

# Machine Learning in the Frequency Domain

Moritz Wolter

*wolter@cs.uni-bonn.de*

July 24, 2019



# The Fourier transform

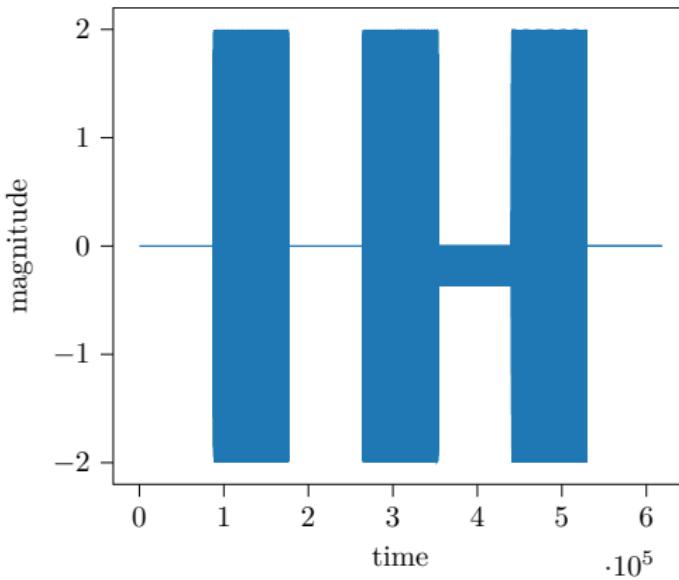


## An analog dial pad

	1209 Hz	1336 Hz	1477 Hz
697 Hz	1	2	3
770 Hz	4	5	6
852 Hz	7	8	9
941 Hz	*	0	#

dial

## The time domain signal



**Figure:** The time domain signal sampled at 44.1khz

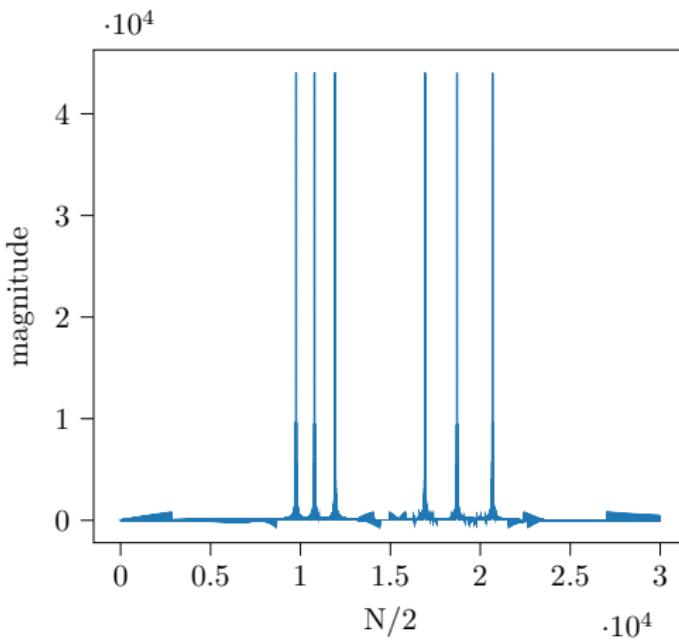
# The Fourier Transform

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

## Euler's formula:

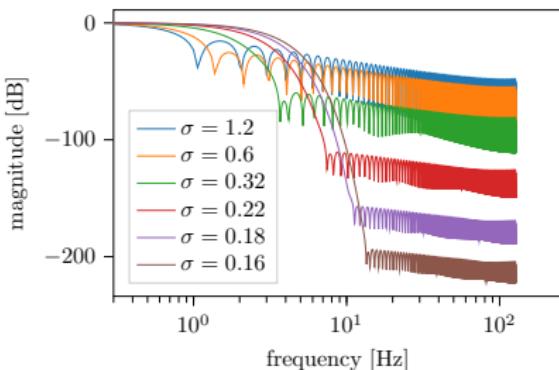
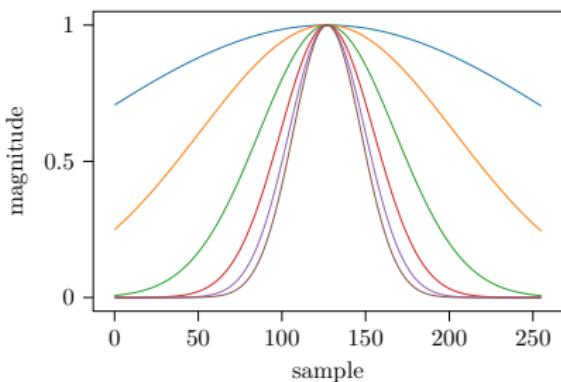
$$e^{j\omega} = \cos(\omega) + j\sin(\omega) \quad (2)$$

# The fourier transfrom applied

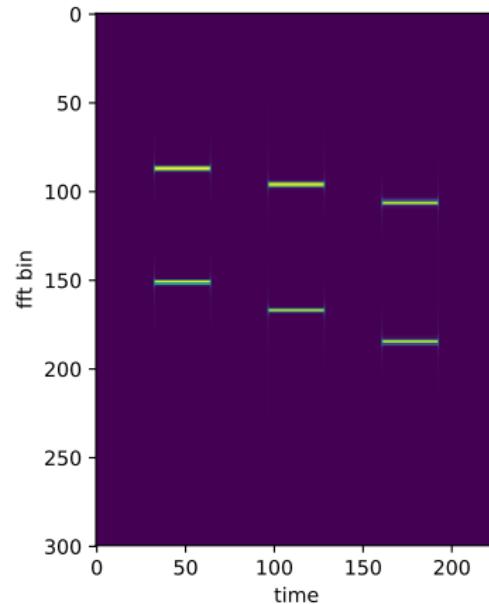


# The Short time Fourier transform

$$\mathbf{X}[\omega, Sm] = \mathcal{F}_s(\mathbf{x}) = \mathcal{F}(\mathbf{w}[Sm-l]\mathbf{x}[l]) = \sum_{l=-\infty}^{\infty} \mathbf{w}[Sm-l]\mathbf{x}[l]e^{-j\omega l}, \quad (3)$$

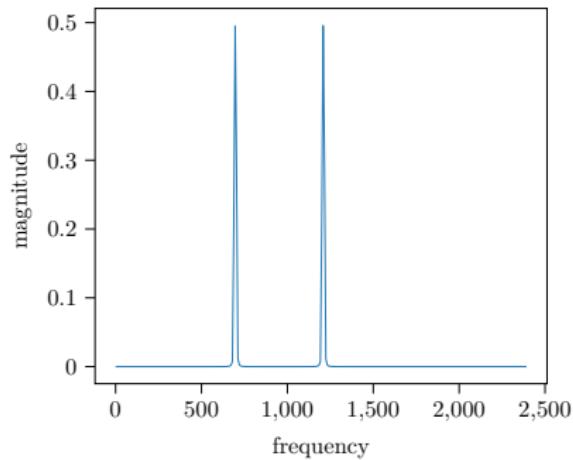


# Short Time Fourier Transform Magnitude



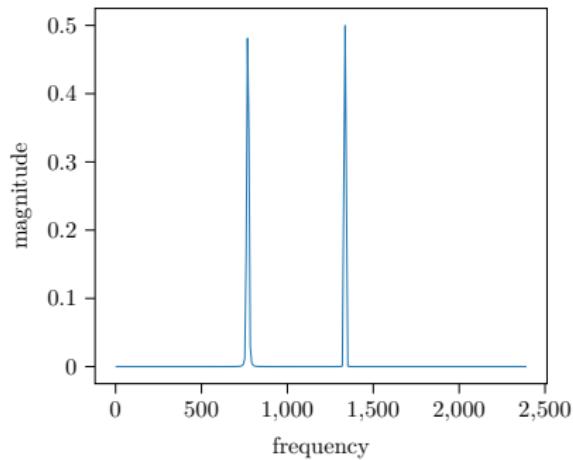
# Key 1

STFT during key 1



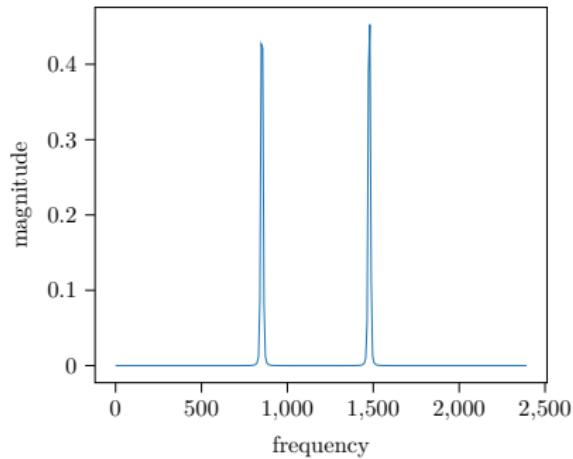
## Key 2

STFT during key 2



# Key 3

STFT during key 3



# The uncertainty principle <sup>1</sup>

Example 1:

- 80ms
- 170ms
- 330ms
- 670ms
- 1s
- 2s
- 5s

---

<sup>1</sup><http://newt.phys.unsw.edu.au/jw/uncertainty.html>

# The uncertainty principle 2<sup>2</sup>

Example 2:

- 80ms
- 170ms
- 330ms
- 670ms
- 1s
- 2s
- 5s

---

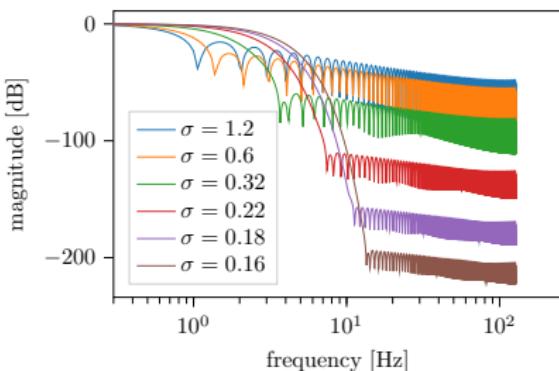
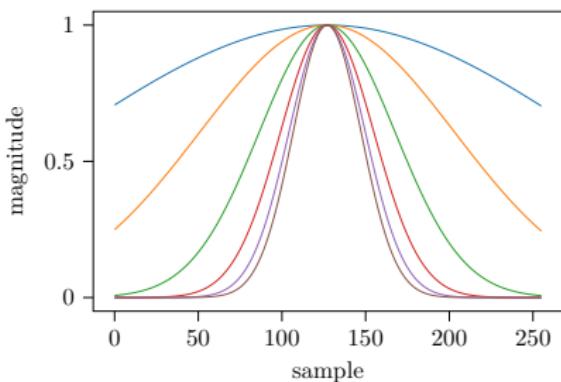
<sup>2</sup><http://newt.phys.unsw.edu.au/jw/uncertainty.html>

## Learning through the STFT

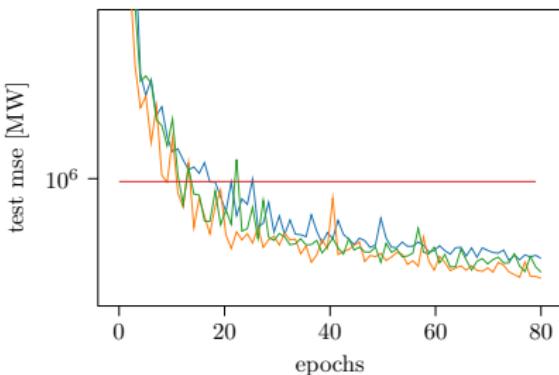
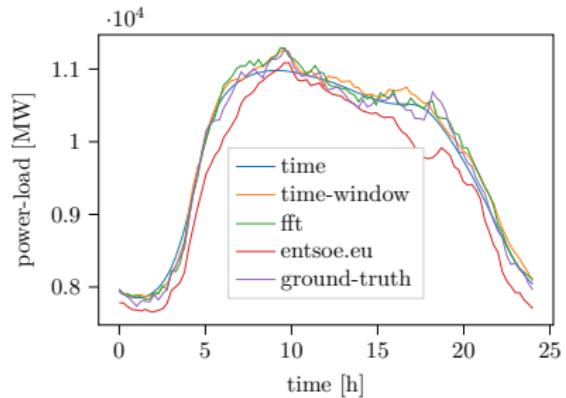
- The first example consisted of a 400 and 403Hz Sine wave.
- The second of a 400 and 401Hz Sine wave.
- Time and frequency resolution are coupled through the uncertainty principle.
- Working in the time domain can overload RNNs for long sequences.
- Transfer a signal into the Frequency domain do the prediction and compute the inverse transform.
- IDEA: Learn the window shape.

# Learning through the STFT

$$\mathbf{X}[\omega, Sm] = \mathcal{F}_s(\mathbf{x}) = \mathcal{F}(\mathbf{w}[Sm - l]\mathbf{x}[l]) = \sum_{l=-\infty}^{\infty} \mathbf{w}[Sm - l]\mathbf{x}[l]e^{-j\omega l}, \quad (4)$$

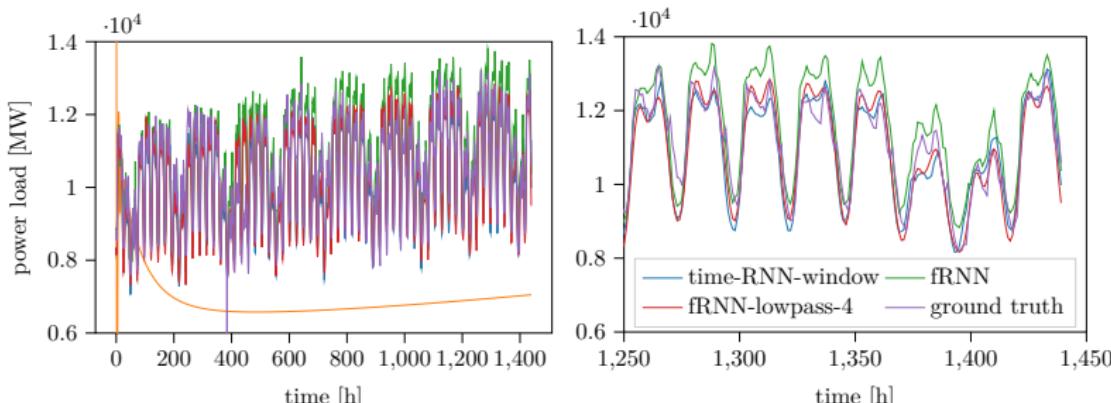


# Learning through the STFT [WY18b]



# Learning through the STFT [WY18b]

Network	mse [MW]	weights	run [min]
time-RNN	$1.3 \cdot 10^7$	13k	772
time-RNN-windowed	$8.8 \cdot 10^5$	28k	12
fRNN	$8.3 \cdot 10^5$	44k	13
fRNN-lowpass-1/4	$7.6 \cdot 10^5$	20k	13
fRNN-lowpass-1/8	$1.3 \cdot 10^6$	16k	13



# Complex Machine learning

- Quantum computers require complex unitary weights.
- Fourier transforms produce complex representations.
- Encoding data in magnitude and phase may enable us create a richer representation.
- Complex analysis is a well studied (and very interesting!) subject, lets merge it with machine learning and see what happens.

# Memory and adding benchmark problems for RNNs

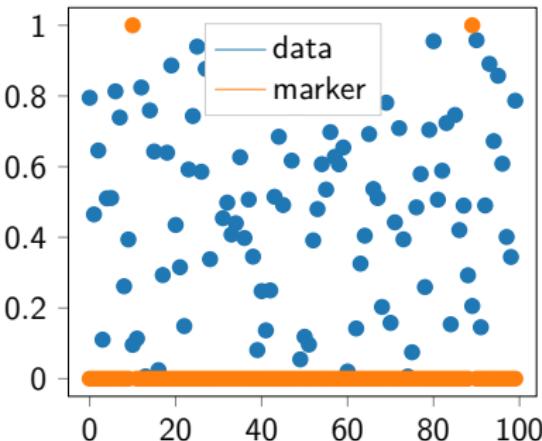
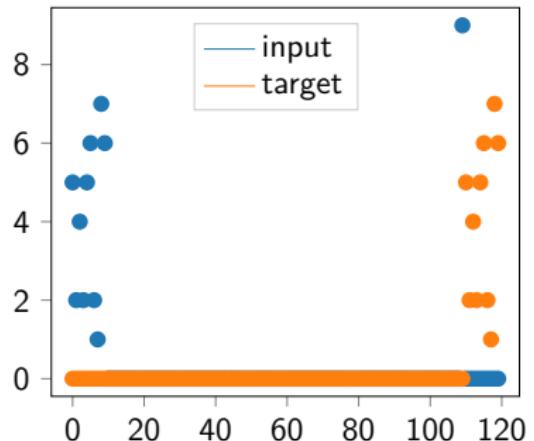


Figure: Illustrations of the memory problem on the left and the adding problem on the right.

## Wirtinger-Calculus [Wir27][MG09][KD09]

For a complex function  $f(z) = u(x, y) - iv(x, y)$  we have:

$$\mathbb{R}\text{-derivative} \triangleq \frac{\partial f}{\partial z}|_{\bar{z}=\text{const}} = \frac{1}{2}\left(\frac{\partial f}{\partial x} - i\frac{\partial f}{\partial y}\right), \quad (5)$$

$$\overline{\mathbb{R}}\text{-derivative} \triangleq \frac{\partial f}{\partial \bar{z}}|_{z=\text{const}} = \frac{1}{2}\left(\frac{\partial f}{\partial x} + i\frac{\partial f}{\partial y}\right). \quad (6)$$

Based on these derivatives, one can define the chain rule for a function  $g(f(z))$  as follows:

$$\frac{\partial g(f(z))}{\partial z} = \frac{\partial g}{\partial f} \frac{\partial f}{\partial z} + \frac{\partial g}{\partial \bar{f}} \frac{\partial \bar{f}}{\partial z} \text{ where } \bar{f} = u(x, y) - iv(x, y). \quad (7)$$

Theoretical tool to convince ourselves, that it's ok to work with equivalent real networks.

## Unitary Evolution matrix RNN-Motivation [ASB16][Pas13]

$$\mathbf{x}_t = \mathbf{W}_{\text{rec}} f(\mathbf{x}_{t-1}) + \mathbf{W}_{\text{in}} \mathbf{u}_t + \mathbf{b}. \quad (8)$$

$$\frac{\partial \mathcal{E}}{\partial \theta} = \sum_{1 \leq t \leq T} \frac{\mathcal{E}_t}{\partial \theta}, \quad (9)$$

$$\frac{\partial \mathcal{E}_t}{\partial \theta} = \sum_{1 \leq k \leq t} \left( \frac{\mathcal{E}_t}{\partial \mathbf{x}_t} \frac{\partial \mathbf{x}_t}{\mathbf{x}_k} \frac{\partial^+ \mathbf{x}_k}{\partial \theta} \right), \quad (10)$$

$$\frac{\partial \mathbf{x}_t}{\partial \mathbf{x}_k} = \prod_{t \geq i > k} \frac{\partial \mathbf{x}_i}{\partial \mathbf{x}_{i-1}} = \prod_{t \geq i > k} W_{\text{rec}}^T \text{diag}(f'(\mathbf{x}_{i-1})). \quad (11)$$

# Stiefel Manifold Weight Updates [WPH<sup>+</sup>16]

$$\mathbf{W}_{k+1} = (\mathbf{I} + \frac{\lambda}{2} \mathbf{A}_k)^{-1} (\mathbf{I} - \frac{\lambda}{2} \mathbf{A}_k) \mathbf{W}_k, \quad (12)$$

$$\text{where } \mathbf{A} = \mathbf{W} \bar{\nabla}_{\mathbf{w}} F^T - \bar{\mathbf{W}}^T \nabla_{\mathbf{w}} F. \quad (13)$$

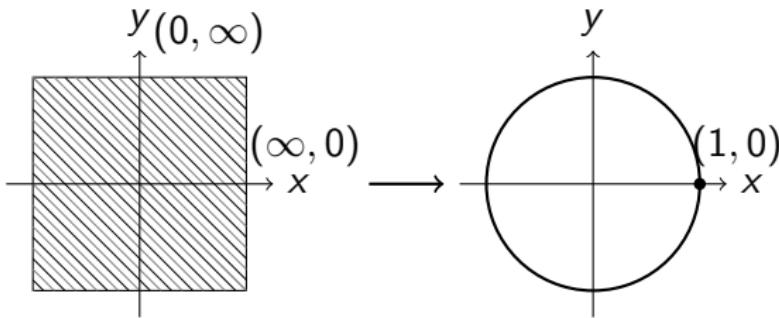
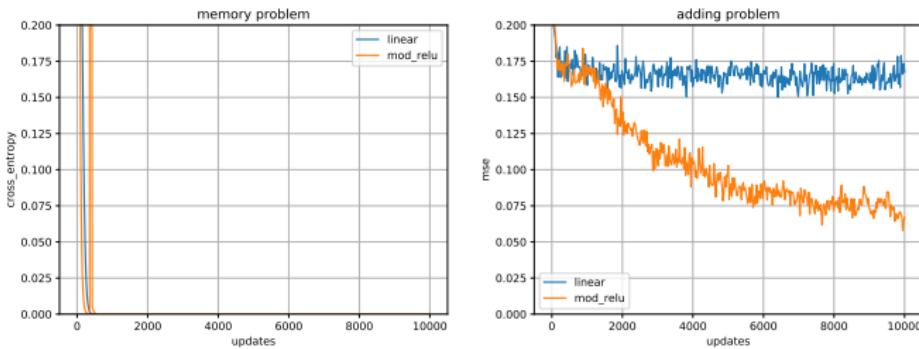


Figure: Fix the optimized matrix eigenvalues onto the unit circle. The key idea behind stiefel-manifold optimization.

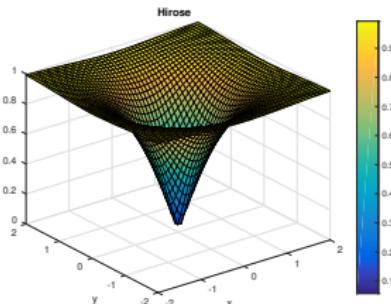
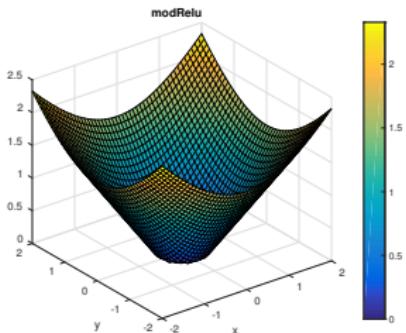
# The linear unitary case

$$\mathbf{x}_t = \mathbf{W}_{\text{rec}} \mathbf{x}_{t-1} + \mathbf{W}_{\text{in}} \mathbf{u}_t + \mathbf{b}. \quad (14)$$



**Figure:** Performance of linear and mod-Relu activated unitary RNNs on the memory (left) and adding (right) problems for T=50. All networks have approx. 40k weights.

# Complex equivalents of tanh and Relu



$$f_{\text{Hirose}}(z) = \tanh\left(\frac{|z|}{m^2}\right) e^{-i \cdot \theta_z} = \tanh\left(\frac{|z|}{m^2}\right) \frac{z}{|z|}, \quad (15)$$

$$f_{\text{modReLU}}(z) = \text{ReLU}(|z| + b) e^{-i \cdot \theta_z} = \text{ReLU}(|z| + b) \frac{z}{|z|}. \quad (16)$$

We will compare their performance as state-to-state non-linearities.

# Unitary evolution network performance

$$\mathbf{x}_t = \mathbf{U}_{\text{rec}} f(\mathbf{x}_{t-1}) + \mathbf{W}_{\text{in}} \mathbf{u}_t + \mathbf{b}. \quad (17)$$

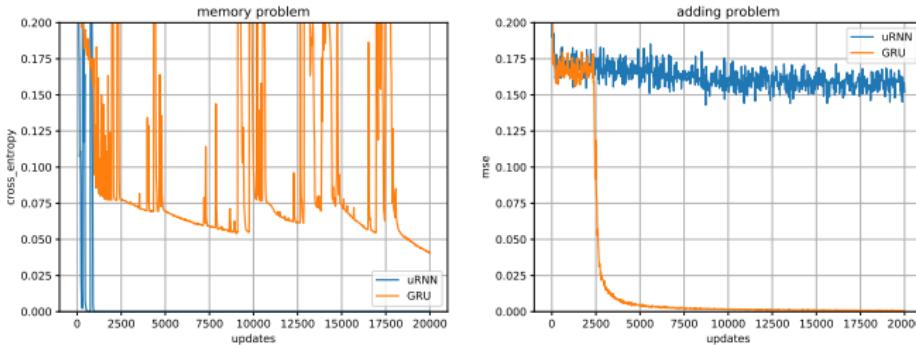
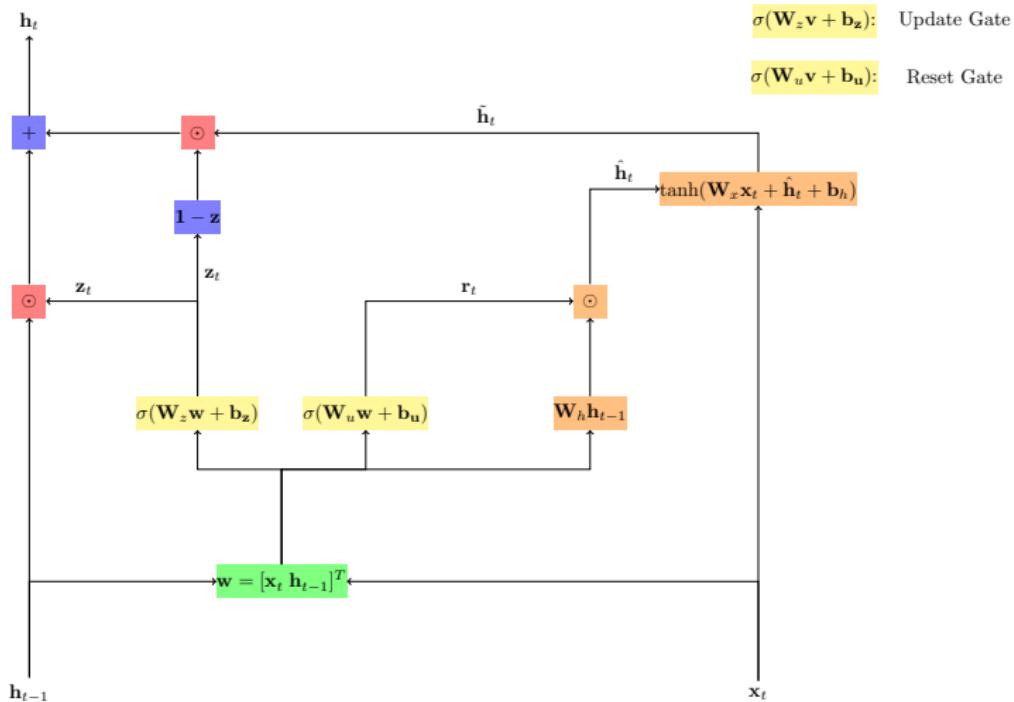


Figure: Current state of the art performance on memory and adding problem for  $T=250$ . Models have approximately 40k weights.

# The gated recurrent unit



## Complex gated Recurrent Recurrent Nets [WY18a]

Gate equation:

$$\mathbf{g}_r = f_g(\mathbf{z}_r), \quad \text{where} \quad \mathbf{z}_r = \mathbf{W}_r \mathbf{h} + \mathbf{V}_r \mathbf{x}_t + \mathbf{b}_r, \quad (18)$$

$$\mathbf{g}_z = f_g(\mathbf{z}_z), \quad \text{where} \quad \mathbf{z}_z = \mathbf{W}_z \mathbf{h} + \mathbf{V}_z \mathbf{x}_t + \mathbf{b}_z, \quad (19)$$

Update equations:

$$\tilde{\mathbf{z}}_t = \mathbf{W}(\mathbf{g}_r \odot \mathbf{h}_{t-1}) + \mathbf{V}\mathbf{x}_t + \mathbf{b}, \quad (20)$$

$$\mathbf{h}_t = \mathbf{g}_z \odot f_a(\tilde{\mathbf{z}}_t) + (1 - \mathbf{g}_z) \odot \mathbf{h}_{t-1}, \quad (21)$$

$\mathbb{C} \rightarrow \mathbb{R}$ , mapping:

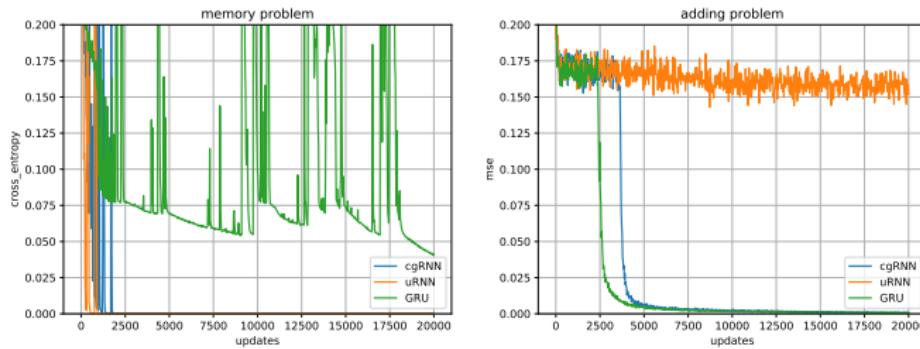
$$\mathbf{o}_r = \mathbf{W}_o[\Re(\mathbf{h}) \ \Im(\mathbf{h})] + \mathbf{b}_o. \quad (22)$$

# Complex gate activations

$$f_{\text{mod sigmoid}}(\mathbf{z}) = \sigma(\alpha \Re(\mathbf{z}) + \beta \Im(\mathbf{z})). \quad (23)$$

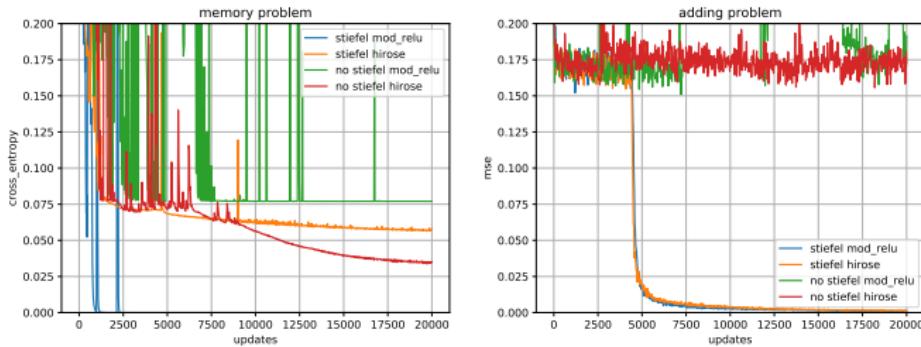
With  $\alpha \in [0, 1]$  and  $\beta = (1 - \alpha)$ .

# Comparison to state of the art



**Figure:** Comparison of our complex gated RNN (cgRNN, blue,  $n_h=80$ ) with the unitary RNN [ASB16](uRNN, orange,  $n_h=140$ ) and standard GRU [CvMG<sup>+</sup>14](orange,  $n_h=112$ ) on the memory (left) and adding (right) problem for  $T=250$ .

# Stiefel optimization and activations



**Figure:** Comparison of non-linearities and norm preserving state transition matrices on the complex gated RNNs for the memory (a) and adding (b) problems for  $T=250$ . We use  $n_h = 80$  for all experiments.

# Weight reductions on mocap data

Table 2: Comparison of our cgRNN with the GRU [28] on human motion prediction.

Action	cgRNN				GRU[28]			
	80ms	160 ms	320ms	400ms	80ms	160ms	320ms	400ms
walking	0.29	0.48	0.74	0.84	<b>0.27</b>	<b>0.47</b>	<b>0.67</b>	<b>0.73</b>
eating	0.23	<b>0.38</b>	0.66	0.82	0.23	0.39	<b>0.62</b>	<b>0.77</b>
smoking	<b>0.31</b>	<b>0.58</b>	<b>1.01</b>	<b>1.1</b>	0.32	0.6	1.02	1.13
discussion	0.33	0.72	<b>1.02</b>	<b>1.08</b>	<b>0.31</b>	<b>0.7</b>	1.05	1.12
directions	0.41	<b>0.65</b>	<b>0.83</b>	<b>0.93</b>	<b>0.41</b>	0.65	0.83	0.96
greeting	0.53	0.87	<b>1.26</b>	<b>1.43</b>	<b>0.52</b>	<b>0.86</b>	1.30	1.47
phoning	<b>0.58</b>	1.09	1.57	1.72	0.59	<b>1.07</b>	<b>1.50</b>	<b>1.67</b>
posing	<b>0.37</b>	<b>0.72</b>	<b>1.38</b>	<b>1.65</b>	0.64	1.16	1.82	2.1
purchases	0.61	0.86	<b>1.21</b>	1.31	<b>0.6</b>	<b>0.82</b>	1.13	<b>1.21</b>
sitting	0.46	0.75	1.22	<b>1.44</b>	<b>0.44</b>	<b>0.73</b>	<b>1.21</b>	1.45
sitting down	0.55	1.02	1.54	1.73	<b>0.48</b>	<b>0.89</b>	<b>1.36</b>	<b>1.57</b>
taking photo	<b>0.29</b>	<b>0.59</b>	<b>0.92</b>	<b>1.07</b>	0.29	0.59	0.95	1.1
waiting	0.35	0.68	1.16	<b>1.36</b>	<b>0.33</b>	<b>0.65</b>	<b>1.14</b>	1.37
walking dog	0.57	1.09	1.45	1.55	<b>0.54</b>	<b>0.94</b>	<b>1.32</b>	<b>1.49</b>
walking together	<b>0.27</b>	<b>0.53</b>	<b>0.77</b>	<b>0.86</b>	0.28	0.56	0.8	0.88
average	<b>0.41</b>	<b>0.73</b>	<b>1.12</b>	<b>1.26</b>	0.42	0.74	<b>1.12</b>	1.27

Our cgRNN ( $n_h = 512$ , 1.8M params) predicts human motions which are either comparable or slightly better than the real-valued GRU [28] ( $n_h = 1024$ , 3.4M params) despite having only approximately half the parameters.

## Upcoming: Temporal Convolutions in the frequency domain

- The STFT turns sequences of numbers into images.
- What happens if we convolve these in time and frequency?
- What is the best way to convolve in frequency and time?
- How can recurrent connections and convolutions best be used together?

## References I

-  Martin Arjovsky, Amar Shah, and Yoshua Bengio, *Unitary evolution recurrent neural networks*, ICML, 2016.
-  Kyunghyun Cho, Bart van Merriënboer, Çağlar Gülcühre, Dzmitry Bahdanau, Fethi Bougares, Holger Schwenk, and Yoshua Bengio, *Learning phrase representations using RNN encoder–decoder for statistical machine translation*, EMNLP, October 2014.
-  Ken Kreutz-Delgado, *The complex gradient operator and the cr-calculus*, arXiv preprint arXiv:0906.4835 (2009).

## References II

-  Danilo P Mandic and Vanessa Su Lee Goh, *Complex valued nonlinear adaptive filters: noncircularity, widely linear and neural models*, vol. 59, John Wiley & Sons, 2009.
-  Pascanu, *On the difficulty of training recurrent neural networks*, Journal of Machine Learning Research (2013).
-  W. Wirtinger, *Zur formalen theorie der funktionen von mehr komplexen veränderlichen*, 1927.
-  Scott Wisdom, Thomas Powers, John R. Hershey, Jonathan Le Roux, , and Les Atlas, *Full-capacity unitary recurrent neural networks*, Advances in Neural Information Processing Systems, 2016.

## References III

-  Moritz Wolter and Angela Yao, *Complex gated recurrent neural networks*, 32nd Conference on Neural Information Processing Systems, 2018.
-  \_\_\_\_\_, *Fourier rnns for sequence prediction*, arXiv preprint arXiv:1812.05645, 2018.

# Discussion

Thanks for your attention and feedback.  
Feel free to contact me at: [wolter@cs.uni-bonn.de](mailto:wolter@cs.uni-bonn.de)