

Temporal Overdrive Recurrent Neural Network

Filippo Maria Bianchi¹, Michael Kampffmeyer¹, Enrico Maiorino², and Robert Jenssen¹

¹Machine Learning Group, Dept. of Physics and Technology, University of Tromsø, Norway

²Dept. of Information, Electronic and Communication Engineering, Sapienza University, Rome, Italy

Abstract—In this work we present a novel recurrent neural network architecture designed to model systems characterized by multiple characteristic timescales in their dynamics. The proposed network is composed by several recurrent groups of neurons that are trained to separately adapt to each timescale, in order to improve the system identification process. We test our framework on time series prediction tasks and we show some promising, preliminary results achieved on synthetic data. To evaluate the capabilities of our network, we compare the performance with several state-of-the-art recurrent architectures.

Keywords— Recurrent Neural Network; Multiple Timescales; Time Series Analysis.

I. INTRODUCTION

Time series analysis focuses on reconstructing the properties of a dynamical system from a sequence of values, produced from a noisy measurement process acting on the state space of the system. Correct modeling is necessary to understand, characterize the dynamics of the system and, consequently, to predict its future behavior. A Recurrent Neural Network (RNN) is an universal approximator of Lebesgue measurable dynamical systems and is a powerful tool for predicting time series [13]. In its basic formulation, the whole state of a RNN evolves according to precise timing. However, in order to model a system whose internal dynamic variables evolve at different timescales, a RNN is expected to learn both short and long-term dependencies. The capability of modeling multiple time scales is particularly important in the prediction of real-world time series, since they are generally characterized by multiple seasonalities [10].

In this paper we propose a novel RNN architecture, specifically designed for modeling a system characterized by heterogeneous dynamics operating at different timescales. We combine the idea of using bandpass neurons in a randomly connected recurrent layer, an

inherent feature of reservoir computing approaches [17], with the gradient-based optimization, adopted in classic RNN architectures. Contrarily to standard RNNs, where the internal state of the neurons is influenced by all the frequencies of the input signal, in our network specific groups of neurons update their states according to different frequencies in the input signal. Therefore, since the hidden neurons can operate at temporal resolutions that are slower than the fastest frequencies in the input, we named our network Temporal Overdrive RNN (TORNN), after the mechanism in vehicles engine that, under certain conditions, decouples the speed of the wheels from the revolutions of the engine.

TORNN is capable of modeling with high accuracy the timescales in the dynamics of the underlying system. In the proposed architecture, the hidden recurrent layer is partitioned in different groups, each one specialized in modeling dynamics at a single timescale. This is achieved by organizing the recurrent layer in groups of bandpass neurons with shared parameters, which determine the frequency band to operate on. During training, while the neuron connection weights are fixed to their initial random values, the bandpass parameters are learned via backpropagation. In this way, the recurrent layer dynamics are smoothed out by the bandpass filtering in order to adapt to the relevant timescale. The main hyperparameter required by TORNN is the number of neuron groups composing the recurrent layer, that has to be chosen by considering the expected number of separate timescales characterizing the analyzed process. This information can be easily retrieved through an analysis in the frequency domain and/or by *a priori* knowledge of the input data. The network also depends on few other structural hyperparameters, which are not critical as the final performance of the model is not particularly sensitive to their tuning. Since the pair of bandpass parameters are the only weights trained in each group of neurons, the total number of parameters to be learned are considerably less than in other RNN architectures, with a consequent simplification of the training procedure.

filippo.m.bianchi@uit.no
Corresponding author
michael.c.kampffmeyer@uit.no
enrico.maiorino@uniroma1.it
robert.jenssen@uit.no

In this work we focus on the prediction of real-valued time series, generated from a system whose internal dynamics operate at different timescales. In Sect. II we review previous works that dealt with the problem of modeling dependencies across different time scales in the data with RNNs. In Sect. III we provide the details of the proposed architecture. In Sect. IV we present some initial results on synthetic data obtained by the proposed architecture and other state-of-the-art alternatives. Finally, in Sect. V we draw our conclusions.

II. RELATED WORKS

Classic RNN architectures struggle to model different temporal resolutions, especially in presence of long-term relations that are poorly handled, due to the issue of the vanishing gradient [1]. Several solutions have been proposed in the past to deal with long-term dependencies [19] and an important milestone has been reached with the introduction of Long Short Term Memory (LSTM) [14]. This architecture has been shown to be robust to the vanishing gradient problem, thanks to its gated access mechanism to its neurons' state. This results in the LSTM networks' capability to handle long time lags between events and to automatically adapt to multiple time granularity. In particular, the LSTM neurons, referred to as "cells", are composed by an internal state that can be accessed through three gates: the input gate, that controls whether the current input value should replace the existing state; the forget gate, that controls the reset of the cell's state; the output gate, that determines whether the unit should output its current state. An additional component in several implementations of LSTM cells are the peephole connections, which allow gate layers to look directly at the cell state. A recent variant of the traditional LSTM architecture is the Gated Recurrent Unit (GRU) [5], which always exposes its internal state completely, without implementing mechanisms to control the degree of exposure. In GRU, both peephole connections and output activation functions are removed, while the input and the forget gate are coupled into an update gate. While LSTM and GRU excel in learning at the same time long and short time relationships, they are difficult to train due to an objective function landscape characterized by a high dimensionality and several saddle points [9]. Additionally, these networks depend on several hyperparameters, whose tuning is non-trivial, and the utility of their various computational components depends on the application at hand [12].

An initial attempt to model multiple dynamics and timescales was based on the idea that temporal relationships are structured hierarchically, hence the RNN

should be organized accordingly [8]. The resulting architecture managed to improve the modeling of slow-changing contexts. An analogous hierarchical organization has been implemented by stacking multiple recurrent layers [11]. In the same spirit, a more complex stacked architecture called Gated Feedback Recurrent Neural Networks has been proposed in Ref. [7]. In this architecture, the recurrent layers are connected by gated-feedback connections which allow them to operate at different timescales. By referring to the unfolding in time of the recurrent network, the states of consecutive time steps are fully connected and the strength of the connections is trained by gradient descent. The main shortcoming of these layered architectures is the high amount of parameters that must be learned in order to adapt the network to process the right time scales. This results in a long training time, with the possibility of overfitting the training data.

The issue of modeling multiple timescales has also been discussed within the framework of reservoir computing, a quite recent approach to temporal signal processing that leverages on a large, randomly connected and untrained recurrent layer called reservoir. The desired output is computed by a linear memory-less read-out, which maps the instantaneous states of the network and is usually trained by linear regression. The original reservoir computing architectures, known as Echo State Networks [15] and Liquid State Machines [21], are characterized by a fading memory that prevents to model slow periodicities that extend beyond the memory capacity of the reservoir. A solution proposed in Ref. [17] is to slow down the dynamics of the reservoir by using leaky integrator units as neurons. These units have individual state dynamics that can model different temporal granularity in the target task. Leaky integrators behave as lowpass filters, whose cutoff frequency is controlled by the leakage parameter. Precise timing phenomena emerge from the dynamics of a random connected network implementing these processing units, without the requirement of dedicated timing mechanism, such as clocks or tapped delay lines. A different kind of unit that implements a bandpass filter has been proposed to encourage the decoupling of the dynamics within a single reservoir and to create richer and more diverse echoes, which increase the processing power of the reservoir [25]. Improved results have been achieved by means of neurons implementing more advanced digital bandpass filters, with particular frequency-domain characteristics [26]. Differently from other strategies where the optimal timescales for the task at hand are learned from data, reservoir computing follows a "brute force" approach. In fact, not only a large, sparse and randomly

connected reservoir is generated in order to provide rich and heterogeneous dynamics, but also a high amount of filters are initialized with random cutoff frequencies, to encode a wide range of timescales in the recurrent layer. By providing filters that cover the whole spectrum, the correct timings required to model the target dynamics are likely to be provided, at the cost of a considerably redundant representation of the system. An alternative approach has also been proposed, where the filters are tuned manually according to the information of the frequency spectrum of the desired response signal [26].

III. MODEL ARCHITECTURE

In this section we discuss the details of TORNN architecture. Let us consider a time series $\mathbf{x} = \{x[t]\}_{t=1}^T$, whose values are relative to the measurement of the states of a system, characterized by slow and fast dynamics. Each dynamical component is modeled by a specialized group of neurons in the recurrent layer which operates at the same characteristic frequencies. The number of groups K is equal to the number of main dynamical components of the system and can be determined through a frequency analysis on the input signal \mathbf{x} . In this work, by means of a power spectral density estimate, we identify K as the number of peaks in the power spectrum, but alternative approaches are possible.

In order to operate only on the portion of the spectrum located around one of the maxima, each processing unit implements a bandpass filter, whose configuration is identical to the one of the other units in the same group. As it is shown later in Sect. III-B, the transfer function of each bandpass filter depends on two parameters related to the two cutoff frequencies, but also on the connection weights of input and recurrent layers. If the configuration of these connections is modified, the frequency response of each filter will change as well. In TORNN we implement an hybrid approach. Analogously to reservoir computing methodologies, input and recurrent connections are randomly initialized and kept unchanged, while the pair of parameters that define the bandpass in each group of neurons are trained via backpropagation, along with the output layer.

In the following, we first explain how the recurrent layer is structured. Then, we describe how the bandpass filters are implemented in the hidden processing units and, finally, we discuss the learning procedure adopted to train the parameters in the model.

A. Recurrent layer structure

The recurrent layer of the network is randomly generated, under the constraint that the connections form

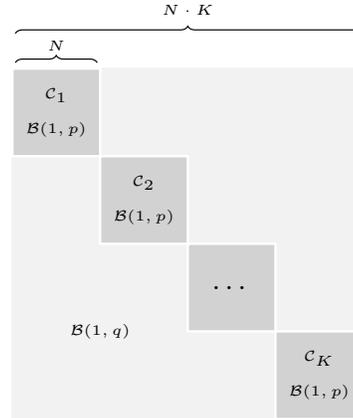


Fig. 1. Structure of the recurrent layer, encoded in the weight matrix \mathbf{W}_r . Inter-group connections are drawn from a binomial distribution $\mathcal{B}(1, p)$, while intra-group connections are drawn from $\mathcal{B}(1, q)$. To ensure inter-group connections to be more dense than intra-group connections, we set $p \gg q$.

a structure with K groups. Each group C_k contains N neurons, which are strongly connected to the other neurons in C_k . Additionally, in order to allow for some weak coupling between different dynamical components in the modeled system, we also generate a small amount of connections between neurons belonging to different groups. This allows synchronization mechanisms between different groups, which are required if the different dynamics of the modeled system are coupled. For intra-group connections, we define a probability p that determines the presence of a link $e_{i,j}$ between the neurons n_i and n_j of the same group C_k . A second probability q determines the presence of an inter-group connection $e_{i,j}$ between the neurons n_i and n_j of different groups. To guarantee higher intra-group connectivity we set p, q such that $p \gg q$.

To define the recurrent layer connections, we first generate a Boolean squared matrix $\mathbf{W}_r \in \mathbb{R}^{NK \times NK}$ by drawing values from a binomial distribution $\mathcal{B}(1, p)$ for elements belonging to one of K blocks on the diagonal, or a binomial distribution $\mathcal{B}(1, q)$, for the remaining elements. The structure of \mathbf{W}_r is depicted in Fig. 1. To obtain a higher degree of heterogeneity among neuron activations and to enrich internal dynamics [23], each non-zero connection $e_{i,j}$ is multiplied by a value, drawn from a normal distribution $\mathcal{N}(0, 1)$. Finally, to guarantee stability in the network [4], we rescale the spectral radius of \mathbf{W}_r to a value lower than 1.

Analogously to the connections in the recurrent layer, the input weights (stored in the matrix $\mathbf{W}_i \in \mathbb{R}^{I \times NK}$, with I the input size) are randomly initialized with

values drawn from a distribution $\mathcal{N}(0, 1)$ and are kept untrained as well.

B. bandpass processing units

A leaky integrator neuron outputs a weighted average – controlled by a leakage rate γ – of the new activation of the neuron with the one from the previous time interval [2, 17]. A leaky neuron acts as a lowpass filter, whose cutoff frequency is proportional to γ . Its state-update equation is defined as

$$x[t+1] = (1 - \gamma)x[t] + \gamma f(\mathbf{W}_r x[t] + \mathbf{W}_i u[t+1]) \quad (1)$$

or, by following an alternative definition [24], as

$$x[t+1] = f((1 - \gamma)x[t] + \gamma \mathbf{W}_r x[t] + \mathbf{W}_i u[t+1]) \quad (2)$$

where $f(\cdot)$ is a *tanh* (or a sigmoid) function. If $\gamma = 0$, the new state at time $t+1$ maintains the value of state t . If $\gamma = 1$, the state-update equation reduces to a classic nonlinear transfer function.

In the integrator of Eq. 2, the non-linearity $f(\cdot)$ is applied also to the leakage term. This guarantees stability, since the poles of the transfer function are constrained to the unit circle, and has the advantage of no computational overhead, since the integrators can be incorporated in \mathbf{W}_r [24]. This integrator always leaks, due to the contracting property of the non-linear mapping of the hyperbolic tangent upon itself. Since this is not an issue in TORNN, these are the filters we chose to implement.

As previously discussed, we want the processing units in each group of the recurrent layer to act as bandpass filters, which allow signals between two specific frequencies to pass, but discriminate against signals at other frequencies. In particular, we want the neurons in the group \mathcal{C}_k to reject the frequency components that are not close to a given peak in the spectrum. A bandpass filter can be generated by combining a lowpass filter with a highpass filter (implemented by negating the lowpass filter). Alternatively, a bandpass filter is obtained by combining two lowpass filters. In this latter case, according to Eq. 2, the state update equation of a bandpass neuron reads as

$$\begin{aligned} x' [t+1] &= f((1 - \gamma_1)x' [t] + \\ &\quad + \gamma_1 \mathbf{W}_r x[t] + \mathbf{W}_i u[t+1]), \\ x'' [t+1] &= (1 - \gamma_2)x'' [t] + \gamma_2 x' [t+1], \\ x[t+1] &= x' [t+1] - x'' [t+1]. \end{aligned} \quad (3)$$

The parameters γ_1 and γ_2 control respectively the high and low frequency cutoffs of the bandpass filter. Interestingly, the filter has a transfer function equivalent to the one of an analogue electronic RCL circuit [25].

C. Learning

The bandpass filters implemented by the units in the group \mathcal{C}_k must specialize to pass only the frequencies in neighborhood of one of the peaks in the power spectrum of \mathbf{x} . According to Eq. 3, to control the band to be passed by the filters in each group, one must tune the parameters γ_1 and γ_2 . However, in practical cases the bandpass width of each filter is not easy to determine in advance, as the neighborhood of each peak in the power spectrum could either include noise or useful information for the problem at hand. Most importantly, γ_1 and γ_2 are related to the effective cutoff frequencies in the spectrum through a highly non-linear dependency and the desired response of the filter is difficult to determine. Therefore, we follow a data-driven approach to automatically learn the control parameters by means of a gradient descent optimization, which accounts the nature of the task at hand.

The last component in the architecture is a dense layer, composed by a set of weights stored in the matrix $\mathbf{W}_o \in \mathbb{R}^{H \times O}$, which combines the output of the neurons in the recurrent layer to produce the O -dimensional output. A schematic depiction of the whole architecture is reported in Fig. 2.

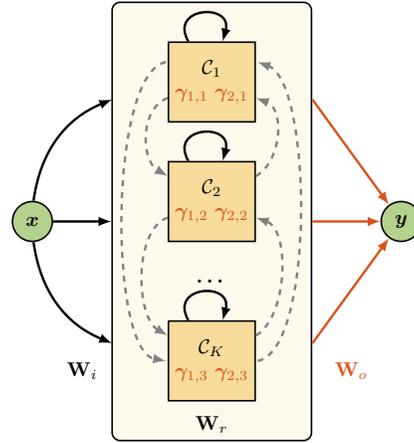


Fig. 2. Input connections (encoded in \mathbf{W}_i) and recurrent layer connections (encoded in \mathbf{W}_r) are initialized pseudo-randomly and they are kept fixed afterwards. In the recurrent layer, the connections between neurons in the same group, depicted as solid black lines, are more dense than the connections between neurons in different groups (dashed gray lines). In red, we marked the elements that are trained by means of a gradient descent optimization, which are the parameters that determine the cutoff frequencies in each group \mathcal{C}_k ($\gamma_{1,k}$ and $\gamma_{2,k}$) and the weights connecting the recurrent layer with the network output (\mathbf{W}_o).

Like the filter parameters, also the connection weights in \mathbf{W}_o are learned through a gradient descent optimization. We opted for a standard backpropagation through

time (BPTT) parametrized by a threshold τ_{inc} , which determines the truncation of the unfolding of the network in time. To constrain the parameters γ_1 and γ_2 to lay within the $[0, 1]$ interval (required from Eq. 3), during the optimization we apply a sigmoid function to the values of γ_1 and γ_2 learned by the gradient descent.

IV. EXPERIMENTS

In this section, we evaluate the capability of TORNN to model multiple dynamics characterized by different timescales. To test our model, we consider the prediction of a time series generated by superimposed sine waves, whose frequencies are incommensurable. Since the ratio of the frequencies is an irrational number, the periods of the sine waves do not cross with each other and hence the period of the resulting signal is, in theory, infinite. This academic, yet important task, has been previously considered to test the memory capacity of a recurrent neural network [16]. Indeed, to accurately predict the unseen values of the time series, the network would require a huge amount of memory. Here, we tackle the problem from a different perspective. We show that if the network is capable of learning how to model the dynamics of each single oscillator separately, the prediction task can be solved accurately with a limited amount of computational resources.

In the following, we try to solve several time series prediction tasks of increasing difficulty. Each time series is obtained by superimposing sinusoids, whose frequencies are pairwise incommensurable. The sinusoids are generated by multiplying a frequency ϕ by distinct integer powers of a transcendental number, such as e (Euler number) or π . In our experiments we chose e as the base number and the considered time series have the form

$$x_K[t] = \sum_{k=1}^K \sin(e^k \phi)[t], \quad (4)$$

where K is the number of superimposed sine waves. The difficulty of the prediction task is controlled by increasing the number of components K and by adding to the time series a white noise $n[t] \sim \mathcal{N}(0, 1)$. In the experiments, we set a noise-to-signal ratio of 0.2 and we evaluated time series with $T = 5000$ time-steps.

To quantify the performance of TORNN on each task, we compare its prediction accuracy with the ones achieved by Elmann-RNN (ERNN), LSTM, GRU and ESN. To initialize the trainable weights in TORNN, ERNN, LSTM and GRU, we draw their initial values from a Gaussian distribution $\mathcal{N}(0, 1/\sqrt{d})$, d being the number of processing units in the successive layer in the network [9]. For each network we used Adam

algorithm as the optimizer of the gradient descent step [18]. ERNN, LSTM and GRU are configured to use a comparable number of parameters (approximately 8500). For TORNN and ESN instead, we instantiate 20 neurons in the recurrent layer for each sinusoidal component in the signal. Therefore, the number of parameters trained in TORNN is $(2 \times K) + (20 \times K) + (1)$, which are, respectively, the filter parameters for all groups, the number of connections from the recurrent layer to the output and the bias. While the number of neurons in ESN reservoirs is usually much larger, the reason for this experimental setup is to show that TORNN manages to handle the same computational resources of an ESN more effectively.

The details of the configuration in each network are reported in Tab. I. The model parameters of ESN are not learned through gradient descent optimization. Indeed, the output weights in the readout are computed through a simple ridge-regression, an optimization problem that can be expressed in close form and solved in a time that grows as the cube of the number of neurons [3]. On the other hand, a correct setup of the hyperparameters in the ESN is in general more critical than in other architectures and they are usually tuned through cross validation procedures. In our experiments, we used a genetic algorithm to search for the optimal hyperparameters values. We followed the same approach described in [20], to which we refer the interested reader for further details. The bounds considered in the search of each hyperparameter are reported in Tab I.

In each task we perform a prediction with a forecast horizon of 15 time steps. Prediction error is expressed in terms of Normalized Root Mean Squared Error (NRMSE), computed as

$$\text{NRMSE} = \sqrt{\langle \|\mathbf{y} - \mathbf{y}^*\|^2 \rangle / \langle \|\mathbf{y} - \langle \mathbf{y}^* \rangle\|^2 \rangle},$$

where \mathbf{y} is the output predicted by the network and \mathbf{y}^* the ground truth. We consider 4 different time series $x_K[t]$, generated according to Eq. 4, with $K = 2, 3, 5, 7$ the number of superimposed oscillators. The power spectral density estimates of the time series are depicted in Fig. 3. As is it possible to see, in the spectra of $x_5[t]$ and $x_7[t]$ there are two frequencies which are very close. This increase the difficulty of the problem as these frequencies are harder to separate. For each time series, we also consider a version with superimposed white noise $n[t]$.

The prediction accuracies are reported in Fig. 4, where the colored boxes represent the average NRMSE error (the shorter, the better) and the error bars are the standard deviations, computed over 10 different trials with independent random initializations. TORNN achieves better performance in every test. In order to provide

TABLE I

SUMMARY OF THE HYPERPARAMETERS CONFIGURATION FOR EACH NETWORK. ERNN, LSTM, GRU: NUMBER OF PROCESSING UNITS (N_r), NUMBER OF TRAINABLE PARAMETERS (# PARAMS), L_2 REGULARIZATION OF \mathbf{W}_o WEIGHTS (λ), GRADIENT TRUNCATION IN BPPT (τ_{tnc}). ESN (ADMISSIBLE RANGE OF HYPERPARAMETERS, OPTIMIZED WITH A GENETIC ALGORITHM [20]): NEURONS OF THE RESERVOIR (N_r – NOT OPTIMIZED), SPECTRAL RADIUS (ρ), RESERVOIR CONNECTIVITY (r), NOISE IN STATE UPDATE (ξ), INPUT SCALING (ω_i), TEACHER SCALING (ω_o), FEEDBACK SCALING (ω_f), L_2 RIDGE REGRESSION NORMALIZATION (λ). TORNN: PROCESSING UNITS PER GROUP (N), INTRA-GROUP CONNECTIONS PROBABILITY (p), INFRA-GROUP CONNECTION PROBABILITY (q), SPECTRAL RADIUS OF \mathbf{W}_r (ρ), L_2 REGULARIZATION OF \mathbf{W}_o WEIGHTS (λ) AND GRADIENT TRUNCATION IN BPPT (τ_{tnc}).

ERNN	N_r 91	# params 8555	λ 1E-5	τ_{tnc} 10				
LSTM	N_r 45	# params 8506	λ 1E-5	τ_{tnc} 10				
GRU	N_r 52	# params 8477	λ 1E-5	τ_{tnc} 10				
ESN	N_r $20 \times K$	ρ [0.1, 1.5]	r [0.1, 0.5]	ξ [0, 0.1]	ω_i [0.1, 1]	ω_o [0.1, 1]	ω_f [0.1, 1]	λ [0, 0.5]
TORNN	N 20	p 0.4	q 0.1	ρ 0.95	λ 1E-5	τ_{tnc} 10		

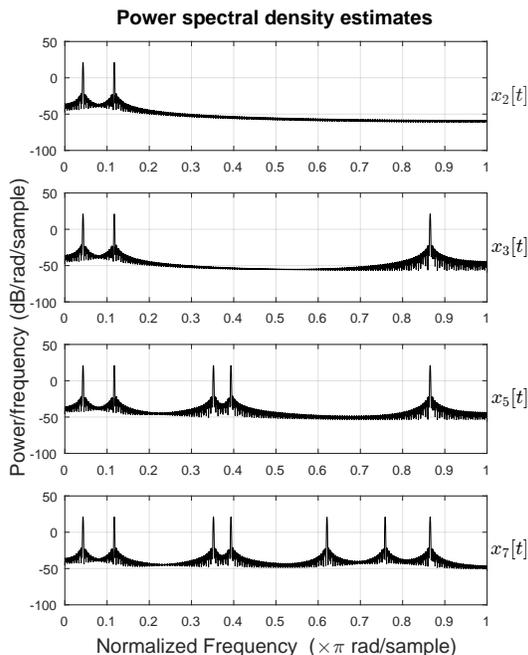


Fig. 3. Power spectral densities of the four time series $x_2[t]$, $x_3[t]$, $x_5[t]$ and $x_7[t]$ considered in the experiments. Here, we report the spectrum of the versions without the superimposed white noise.

a qualitative description of the forecast results, in Fig. 5 we show a portion of the prediction of the time series $x_7[t] + n[t]$ by TORNN and the other considered networks, respectively RNN, LSTM, GRU and ESN. In every plot the gray line represents the ground truth, the black line the considered network’s forecast and the light red area the residual between the two lines. By looking

at the residual areas of each time series, it is possible to notice an overall better prediction by TORNN with respect to the other RNNs.

These results provide empirical evidence of the effectiveness of the proposed framework, which is capable of modeling the different components and thus performing a more accurate prediction. As expected, when time series are corrupted by noise the prediction forecast accuracy heavily decreases in every model.

The state-of-the-art architectures based on gradient optimization, which are ERNN, LSTM and GRU, obtain on average a larger prediction error with respect to ESN and TORNN. Furthermore, ERNN performs worse than LSTM and GRU in every test. This is a consequence of the simplicity in the architecture, which is unable to handle properly long-term dependencies in the analyzed time series. The performance of LSTM and GRU are comparable and it is not possible to conclude which one is the best performing. This result is expected and is in agreement with previous studies which evidenced that the selection of a specific gated architecture should depend on the dataset and the corresponding task at hand [6]. However, even if the GRU structure is simpler than the plain LSTM architecture (without peephole connections), the training time of GRU is higher, due to the larger number of operations required to compute the forward and the backward propagation steps [5].

The performance of TORNN are always better or at least equal to ESN. It is worth noticing that TORNN can be considered as a general case of an ESN, since its architecture reduces to ESN when $\gamma_1 = 1$ and $\gamma_2 = 0$ in every group of units. Indeed, with such a configuration the bandpass filters reduce to regular neurons. Finally, we underline that in ESN the training of the model is

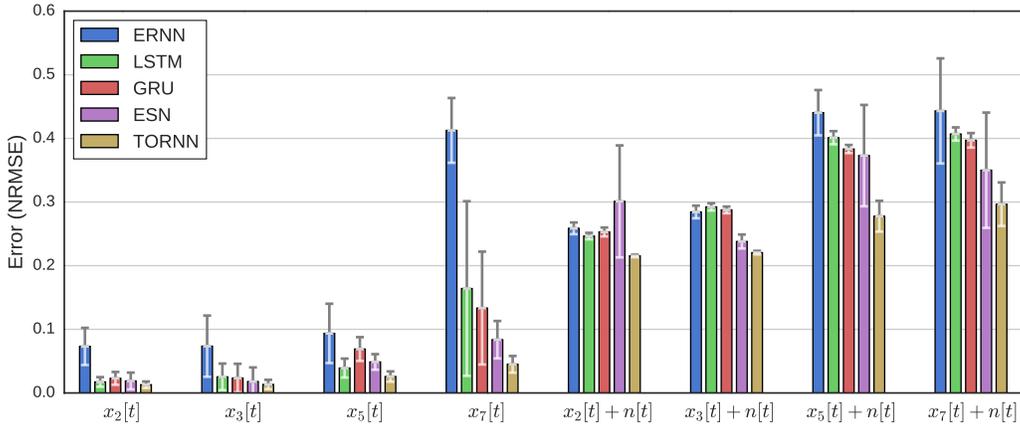


Fig. 4. Forecast accuracy (NRMSE) achieved by TORNN and other RNN architectures on the prediction of different time series, relative to a forecast horizon of 15 time intervals. Colored boxes represent the average error on 10 independent trials, the error bars the standard deviation. Each time series $x_K[t]$ is the combination of K sinusoids with incommensurable frequencies and they are evaluated with and without a superimposed white noise $n[t]$.

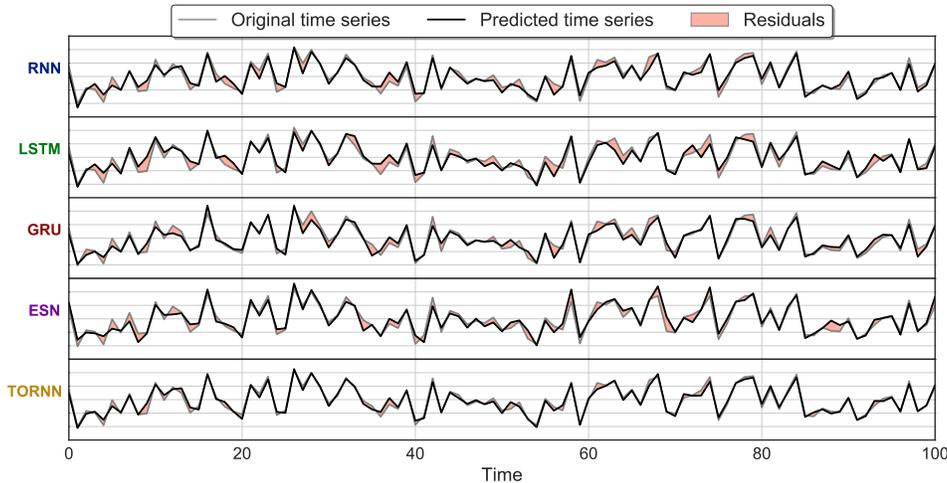


Fig. 5. Prediction results for a fraction of the time series $x_7[t] + n[t]$ with a prediction step-ahead of 15 time steps. For each network, the gray line represents the ground truth, the black line the predicted time series and the shaded red areas are the error residuals of the prediction.

faster than in the other architectures. In fact, rather than a slow gradient descent optimization, the weights of the linear readout in a ESN can be evaluated by performing a simple linear regression. However, ESNs trade the precision of gradient descent with the “brute force” redundancy of random reservoirs and this, inevitably, makes the models more sensitive to the selection of the hyperparameters (like spectral radius). Hence, the computational resources required for an accurate search of the optimal hyperparameters should be accounted in

the comparison with other architectures.

V. CONCLUSIONS

In this work we presented the Temporal Overdrive Recurrent Neural Network, a novel RNN architecture designed to model multiple dynamics that are characterized by different time scales. The proposed model is easy to configure, as it only requires to specify the number of different time scales that should be accounted for, an information that can be easily obtained from a

rough frequency analysis. Each dynamical component is handled by a group of neurons implementing a bandpass filter, whose behavior is determined by two parameters that determine the cutoff frequencies.

The proposed methodology follows the strategy of reservoir computing approaches, as the input and recurrent connections are randomly initialized and they are kept untrained. However, while reservoir computing implements a “brute force” approach to generate a high, yet redundant, number of internal dynamics, TORNN learns from data the optimal configuration of its dynamics through a gradient descent optimization. Furthermore, with respect to other gradient-based RNNs, the number of trainable parameters is significantly lower, with a consequent simplification of the learning procedure.

We performed some preliminary tests on synthetic data, which showed promising results and demonstrated that our network achieves superior performance in prediction with respect to other state-of-the-art RNNs. TORNN can also be used for system identification [22]. In this case, the network must operate in a generative setup and the output of the network has to be fed back into the recurrent layer.

We are currently working on real-world data, relative to load forecast, which will be presented in a future extension of this work. We also plan to explore the implementation of more efficient bandpass filters, with sharper frequency cutoffs. We believe that this will help in separating the dynamical components more effectively. In the future we also plan to investigate further the internal dynamics of the network.

REFERENCES

- [1] Y. Bengio, P. Simard, and P. Frasconi. Learning long-term dependencies with gradient descent is difficult. *IEEE Transactions on Neural Networks*, 5(2):157–166, Mar 1994. ISSN 1045-9227. doi: 10.1109/72.279181.
- [2] Y. Bengio, N. Boulanger-Lewandowski, and R. Pascanu. Advances in optimizing recurrent networks. In *2013 IEEE International Conference on Acoustics, Speech and Signal Processing*, pages 8624–8628. IEEE, 2013.
- [3] F. M. Bianchi, S. Scardapane, A. Uncini, A. Rizzi, and A. Sadeghian. Prediction of telephone calls load using Echo State Network with exogenous variables. *Neural Networks*, 71: 204–213, 2015. doi: 10.1016/j.neunet.2015.08.010.
- [4] F. M. Bianchi, L. Livi, and C. Alippi. Investigating echo state networks dynamics by means of recurrence analysis. *arXiv preprint arXiv:1601.07381*, 2016.
- [5] K. Cho, B. Van Merriënboer, C. Gulcehre, D. Bahdanau, F. Bougares, H. Schwenk, and Y. Bengio. Learning phrase representations using rnn encoder-decoder for statistical machine translation. *arXiv preprint arXiv:1406.1078*, 2014.
- [6] J. Chung, Ç. Gülçehre, K. Cho, and Y. Bengio. Empirical evaluation of gated recurrent neural networks on sequence modeling. *CoRR*, abs/1412.3555, 2014. URL <http://arxiv.org/abs/1412.3555>.
- [7] J. Chung, C. Gülçehre, K. Cho, and Y. Bengio. Gated feedback recurrent neural networks. *CoRR*, abs/1502.02367, 2015.
- [8] S. El Hahi and Y. Bengio. Hierarchical recurrent neural networks for long-term dependencies. In *NIPS*, volume 400, page 409. Citeseer, 1995.
- [9] X. Glorot and Y. Bengio. Understanding the difficulty of training deep feedforward neural networks. In *Aistats*, volume 9, pages 249–256, 2010.
- [10] P. G. Gould, A. B. Koehler, J. K. Ord, R. D. Snyder, R. J. Hyndman, and F. Vahid-Araghi. Forecasting time series with multiple seasonal patterns. *European Journal of Operational Research*, 191(1):207 – 222, 2008. ISSN 0377-2217. doi: <http://dx.doi.org/10.1016/j.ejor.2007.08.024>.
- [11] A. Graves. Generating sequences with recurrent neural networks. *arXiv preprint arXiv:1308.0850*, 2013.
- [12] K. Greff, R. K. Srivastava, J. Koutnik, B. R. Steunebrink, and J. Schmidhuber. LSTM: A search space odyssey. *CoRR*, abs/1503.04069, 2015. URL <http://arxiv.org/abs/1503.04069>.
- [13] B. Hammer, A. Micheli, A. Sperduti, and M. Strickert. Recursive self-organizing network models. *Neural Networks*, 17(8-9):1061–1085, 2004. ISSN 0893-6080. doi: 10.1016/j.neunet.2004.06.009.
- [14] S. Hochreiter and J. Schmidhuber. Long short-term memory. *Neural computation*, 9(8):1735–1780, 1997.
- [15] H. Jaeger. The “echo state” approach to analysing and training recurrent neural networks—with an erratum note. *Bonn, Germany: German National Research Center for Information Technology GMD Technical Report*, 148:34, 2001.
- [16] H. Jaeger and H. Haas. Harnessing nonlinearity: Predicting chaotic systems and saving energy in wireless communication. *science*, 304(5667):78–80, 2004.
- [17] H. Jaeger, M. Lukoševičius, D. Popovici, and U. Siewert. Optimization and applications of echo state networks with leaky-integrator neurons. *Neural Networks*, 20(3):335 – 352, 2007. ISSN 0893-6080. doi: <http://dx.doi.org/10.1016/j.neunet.2007.04.016>. Echo State Networks and Liquid State Machines.
- [18] D. Kingma and J. Ba. Adam: A method for stochastic optimization. *arXiv preprint arXiv:1412.6980*, 2014.
- [19] T. Lin, B. G. Horne, P. Tiño, and C. L. Giles. Learning long-term dependencies in narx recurrent neural networks. *Neural Networks, IEEE Transactions on*, 7(6):1329–1338, 1996.
- [20] S. Løkse, F. M. Bianchi, and R. Jenssen. Training echo state networks with regularization through dimensionality reduction. *CoRR*, abs/1608.04622, 2016. URL <http://arxiv.org/abs/1608.04622>.
- [21] Maass Wolfgang, Natschläger Thomas, and Markram Henry. Real-Time Computing Without Stable States: A New Framework for Neural Computation Based on Perturbations. *Neural Computation*, 14(11):2531–2560, 2002. ISSN 0899-7667. doi: <http://dx.doi.org/10.1162/089976602760407955>. doi: 10.1162/089976602760407955.
- [22] O. Nelles. *Nonlinear system identification: from classical approaches to neural networks and fuzzy models*. Springer Science & Business Media, 2013.
- [23] M. C. Ozturk, D. Xu, and J. C. Principe. Analysis and design of echo state networks. *Neural Computation*, 19(1):111–138, 2007.
- [24] B. Schrauwen, J. Defour, D. Verstraeten, and J. Van Campenhout. *The Introduction of Time-Scales in Reservoir Computing, Applied to Isolated Digits Recognition*, pages 471–479. Springer Berlin Heidelberg, Berlin, Heidelberg, 2007. ISBN 978-3-540-74690-4. doi: 10.1007/978-3-540-74690-4_48.
- [25] U. Siewert and W. Wustlich. Echo-state networks with band-pass neurons: Towards generic time-scale-independent reservoir structures. *Internal status report, PLANET intelligent systems GmbH*, 2007.
- [26] F. Wyffels, B. Schrauwen, D. Verstraeten, and D. Stroobandt. Band-pass reservoir computing. In *2008 IEEE International Joint Conference on Neural Networks (IEEE World Congress on Computational Intelligence)*, pages 3204–3209, June 2008. doi: 10.1109/IJCNN.2008.4634252.