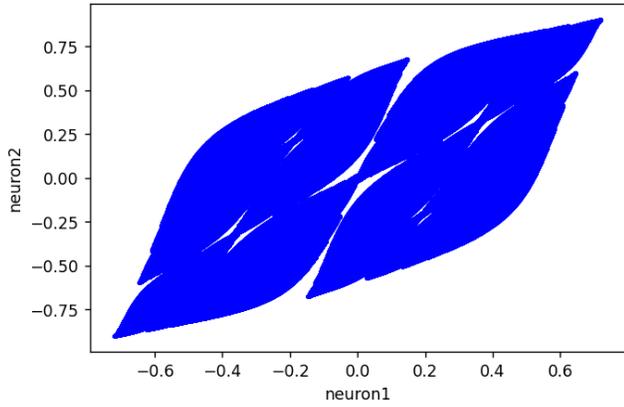
The slope and the 3D illustration in Fig. 3 and Fig. 4, respectively, were produced by the box-counting method with

$$\epsilon = \{300, 500, 600, 650, 700, 750, 800, 850, 900, 950, 1000\}$$

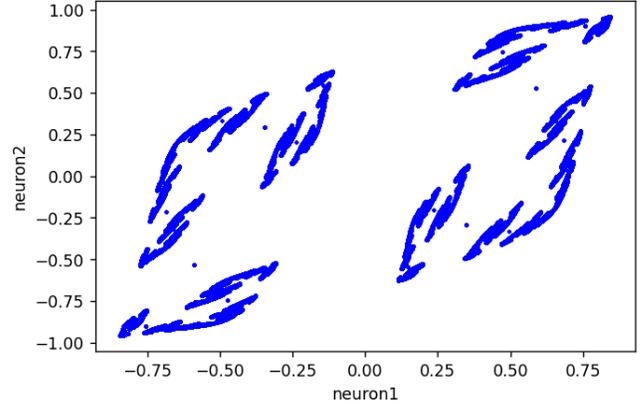
The simulation of Fig. 5 consists of 3000 training data generated from network states in SVM model. The weight matrix used for Fig. 5 was same as the one used for Fig. 3. The values  $\alpha$  and  $\beta$  are tuned between 0 and 1. The classifier to explore separation is a radial basis function kernel SVM, with the Gaussian kernel:

$$K(\mathbf{x}_a, \mathbf{x}_b) = \exp\left(-\frac{\|\mathbf{x}_a - \mathbf{x}_b\|^2}{2\sigma^2}\right)$$

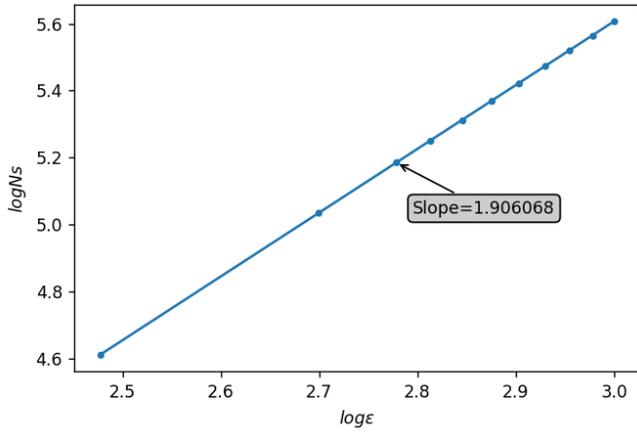
where  $\sigma$  is a tunable parameter. This parameter influences the "reach" of each single example: when  $\sigma$  is small, the SVM model can capture more complex data. We chose  $\sigma = 0.05$ ,



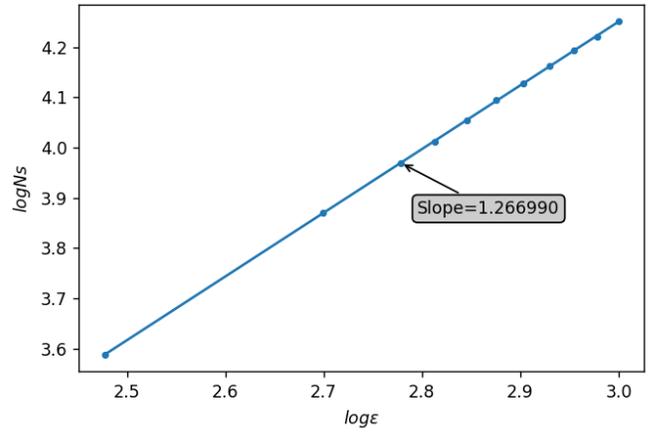
(a) The state space generated by the network for  $\beta = 0.45$ . Plot for  $10^8$  sample points on 2 neurons with  $\alpha = 1$  and  $\beta = 0.45$ .



(b) The state space generated by the network for  $\beta = 0.8$ . Plot for  $10^8$  sample points on 2 neurons with  $\alpha = 1$  and  $\beta = 0.8$ .



(c) Fractal dimension obtained by Richardson effect for  $\beta = 0.45$ .



(d) Fractal dimension obtained by Richardson effect for  $\beta = 0.8$ .

Figure 3: Box-counting results for 2 different beta values.

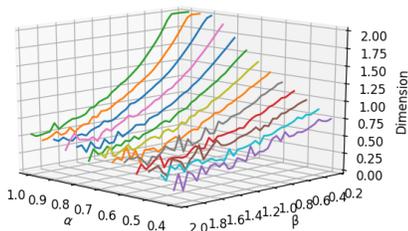


Figure 4: Fractal dimension on multiple scales of  $\alpha$  and  $\beta$ .

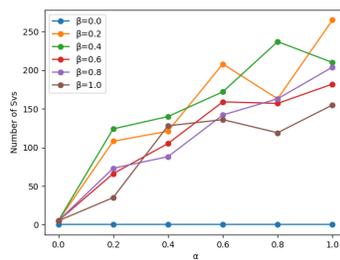


Figure 5: Number of support vectors for different  $\beta$  and fixed  $\alpha$ .

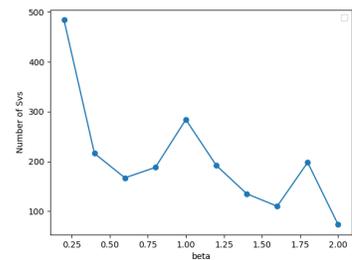


Figure 6: Number of support vectors for different  $\beta$  and  $\alpha = 1$ .

to have a good ability of classifying the overlapping samples, which means the shape of decision boundary may be more complicated. An additional hyper-parameter  $C$  trades off correctness and complexity of the model (essentially a regularization parameter). We set  $C = 2.0$ .

## V. DISCUSSION

On one hand, the arithmetic encoding approach has a strong relationship to ESNs with one neuron, and complies with the ESP. On the other hand, it is well established that arithmetic encoders are optimal encoders of an (input) sequence. In the ideal case, for correctly applied arithmetic encoders and appropriate input sequences, this will lead to a compact, dense state representation of input sequences, resulting in a FD of 1 (per neuron). At the same time overlapping representations have to be avoided. Heuristics hint towards the notion that both conditions are fulfilled exactly at the limit point where the FD reaches the reservoir size.

For input sequences composed from discrete sets, we see arithmetic encoders as a very valuable guideline to understand optimal representations in reservoirs, leading to new criteria for reservoir and RNN initialization. We suggest the following new concepts with regard to optimizing reservoirs:

- The FD of the mapping  $\mathbb{U} \rightarrow \mathbb{X}$  shall be near the number of neurons. This is an important requirement, since the requirement of the complete coverage of the layer of input neurons appears to be difficult in the case of the ESN formulation.
- The modified ESN shall work like a non-linear vector arithmetic encoder.
- Reservoirs with overlapping representations have indistinguishable states for readouts by means of regression or classification. We have to avoid those overlapping representations of different input histories.

This result has implications for a previously formalized idea of power law forgetting [37] and the edge of chaos heuristics: for input sequences with non-zero entropy, this idea would lead to overlapping representations. Overlapping representations of several recently variant input time series have clearly an impact on the performance of all ESNs including those with hundreds of recurrent neurons. Since this phenomenon occurs at a lower bound than the heuristically found edge of chaos it seems plausible that future discussions might lead to a new bound for the ESP, that is the ESP of non-overlapping representations with regard to a given input statistics.

## VI. CONCLUSION

The fractal structure of representations in RNNs is to our best knowledge an unnoticed aspect with implications for reservoir computing. It opens a new field of research which can bring together aspects of fractal analysis, information theory, and representation learning. Our results have been combined with the insight that overlapping representations of variant input time series better shall be avoided. This leads us to still preliminary proposed ESP II, for which a less heuristic and

better analytically founded formulation is still ongoing work in progress.

## VII. ACKNOWLEDGMENT

We thank Bo Ruei Jiang, Gorri Anil Kumar, and Ming Jie Li for their help. We thank Advanced Institute of Manufacturing with High-tech Innovations (AIM-HI) for their financial support.

## REFERENCES

- [1] H. Jaeger, "The "echo state" approach to analysing and training recurrent neural networks – with an erratum note," German National Research Centre for Information Technology, GMD report 148, 2001.
- [2] H. Jaeger, "Short term memory in echo state networks," German National Research Centre for Information Technology, GMD report 152, 2002.
- [3] H. Jaeger, "Tutorial on training recurrent neural networks, covering BPPT, RTRL, EKF and the "echo state network" approach," German National Research Centre for Information Technology, GMD report 159, 2002.
- [4] O. L. White, D. D. Lee, and H. Sompolinsky, "Short-term memory in orthogonal neural networks," *Phys. Rev. Lett.*, vol. 92, p. 148102, Apr 2004.
- [5] N. M. Mayer and Y.-H. Yu, "Orthogonal echo state networks and stochastic evaluations of likelihoods," *Cognitive Computation*, vol. 9, no. 3, pp. 379–390, 2017.
- [6] J. Boedecker, O. Obst, J. T. Lizier, N. M. Mayer, and M. Asada, "Information processing in echo state networks at the edge of chaos," *Theory in Biosciences*, vol. 131, no. 3, pp. 205–213, 2012.
- [7] I. B. Yildiz, H. Jaeger, and S. J. Kiebel, "Re-visiting the echo state property," *Neural networks*, vol. 35, pp. 1–9, 2012.
- [8] M. Buehner and P. Young, "A tighter bound for the echo state property," *IEEE Transactions on Neural Networks*, vol. 17, no. 3, pp. 820–824, 2006.
- [9] N. M. Mayer and M. Browne, "Echo state networks and self-prediction," in *International workshop on biologically inspired approaches to advanced information technology*. Springer, 2004, pp. 40–48.
- [10] M. A. Hajnal and A. Lőrincz, "Critical echo state networks," in *International Conference on Artificial Neural Networks*. Springer, 2006, pp. 658–667.
- [11] A. P. Pentland, "Fractal-based description of natural scenes," *IEEE transactions on pattern analysis and machine intelligence*, no. 6, pp. 661–674, 1984.
- [12] H.-O. Peitgen, H. Jürgens, and D. Saupe, *Chaos and fractals: new frontiers of science*. Springer Science & Business Media, 2006.
- [13] A. S. Balghonaim and J. M. Keller, "A maximum likelihood estimate for two-variable fractal surface," *IEEE Transactions on Image Processing*, vol. 7, no. 12, pp. 1746–1753, 1998.
- [14] J. Gagnepain and C. Roques-Carmes, "Fractal approach to two-dimensional and three-dimensional surface roughness," *wear*, vol. 109, no. 1–4, pp. 119–126, 1986.
- [15] N. Sarkar and B. B. Chaudhuri, "An efficient differential box-counting approach to compute fractal dimension of image," *IEEE Transactions on systems, man, and cybernetics*, vol. 24, no. 1, pp. 115–120, 1994.
- [16] W.-S. Chen, S.-Y. Yuan, and C.-M. Hsieh, "Two algorithms to estimate fractal dimension of gray-level images," *Optical Engineering*, vol. 42, no. 8, pp. 2452–2465, 2003.
- [17] S. Novianto, Y. Suzuki, and J. Maeda, "Near optimum estimation of local fractal dimension for image segmentation," *Pattern Recognition Letters*, vol. 24, no. 1–3, pp. 365–374, 2003.
- [18] H. Jaeger, "Echo state network," *scholarpedia*, vol. 2, no. 9, p. 2330, 2007.
- [19] H. Jaeger, "Reservoir riddles: Suggestions for echo state network research," in *Proceedings. 2005 IEEE International Joint Conference on Neural Networks, 2005.*, vol. 3. IEEE, 2005, pp. 1460–1462.
- [20] Y. Xue, L. Yang, and S. Haykin, "Decoupled echo state networks with lateral inhibition," *Neural Networks*, vol. 20, no. 3, pp. 365–376, 2007.
- [21] J. J. Steil, "Online reservoir adaptation by intrinsic plasticity for backpropagation-decorrelation and echo state learning," *Neural Networks*, vol. 20, no. 3, pp. 353–364, 2007.

- [22] Q. Song, Z. Feng, and M. Lei, "Stable training method for echo state networks with output feedbacks," in *2010 International Conference on Networking, Sensing and Control (ICNSC)*. IEEE, 2010, pp. 159–164.
- [23] Y. Liu, J. Zhao, W. Wang, Y. Wu, and W. Chen, "Improved echo state network based on data-driven and its application to prediction of blast furnace gas output," *Acta Automatica Sinica*, vol. 35, no. 6, pp. 731–738, 2009.
- [24] O. Obst, J. Boedecker, and M. Asada, "Improving recurrent neural network performance using transfer entropy," in *ICONIP 2010: Neural Information Processing, Models and Applications*, 2010, pp. 193–200.
- [25] C. Sheng, J. Zhao, Y. Liu, and W. Wang, "Prediction for noisy nonlinear time series by echo state network based on dual estimation," *Neurocomputing*, vol. 82, pp. 186–195, 2012.
- [26] J. Rissanen and G. G. Langdon, "Arithmetic coding," *IBM Journal of research and development*, vol. 23, no. 2, pp. 149–162, 1979.
- [27] N. M. Mayer, "Arithmetic encoding, reservoir computing, criticality and biological implications," in *Bernstein Conference, Göttingen*, 2017, (Abstract). [Online]. Available: <https://abstracts.g-node.org/abstracts/bfcd6994-61b6-4b25-83bb-e3ce10217745>
- [28] M. F. Barnsley, *Fractals everywhere*. Academic press, 2014.
- [29] G. Larsson, M. Maire, and G. Shakhnarovich, "FractalNet: Ultra-deep neural networks without residuals," in *ICLR*, 2017.
- [30] N. Brady, F. Bossen, and N. Murphy, "Context-based arithmetic encoding of 2d shape sequences," in *Proceedings of International Conference on Image Processing*, vol. 1. IEEE, 1997, pp. 29–32.
- [31] N. Brady and F. Bossen, "Shape compression of moving objects using context-based arithmetic encoding," *Signal Processing: Image Communication*, vol. 15, no. 7-8, pp. 601–617, 2000.
- [32] M. Li, S. Gu, D. Zhang, and W. Zuo, "Efficient trimmed convolutional arithmetic encoding for lossless image compression," *arXiv preprint arXiv:1801.04662*, pp. 107–120, 2018.
- [33] M. Li, S. Gu, D. Zhang, and W. Zuo, "Enlarging context with low cost: Efficient arithmetic coding with trimmed convolution," *arXiv preprint arXiv:1801.04662*, 2018.
- [34] J. Boedecker, O. Obst, N. M. Mayer, and M. Asada, "Studies on reservoir initialization and dynamics shaping in echo state networks." in *ESANN*, 2009.
- [35] N. M. Mayer, "Echo state condition at the critical point," *Entropy*, vol. 19, no. 1, p. 3, 2017.
- [36] C. Weber, K. Masui, N. M. Mayer, J. Triesch, and M. Asada, "Reservoir computing for sensory prediction and classification in adaptive agents," *Machine Learning Research Progress: Nova publishers: Hauppauge, NY, USA*, 2008.
- [37] N. M. Mayer, "Input-anticipating critical reservoirs show power law forgetting of unexpected input events," *Neural Computation, MIT Press*, vol. 27, pp. 1102–1119, May 2015.