# Sequence Prediction using Spectral RNNs

**Moritz Wolter**
Institute for Computer Science, University of Bonn
Fraunhofer Center for Machine Learning and SCAI
wolter@cs.uni-bonn.de

**Jürgen Gall**
Institute for Computer Science, University of Bonn
gall@cs.uni-bonn.de

**Angela Yao**
National University of Singapore
ayao@comp.nus.edu.sg

August 17, 2020

## Abstract

Fourier methods have a long and proven track record as an excellent tool in data processing. As memory and computational constraints gain importance in embedded and mobile applications, we propose to combine Fourier methods and recurrent neural network architectures. The short-time Fourier transform allows us to efficiently process multiple samples at a time. Additionally, weight reductions trough low pass filtering is possible. We predict time series data drawn from the chaotic Mackey-Glass differential equation and real-world power load and motion capture data. [1]

***Keywords*** Sequence Modelling · Frequency Domain · Short Time Fourier Transform

## 1 Introduction

Deployment of machine learning methods in embedded and real-time scenarios leads to challenging memory and computational constraints. We propose to use the short time Fourier transform (STFT) to make sequential data forecasting using recurrent neural networks (RNNs) more efficient. The STFT windows the data and moves each window into the frequency domain. A subsequent recurrent network processes multi sample windows instead of single data points and therefore runs at a lower clock rate, which reduces the total number of RNN cell evaluations per time step. Additionally working in the frequency domain allows parameter reductions trough low-pass filtering. Combining both ideas reduces the computational footprint significantly. We hope this reduction will help deploying RNNs on VR-devices for human pose prediction or on embedded systems used for load forecasting and management.

We show in our work that it is possible to propagate gradients through the STFT *and* its inverse. This means that we can use time domain-based losses that match the (temporal) evaluation measures, but still apply the RNN in the frequency domain, which imparts several computation and learning benefits. Firstly, representations can often be more compact and informative in the frequency domain. The Fourier basis clusters the most important information in the low-frequency coefficients. It conveniently allows us to integrate a low-pass filter to not only reduce the representation size but also remove undesirable noise effects that may otherwise corrupt the learning of the RNN.

One of the challenges of working with the Fourier transform is that it requires handling of complex numbers. While not yet mainstream, complex-valued representations have been integrated into deep convolutional architectures [1], RNNs [2, 3, 4] and Hopfield like networks [5]. Complex-valued networks require some additional overhead for book-keeping between the real and imaginary components of the weight matrices and vector states, but tend to be more parameter efficient. We compare fully complex and real valued-concatenation formulations.

---

[1]Source code available at https://github.com/v0lta/Spectral-RNN

Naively applying the RNN on a frame-by-frame basis does not always lead to smooth and coherent predictions [6]. When working with short time Fourier transformed (STFT) windows of data, smoothness can be built-in by design through low-pass filtering before the inverse short time Fourier transformation (iSTFT). Furthermore, by reducing the RNN's effective clockrate, we extend the memory capacity (which usually depends on the number of recurrent steps taken) and achieve computational gains.

In areas such as speech recognition [7] and audio processing [8], using the STFT is a common pre-processing step, *e.g.* as a part of deriving cepstral coefficients. In these areas, however, the interest is primarily for classification based on the spectrum coefficients and the complex phase is discarded, so there is no need for the inverse transform and the recovery of a temporal signal as output. We advocate the use of the forwards and inverse STFT directly before and after recurrent memory cells and propagate gradients through both the forwards and inverse transform. In summary,

- we propose a novel RNN architecture for analyzing temporal sequences using the STFT and its inverse.
- We investigate the effect of real and complex valued recurrent cells.
- We demonstrate how our windowed formulation of sequence prediction in the spectral domain, based on the reduction in data dimensionality and rate in RNN recursion can, significantly improve efficiency and reduce overall training and inference time.

## 2   Related Works

Prior to the popular growth of deep learning approaches, analyzing time series and sequence data in the spectral domain, *i.e.* via a Fourier transform was a standard approach [9]. In machine learning, Fourier analysis is mostly associated with signal-processing heavy domains such as speech recognition [7], biological signal processing [10], medical image [11] and audio-processing [8]. The Fourier transform has already been noted in the past for improving the computational efficiency of convolutional neural networks (CNNs) [12, 13]. In CNNs, the efficiency comes from the duality of convolution and multiplication in space and frequency. Fast GPU implementations of convolution, *e.g.* in the NVIDIA cuDNN library [14] use Fourier transforms to speed up their processing. Such gains are especially relevant for convolutional neural networks in 2D [12, 13] and even more so in 3D [15]. However, the improvement in computational efficiency for the STFT-RNN combination comes from reductions in data dimensionality and RNN clock rate recursion. The Fourier transform has been used in neural networks in various contexts [16]. Particularly notable are networks, which have Fourier-like units [17] that decompose time series into constituent bases. Fourier-based pooling has been explored for CNNs as an alternative to standard spatial max-pooling [18].

Various methods for adjusting the rate of RNN recurrency have been explored in the past. One example [19] introduces additional recurrent connections with time lags. Others apply different clock rates to different parts of the hidden layer [20, 21], or apply recursion at multiple (time) scales and or hierarchies [22, 23]. All of these approaches require changes to the basic RNN architecture. Our proposal, however, is a simple alternative which does not require adding any new structures or connections. Among others [24] approached mocap prediction using RNNs, while [6] applied a cosine transform CNN combination. In this paper we explore the combination of the complex valued windowed STFT and RNNs on mackey-glass, power-load and mocap time series data.

## 3   Short Time Fourier Transform

### 3.1   Forwards STFT

The Fourier transform maps a signal into the spectral or frequency domain by decomposing it into a linear combination of sinusoids. In its regular form the transform requires infinite support of the time signal to estimate the frequency response. For many real-world applications, however, including those which we wish to model with RNNs, this requirement is not only infeasible, but it may also be necessary to obtain insights into changes in the frequency spectrum as function of time or signal index. To do so, one can partition the signal into overlapping segments and approximate the Fourier transform of each segment separately. This is the core idea behind the short time Fourier transform (STFT), it determines a signal's frequency domain representation as it changes over time.

More formally, given a signal $\mathbf{x}$, we can partition it into segments of length $T$, extracted every $S$ time steps. The STFT $\mathcal{F}_s$ of $\mathbf{x}$ is defined by [25] as the discrete Fourier transform of $\mathbf{x}$, *i.e.*

$$\mathbf{X}[\omega, Sm] = \mathcal{F}_s(\mathbf{x})$$

$$= \mathcal{F}(\mathbf{w}[Sm - l]\mathbf{x}[l]) = \sum_{l=-\infty}^{\infty} \mathbf{w}[Sm - l]\mathbf{x}[l]e^{-j\omega l} \tag{1}$$

where $\mathcal{F}$ denotes the classic discrete fast Fourier transform. Segments of $\mathbf{x}$ are multiplied with the windowing function $\mathbf{w}$ and transformed afterwards.

Historically, the shape and width of the window function has been selected by hand. To hand the task of window selection to the optimizer, we work with a truncated Gaussian window [26]

$$w[n] = \exp\Big(-\frac{1}{2}\Big(\frac{n - T/2}{\sigma T/2}\Big)^2\Big) \tag{2}$$

of size T and learn the standard deviation $\sigma$. The larger that $\sigma$ is, the more the window function approaches a rectangular window; the smaller the sigma, the more extreme the edge tapering and as a result, the narrower the window width.

## 3.2 Inverse STFT

Supposing that we are given some frequency signal $\mathbf{X}$; the time signal $\hat{\mathbf{x}}$ represented by $\mathbf{X}$ can be recovered with the inverse short time Fourier transform (iSTFT) $\mathcal{F}_s^{-1}$ and is defined by [25] as:

$$\hat{\mathbf{x}}[n] = \mathcal{F}_s^{-1}(\mathbf{X}[n, Sm]) = \frac{\sum_{m=-\infty}^{\infty} \mathbf{w}[Sm - n]\hat{\mathbf{x}}_w[n, Sm]}{\sum_{m=-\infty}^{\infty} \mathbf{w}^2[Sm - n]} \tag{3}$$

where the signal $\hat{\mathbf{x}}_w$ is the inverse Fourier transform of $\mathbf{X}$:

$$\hat{\mathbf{x}}_w = \frac{1}{T}\sum_{l=-\infty}^{\infty} \mathbf{X}[l, Sm]e^{j\omega l} \tag{4}$$

and $l$ indexing the frequency dimension of $\mathbf{X}_m$. Eq. 3 reverses the effect of the windowing function, but implementations require careful treatment of the denominator to prevent division by near-zero values[2]. In Eq. 3, $Sm$ generally evaluates to an integer smaller than the window size $T$ and subsequent elements in the sum overlap, hence the alternative naming of it being an *"overlap and add"* method [27].

# 4 Complex Spectral Recurrent Neural Networks

## 4.1 Network Structure

We can now move RNN processing into the frequency domain by applying the STFT to the input signal $x$. If the output or projection of the final hidden vector is also a signal of the temporal domain, we can also apply the iSTFT to recover the output $y$. This can be summarized by the following set of RNN equations:

$$\mathbf{X}_\tau = \mathcal{F}(\{\mathbf{x}_{S\tau - T/2}, \dots, \mathbf{x}_{S\tau + T/2}\}) \tag{5}$$

$$\mathbf{z}_\tau = \mathbf{W}_c \mathbf{h}_{\tau-1} + \mathbf{V}_c \mathbf{X}_\tau + \mathbf{b}_c \tag{6}$$

$$\mathbf{h}_\tau = f_a(\mathbf{z}_\tau) \tag{7}$$

$$\mathbf{y}_\tau = \mathcal{F}^{-1}(\{\mathbf{W}_{pc}\mathbf{h}_0, \dots, \mathbf{W}_{pc}\mathbf{h}_\tau\}) \tag{8}$$

where $\tau = [0, n_s]$, *i.e.* from zero to the total number of segments $n_s$. The output $y_\tau$ may be computed based on the available outputs $\{\mathbf{W}_p\mathbf{h}_0, \dots, \mathbf{W}_p\mathbf{h}_\tau\}$ at step $\tau$. Adding the STFT-iSTFT pair has two key implications. First of all, because $\mathbf{X}_\tau \in \mathbb{C}^{n_f \times 1}$ is a complex signal, the hidden state as well as subsequent weight matrices all become complex, *i.e.* $\mathbf{h}_\tau \in \mathbb{C}^{n_h \times 1}$, $\mathbf{W}_c \in \mathbb{C}^{n_h \times n_h}$, $\mathbf{V}_c \in \mathbb{C}^{n_h \times n_f}$, $\mathbf{b}_c \in \mathbb{C}^{n_h \times 1}$ and $\mathbf{W}_{pc} \in \mathbb{C}^{n_h \times n_f}$, where $n_h$ is the hidden size of the network as before and $n_f$ is the number of frequencies in the STFT.

The second implication to note is that the step index changes from frame $t$ to $\tau$, which means that the spectral RNN effectively covers $S$ time steps of the standard RNN per step. This has significant bearing on the overall memory consumption as well as the computational cost, both of which influence the overall network training time. Considering only the multiplication of the state matrix $\mathbf{W}_c$ and the state vector $\mathbf{h}_\tau$, which is the most expensive operation, the basic RNN requires $N \cdot O(n_h^3)$ operations for $N$ total time steps. When using the Fourier RNN with an FFT implementation of the STFT, one requires only

$$N/S \cdot (O(T \log T) + O(n_h^3)), \tag{9}$$

---

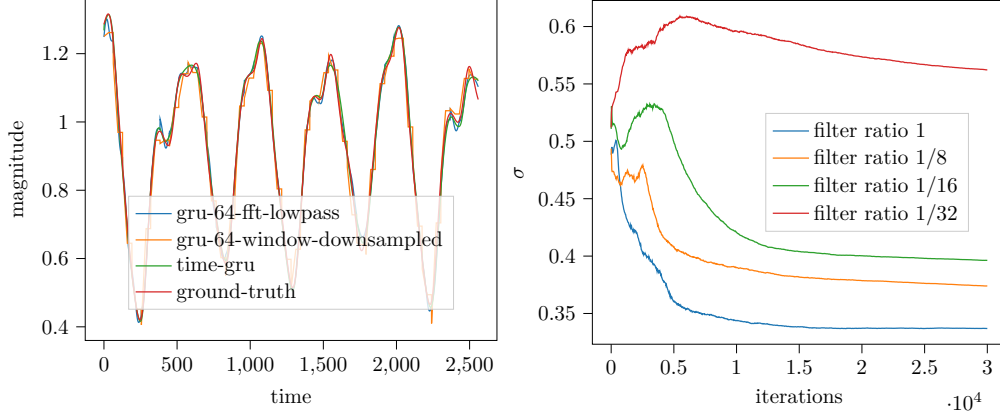[2]We adopt the common strategy of adding a small tolerance $\epsilon = 0.001$

Figure 1: Mackey-Glass series predictions in for different RNN methods. As gradients flow through the STFT we can optimize the width of the gaussian $\sigma$. The learned window width for increasing degrees of low-pass filtering is shown here. Figure best viewed in colour.

where the $T \log T$ term comes from the FFT operation. The architectural changes lead to larger input layers and fewer RNN iterations. $\mathbf{X}$ is higher dimensional than $\mathbf{x}$, but we save on overall computation if the step size is large enough which will make $N/S$ much smaller than $N$. We can generalise the approach described above into:

$$\mathbf{X}_\tau = \mathcal{F}(\{\mathbf{x}_{S\tau-T/2}, \ldots, \mathbf{x}_{S\tau+T/2}\}) \tag{10}$$

$$\mathbf{h}_t = \text{RNN}_\mathbb{C}(\mathbf{X}_\tau, \mathbf{h}_{t-1}) \tag{11}$$

$$\mathbf{y}_t = \mathcal{F}^{-1}(\{\mathbf{W}_{pc}\mathbf{h}_0, \ldots, \mathbf{W}_{pc}\mathbf{h}_\tau\}), \tag{12}$$

where instead of the basic formulation outlined above, more sophisticated complex-valued RNN-architectures [2, 3, 4] represented by $\text{RNN}_\mathbb{C}$ may be substituted. We experiment with a complex-valued GRU, as proposed in [4]. An alternative to a complex approach is to concatenate the real and imaginary components into one (real-valued) hidden vector. The experimental section compares both methods in Table 2.

### 4.2 Loss Functions

In standard sequence prediction, the loss function applied is an L2 mean squared error, applied in the time domain:

$$\mathcal{L}_{mse}(\mathbf{y}_t, \mathbf{y}_{gt}) = \frac{1}{n_y} \sum_{l=0}^{n_y} (\mathbf{y}_t - \mathbf{y}_{gt})^2. \tag{13}$$

We experiment with a similar variation applied to the STFT coefficients applied in the frequency domain (see Table 2) but find that it performs on par but usually a little bit worse than the time-domain MSE. This is not surprising, as the evaluation measure applied is still in the time domain for sequence prediction so it works best to use the same function as the loss.

## 5 Mackey-Glass Chaotic Sequence Prediction

Initially we study our method by applying it to make predictions on the Mackey-Glass series [28]. The Mackey-Glass equation is a non-linear time-delay differential and defines a chaotic system that is semi-periodic: We first evaluate different aspects of our model on the chaotic Mackey-Glass time series [28]:

$$\frac{dx}{dt} = \frac{\beta x_\tau}{1 + x_\tau^n} - \gamma x, \tag{14}$$

with $\gamma = 0.1, \beta = 0.2$ and the power parameter to $n = 10$. $x_\tau$ denotes the value from $\tau$ time steps ago; we use a delay of $\tau = 17$ and simulate the equation in the interval $t \in [0, 512]$, using a forward Euler scheme with a time step of size 0.1. During the warm-up, when true values are not yet known, we use a uniform distribution and randomly draw values from $1 + U[-0.1, 0.1]$. An example of the time series can be found in Figure 1; we split the signal in half, conditioning on the first half as input and predicting the second half.

Table 1: Short time Fourier, Windowed and Time Domain results obtained using GRU cells of size 64. Windowed experiments process multiple samples of data without computing the STFT. Additionally we compare low-pass filtering the spectrum and downsampling the time domain windows. All models where trained for 30k iterations. We downsample and lowpass-filter with a factor of 1/32.

| net | weights | mse | training-time [min] |
|---|---|---|---|
| time-GRU | 13k | $3.8 \cdot 10^{-4}$ | 355 |
| time-GRU-window | 29k | $6.9 \cdot 10^{-4}$ | 53 |
| time-GRU-window-down | 13k | $12 \cdot 10^{-4}$ | 48 |
| STFT-GRU | 46k | $3.5 \cdot 10^{-4}$ | 57 |
| STFT-GRU-lowpass | 14k | $2.7 \cdot 10^{-4}$ | 56 |

Table 2: Real and complex valued architecture comparison on the mackey-glass data, with increasing complex cell size. The complex architectures take longer to run but are more parameter efficient. The last row shows a complex RNN cell in STFT space without iSTFT backpropagation.

| net | weights | mse | training-time [min] |
|---|---|---|---|
| STFT-GRU-64 | 46k | $3.5 \cdot 10^{-4}$ | 57 |
| STFT-cgRNN-32 | 23k | $2.1 \cdot 10^{-4}$ | 63 |
| STFT-cgRNN-54 | 46k | $1.6 \cdot 10^{-4}$ | 63 |
| STFT-cgRNN-64 | 58k | $1.1 \cdot 10^{-4}$ | 64 |
| STFT-cgRNN-64-$\mathcal{L}_{\mathcal{C}}$ | 58k | $210 \cdot 10^{-4}$ | 64 |

### 5.1 Implementation Details

In all experiments, we use RNN architectures based on real [29] and complex [4] GRU-cells with a state size of 64 and a Gaussian window of width 128 initialized at $\sigma = 0.5$ unless otherwise stated. The learning rate was set intially to 0.001 and then reduced with a stair-wise exponential decay with a decay rate of 0.9 every 1000 steps. Training was stopped after 30k iterations. Our run-time measurements include ordinary differential equation simulation, RNN execution as well as back-propagation time.

### 5.2 Experimental Results and Ablation Studies

**Fourier transform ablation**; We first compare against two time-based networks: a standard GRU (time-GRU) and a windowed version (time-GRU-windowed) in which we reshape the input and process windows of data together instead of single scalars. This effectively sets the clock rate of the RNN to be per window rather than per time step. As comparison, we look at the STFT-GRU combination as described in Section 4.1 with a GRU-cell and a low-pass filtered version keeping only the first four coefficients (STFT-GRU-lowpass). Additionally we compare low-pass filtering to time windowed time domain downsampling (time-GRU-window-down). For all five networks, we use a fixed hidden-state size of 64. From Figure 1, we observe that all five architecture variants are able to predict the overall time series trajectory. Results in Table 1 indicate that reducing the RNN clock rate through windowing and parameters through low pass filtering also improves prediction quality. In comparison to time domain downsampling, frequency-domain lowpass filtering allows us to reduce parameters more aggressively.

**Runtime**; As discussed in Equation 9, windowing reduces the computational load significantly. In Table 1, we see that the windowed experiments run much faster than the naive approach. The STFT networks are slightly slower, since the Fourier Transform adds a small amount of computational overhead.

**The window function**; Figure 1 shows the Gaussian width for multiple filtering scenarios. We observe a gradient signal on the window function. For more aggressive filtering the optimizer chooses a wider kernel.

**Complex vs. real cell**; We explore the effect of a complex valued cell in Table 2. We apply the complex GRU proposed in [4], and compare to a real valued GRU. We speculate that complex cells are better suited to process Fourier representations due to their complex nature. Our observations show that both cells solve the task, while the complex cell does so with fewer parameters at a higher computational cost. The increased parameter efficiency of the complex cell could indicate that complex arithmetic is better suited than real arithmetic four Frequency domain machine learning. However due to the increased run-time we proceed using the concatenation approach.

**Time vs. Frequency mse**; One may argue that propagating gradients through the iSTFT is unnecessary. We tested this setting and show a result in the bottom row of Table 2. In comparison to the accuracy of the otherwise equivalent complex network shown above the second horizontal line, computing the error on the Fourier coefficients performs
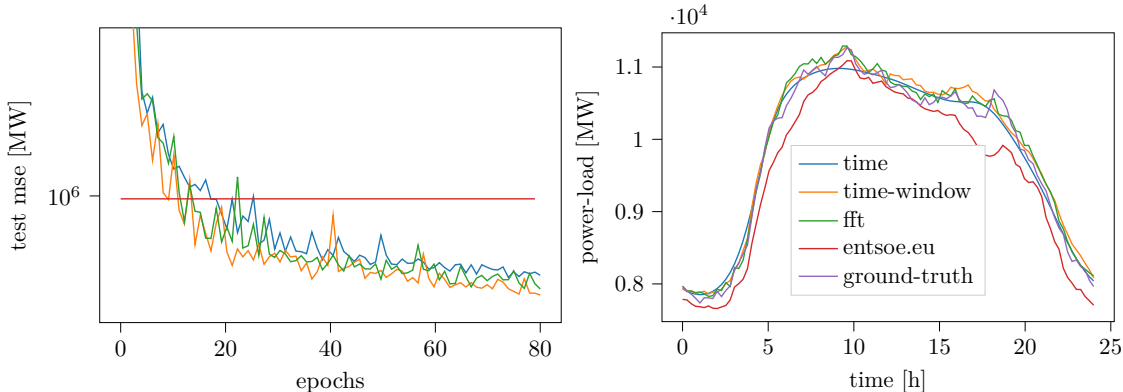
Figure 2: Day ahead prediction results convergence (left) and prediction examples (right). We observe that all deep learning approaches beat the entsoe.eu baseline shown as the red line, which suggests that their approach could benefit from deep learning.

significantly worse. We conclude that if our metric is in the time domain the loss needs to be there too. We therefore require gradient propagation through the STFT.

## 6 Power Load Forecasting

### 6.1 Data

We apply our Fourier RNN to the problem of forecasting power consumption. We use the power load data of 36 EU countries from 2011 to 2019 as provided by the European Network of Transmission System Operators for Electricity[3]. The data is partitioned into two groups; countries reporting with a 15 minute frequency are used for day-ahead predictions, while those with hourly reports are used for longer-term predictions. For testing we hold back the German Tennet load recordings from 2015, all of Belgium's recordings of 2016, Austria's load of 2017 and finally the consumption of the German Ampiron grid in 2018.

### 6.2 Day-Ahead Prediction

We start with the task of day-ahead prediction; using 14 days of context, at 12:00, we are asked to predict 24 hours of load from midnight onwards (00:00 until 24:00 o' clock at 12:00 o'clock on the previous day). We therefore forecast the load from noon until midnight on the prediction day plus the next day and ignore the values from the prediction day. We use the same network architecture as in the previous section. During training, the initial learning rate was set to 0.004 and exponentially decayed using a decay of 0.95 every epoch. We train for 80 epochs overall. We compare time domain, windowed time as well as windowed Fourier approaches. The window size was set to 96 which corresponds to 24 hours of data at a sampling rate of 15 minutes per sample.

In Figure 2 we observe that all approaches we tried produce reasonable solutions and outperform the prediction produced by the European Network of Transmission System Operators for Electricity which suggests that their approach could benefit from deep learning.

### 6.3 Long-Term Forecast

Next we consider the more challenging task of 60 day load prediction using 120 days of context. We use 12 load samples per day from all entenso-e member states from 2015 until 2019. We choose a window size of 120 samples or five days; all other parameters are left the same as the day-ahead prediction setting. Table 3 shows that the windowed methods are able to extract patterns, but the scalar-time domain approach failed. Additionally we observe that we do not require the full set of Fourier coefficients to construct useful predictions on the power-data set. The results are tabulated in Table 3. It turns out that the lower quarter of Fourier coefficients is enough, which allowed us to reduce

---

[3]https://transparency.entsoe.eu/; for reproducibility, we will release our crawled version online upon paper acceptance.

Table 3: 60 day ahead power load prediction using GRUs of size 64. We downsample and lowpass-filter with a factor of 1/4. We observe that windowing leads to large training and inference speed-ups. Our STFT approach performs better in the full spectrum case and with a reduced input-dimensionality.

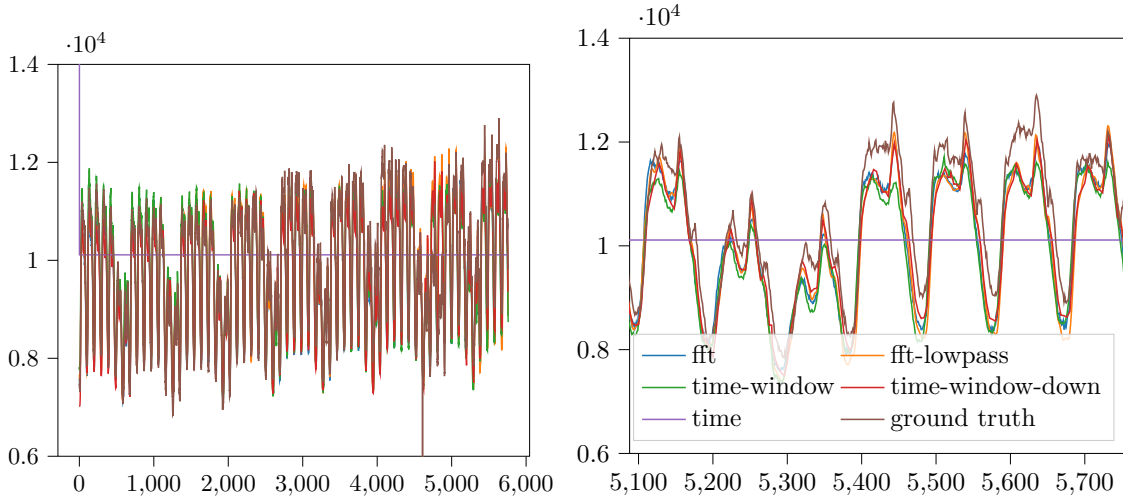| Network | mse [MW]$^2$ | weights | inference [ms] | training-time [min] |
|---|---|---|---|---|
| time-GRU | $261 \cdot 10^5$ | 13k | 1360 | 1472 |
| time-GRU-window | $8.12 \cdot 10^5$ | 74k | 9.25 | 15 |
| time-GRU-window-down | $8.05 \cdot 10^5$ | 28k | 8.2 | 15 |
| STFT-GRU | $7.62 \cdot 10^5$ | 136k | 9.67 | 19 |
| STFT-GRU-lowpass | $7.25 \cdot 10^5$ | 43k | 9.69 | 18 |



Figure 3: A test set sample for showing the 60 day prediction results for all architectures under consideration. Close up for the last week of the 60 day prediction.

the number of parameters considerably. Again we observe that Fourier low-pass filtering outperforms down-sampling and windowing.

# 7 Human Motion Forecasting

## 7.1 Dataset & Evaluation Measure

The Human3.6M data set [30] contains 3.6 million human poses of seven human actors each performing 15 different tasks sampled at 50 Herz. As in previous work [24] we test on the data from actor 5 and train on all others. We work in Cartesian coordinates and predict the absolute 3D-position of 17 joints per frame, which means that we model the skeleton joint movements as well as global translation.

## 7.2 Implementation Details

We use standard GRU-cells with a state size of 3072. We move the data into the frequency domain by computing a short time Fourier transform over each joint dimension separately using a window size of 24. The learning rate was set initially to 0.001 which we reduced using an exponential stair wise decay every thousand iterations by a factor of 0.97. Training was stopped after 19k iterations.

## 7.3 Motion Forecasting Results

Prediction results using our STFT approach are shown in Figure 4. The poses drawn in green and yellow are predictions conditioned on the pose sequences set in blue and red. We observe that the predictions look realistic and smooth. In our experiments low-pass filtering helped us to enforce smoothness in the predictions. Quantitative motion capture prediction results are shown in Table 4. We observe that all windowed approaches run approximately five times faster
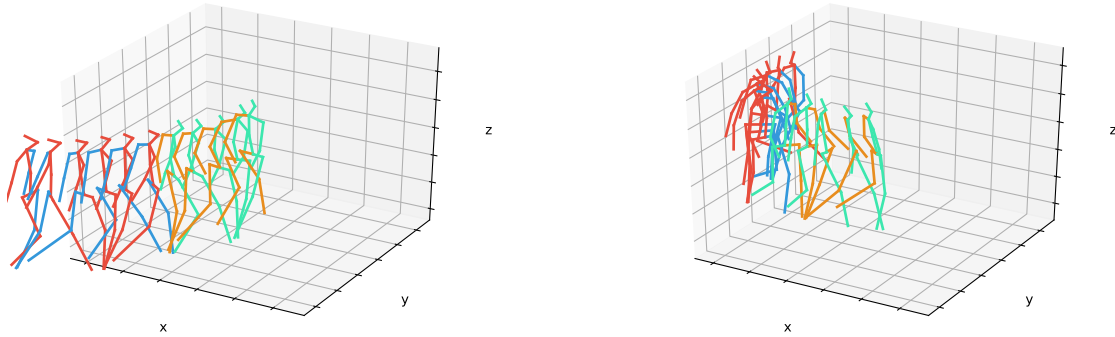
Figure 4: Visualization of input and prediction results using a STFT-RNN combination and low pass filtering. Input is shown in red and blue, predictions in green and yellow.

Table 4: 3d-Human motion forecast of 64 frames or approximately one second. Mean absolute error (mae) is measured in mm. Mean squared errors are reported in mm$^2$. We downsample and lowpass-filter with a factor of 1/4. Windowing runs much faster than the naive time domain approach. Among windowed approaches the STFT allows more aggressive input size reductions.

| Network | mae [mm] | mse [mm]$^2$ | weights | inference [ms] | training-time [min] |
|---|---|---|---|---|---|
| time-GRU [24] | 75.44 | $1.45 \cdot 10^4$ | 29M | 115 | 129 |
| time-GRU-window | 68.25 | $1.47 \cdot 10^4$ | 33M | 18 | 27 |
| time-GRU-window-down | 70.22 | $1.41 \cdot 10^4$ | 30M | 18 | 27 |
| STFT-GRU | 67.88 | $1.25 \cdot 10^4$ | 45M | 25 | 38 |
| STFT-GRU-lowpass | 66.71 | $1.30 \cdot 10^4$ | 32M | 20 | 31 |

than the time domain approach during inference, a significant improvement. In terms of accuracy windowing does comparably well, while the STFT approach does better when the input sampling rate is reduced.

## 8 Conclusion

In this paper we explored frequency domain machine learning using a recurrent neural network. We have proposed to integrate the Short Time Fourier transform and the inverse transform into RNN architectures and evaluated the performance of real and complex valued cells in this setting. We found that complex cells are more parameter efficient, but run slightly longer. Frequency domain RNNs allow us to learn window function parameters and make high-frequency time sequence predictions for both synthetic and real-world data while using less computation time and memory. Low-pass filtering reduced network parameters and outperformed time-domain down-sampling in our experiments.

## References

[1] C. Trabelsi, O. Bilaniuk, Y. Zhang, D. Serdyuk, S. Subramanian, J. F. Santos, S. Mehri, N. Rostamzadeh, Y. Bengio, and C.J. Pal. Deep complex networks. In *ICLR*, 2018.

[2] M. Arjovsky, A. Shah, and Y. Bengio. Unitary evolution recurrent neural networks. In *ICML*, 2016.

[3] S. Wisdom, T. Powers, J.R. Hershey, J. Le Roux, , and L. Atlas. Full-capacity unitary recurrent neural networks. In *NIPS*, 2016.

[4] Moritz Wolter and Angela Yao. Complex gated recurrent neural networks. In *NIPS*, 2018.

[5] Hiroaki Goto and Yuko Osana. Chaotic complex-valued associative memory with adaptive scaling factor independent of multi-values. In *International Conference on Artificial Neural Networks*, pages 76–86. Springer, 2019.

[6] Wei Mao, Miaomiao Liu, Mathieu Salzmann, and Hongdong Li. Learning trajectory dependencies for human motion prediction. In *Proceedings of the IEEE International Conference on Computer Vision*, pages 9489–9497, 2019.

[7] William Chan, Navdeep Jaitly, Quoc Le, and Oriol Vinyals. Listen, attend and spell: A neural network for large vocabulary conversational speech recognition. In *Acoustics, Speech and Signal Processing (ICASSP), 2016 IEEE International Conference on*, pages 4960–4964. IEEE, 2016.

[8] Sander Dieleman and Benjamin Schrauwen. End-to-end learning for music audio. In *Acoustics, Speech and Signal Processing (ICASSP), 2014 IEEE International Conference on*. IEEE, 2014.

[9] Rik Pintelon and Johan Schoukens. *System identification: A frequency domain approach*. John Wiley & Sons, 2012.

[10] K. Minami, H. Nakajima, and T. Toyoshima. Real-time discrimination of ventricular tachyarrhythmia with fourier-transform neural network. *IEEE Transactions on Biomedical Engineering*, 46(2):179–185, Feb 1999.

[11] P. Virtue, S.X. Yu, and M. Lustig. Better than real: Complex-valued neural nets for MRI fingerprinting. In *ICIP*, 2017.

[12] Yoshua Bengio, Yann LeCun, et al. Scaling learning algorithms towards ai. *Large-scale kernel machines*, 34(5):1–41, 2007.

[13] Harry Pratt, Bryan Williams, Frans Coenen, and Yalin Zheng. FCNN: fourier convolutional neural networks. In *Joint European Conference on Machine Learning and Knowledge Discovery in Databases*, pages 786–798. Springer, 2017.

[14] NVIDIA Developer Fourier transforms in convolution. `https://developer.nvidia.com/discover/convolution`. Accessed: 2020-03-25.

[15] Zelong Wang, Qiang Lan, Hongjun He, and Chunyuan Zhang. Winograd algorithm for 3D convolution neural networks. In *ICANN*, 2017.

[16] Wei Zuo, Yang Zhu, and Lilong Cai. Fourier-neural-network-based learning control for a class of nonlinear systems with flexible components. *IEEE transactions on neural networks*, 20(1):139–151, 2008.

[17] L. B. Godfrey and M. S. Gashler. Neural decomposition of time-series data for effective generalization. *IEEE Transactions on Neural Networks and Learning Systems*, 29(7):2973–2985, July 2018.

[18] Jongbin Ryu, Ming-Hsuan Yang, and Jongwoo Lim. Dft-based transformation invariant pooling layer for visual classification. In *ECCV*, 2018.

[19] Tsungnan Lin, Bill G Horne, Peter Tino, and C Lee Giles. Learning long-term dependencies in narx recurrent neural networks. *IEEE Transactions on Neural Networks*, 7(6):1329–1338, 1996.

[20] Tayfun Alpay, Stefan Heinrich, and Stefan Wermter. Learning multiple timescales in recurrent neural networks. In *International Conference on Artificial Neural Networks*, pages 132–139. Springer, 2016.

[21] Jan Koutnik, Klaus Greff, Faustino Gomez, and Juergen Schmidhuber. A clockwork rnn. In *International Conference on Machine Learning*, pages 1863–1871, 2014.

[22] Junyoung Chung, Sungjin Ahn, and Yoshua Bengio. Hierarchical multiscale recurrent neural networks. In *ICLR*, 2017.

[23] Iulian Vlad Serban, Tim Klinger, Gerald Tesauro, Kartik Talamadupula, Bowen Zhou, Yoshua Bengio, and Aaron Courville. Multiresolution recurrent neural networks: An application to dialogue response generation. In *Thirty-First AAAI Conference on Artificial Intelligence*, 2017.

[24] J. Martinez, M.J. Black, and J. Romero. On human motion prediction using recurrent neural networks. In *CVPR*, 2017.

[25] J. Lim D. Griffin. Signal estimation from modified short-time fourier transform. In *IEEE Transactions on Acoustics, Speech, and Signal Processing*, 1984.

[26] Fredric J Harris. On the use of windows for harmonic analysis with the discrete fourier transform. *Proceedings of the IEEE*, 66(1):51–83, 1978.

[27] Karlheinz Gröchenig. *Foundations of time-frequency analysis*. Springer Science & Business Media, 2013.

[28] Felix A Gers, Douglas Eck, and Jürgen Schmidhuber. Applying lstm to time series predictable through time-window approaches. In *Neural Nets WIRN Vietri-01*, pages 193–200. Springer, 2002.

[29] K. Cho, B. van Merriënboer, Ç. Gülçehre, D. Bahdanau, F. Bougares, H. Schwenk, and Y. Bengio. Learning phrase representations using RNN encoder–decoder for statistical machine translation. In *EMNLP*, 2014.

[30] Catalin Ionescu, Dragos Papava, Vlad Olaru, and Cristian Sminchisescu. Human3.6m: Large scale datasets and predictive methods for 3d human sensing in natural environments. *IEEE Transactions on Pattern Analysis and Machine Intelligence*, 36(7):1325–1339, jul 2014.