

The background of the slide features a complex network diagram. It consists of numerous nodes of varying sizes and colors (gray, blue, and dark blue) connected by thin gray lines. Some nodes are highlighted with larger, semi-transparent circles. The overall aesthetic is technical and modern.

# AN EMPIRICAL EVALUATION OF GENERIC CONVOLUTIONAL AND RECURRENT NETWORKS FOR SEQUENCE MODELING

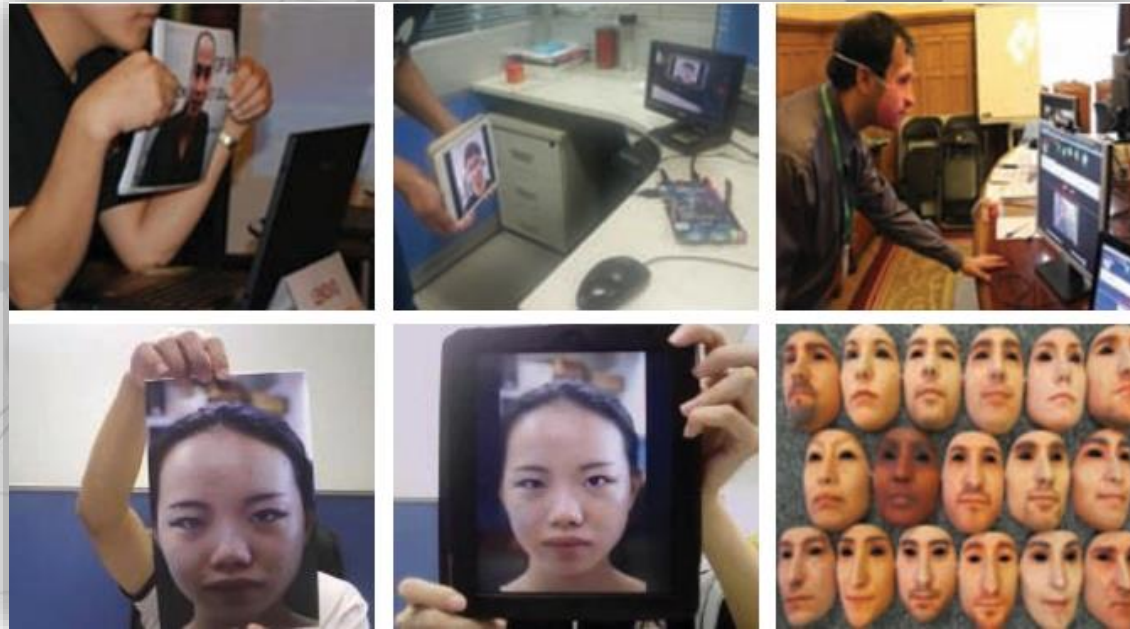
Nova Search Reading Group – June 4th 2020

Ruslan Padnevykh

**Shaojie Bai, J. Zico Kolter and Vladlen Koltun. April 2018**

<https://arxiv.org/pdf/1803.01271.pdf>

# IDENTITY THEFT



**Different identity theft techniques**

# IDENTITY THEFT (CONT.)

Photos taken from Google



In 2011, a passenger boarded the plane in Hong Kong wearing an elderly man's mask and successfully landed in Canada.



# CONTEXT: BORDERS WITHOUT BARRIERS

## Use case



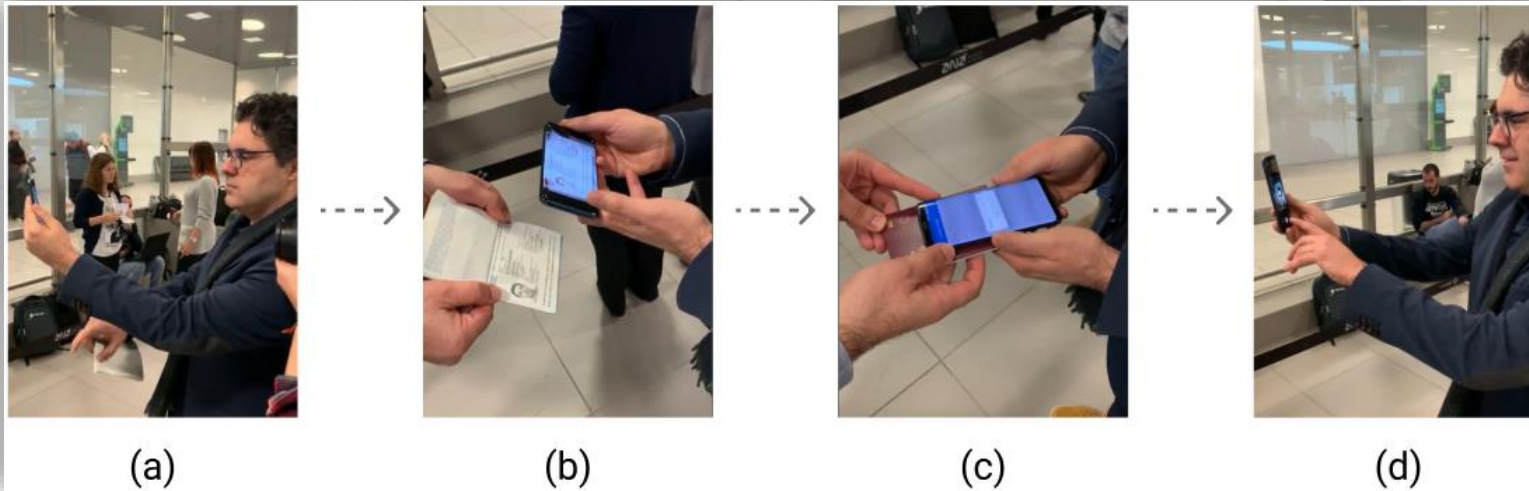
***Biometrics on the Move (SmartyFlow project)***



**Queues to show documents to border guards**

# CONTEXT: BORDERS WITHOUT BARRIERS

## Registration Phase



(a) → **liveliness detection (5s video)**

(b) - passport registration

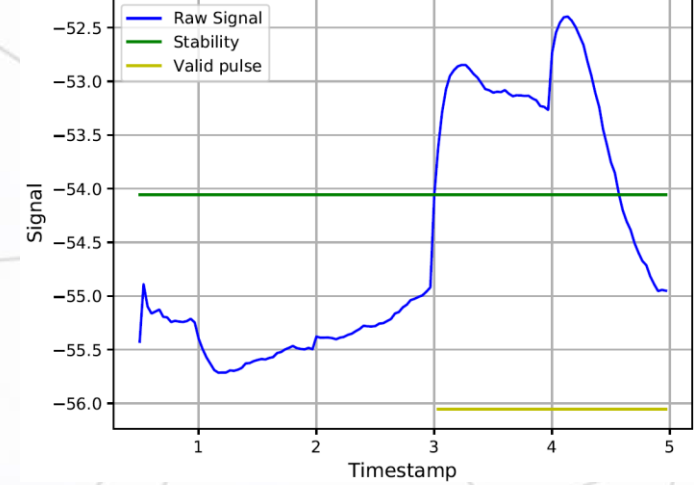
(c) - reading the passport chip

(d) - verification of correspondence

# GOAL

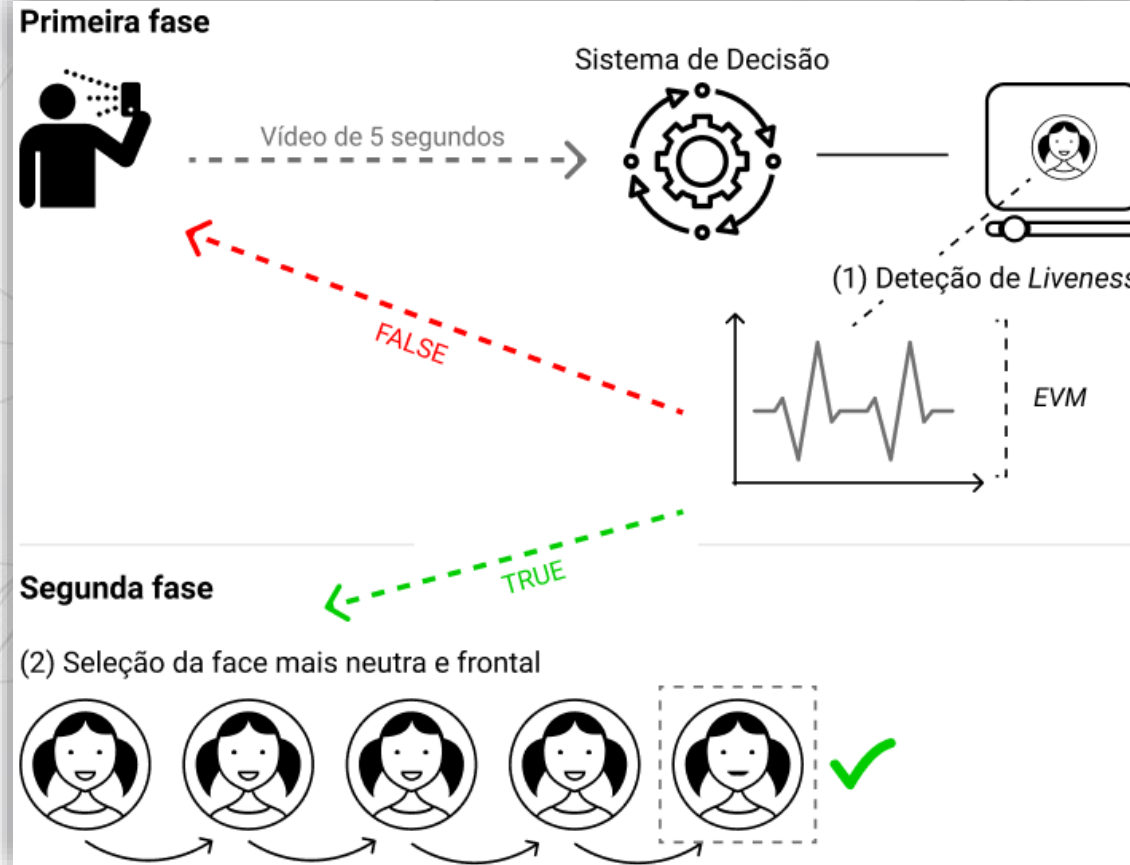
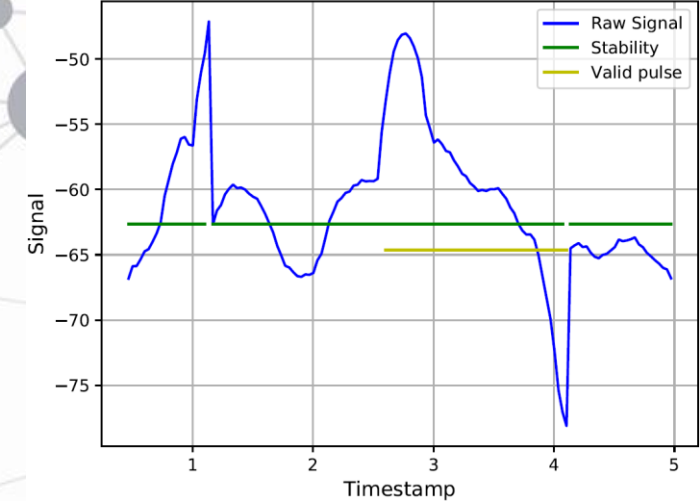
**Fake**

Ambiente\_mais\_claro: 010-015 secs



**Real**

Andar\_rapido: 010-015 secs



Liveness detection

Neutral face selection



# SEQUENCE MODELING

- ❖ Suppose that we are given an input sequence  $x_0, \dots, x_T$ , and wish to predict some corresponding outputs  $y_0, \dots, y_T$  at each time.
- ❖ The key constraint is that to predict the output  $y_t$  for some time  $t$ , we are constrained to only use those inputs that have been previously observed:  $x_0, \dots, x_t$ .
- ❖ Formally, a sequence modeling network is any function:

$$f : \mathcal{X}^{T+1} \rightarrow \mathcal{Y}^{T+1} \quad \text{that produces the mapping} \quad \hat{y}_0, \dots, \hat{y}_T = f(x_0, \dots, x_T)$$

**$y_t$  depends only on  $x_0, \dots, x_t$  and not on any “future” inputs  $x_{t+1}, \dots, x_T$**

# SEQUENCE MODELING (CONT.)

The goal of learning in the sequence modeling setting is to find a network  $f$  that minimizes some expected loss between the actual outputs and the predictions,

$$L(y_0, \dots, y_T ; f(x_0, \dots, x_T)),$$

where the sequences and outputs are drawn according to some distribution.



# BACKGROUND

## Convolutional Networks

- ❖ Used prominently for speech recognition in the 80s and 90s
- ❖ Applied to NLP tasks such as part-of-speech tagging and semantic role labelling
- ❖ Applied to sentence and document classification
- ❖ More recently, applied to **language modeling and machine translation**

## Recurrent Networks

- ❖ Gained tremendous popularity due to prominent applications to **language modeling and machine translation**

**Should Recurrent Networks be regarded as a natural starting point for sequence modeling tasks?**

# RECURRENT NETWORK (RNN)

- ❖ RNN is a class of artificial neural networks where connections between nodes form a directed graph along a temporal sequence (Rumelhart, 1988).
- ❖ Unlike feedforward neural networks, RNNs contain cycles and use an **internal state memory**  $h$  to process sequences of inputs.
- ❖ A basic recurrent neural network is described by the propagation equations:

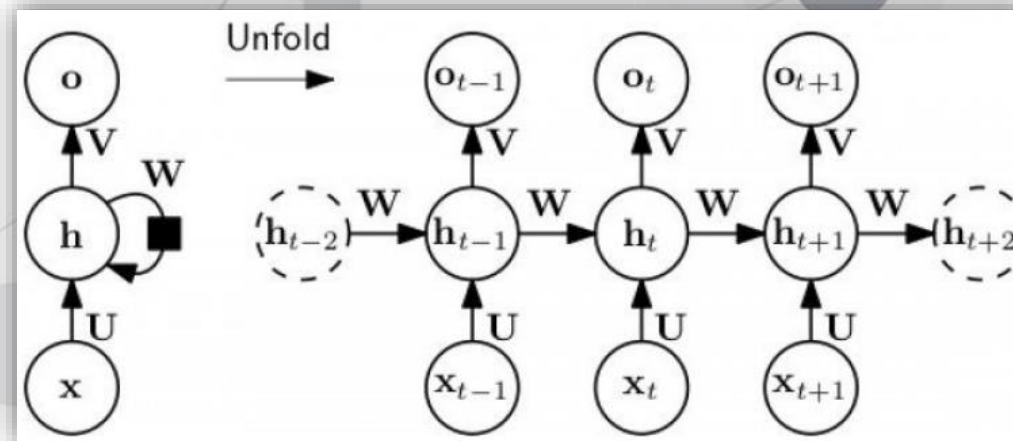
$$\begin{aligned} \mathbf{h}_t &= \sigma(\mathbf{U} \cdot \mathbf{x}_t + \mathbf{W} \cdot \mathbf{h}_{t-1} + \mathbf{b}) \\ \mathbf{o}_t &= \mathbf{V} \cdot \mathbf{h}_t + \mathbf{c} \end{aligned}$$

- ❖ Where the parameters are the bias vectors  $\mathbf{b}$  and  $\mathbf{c}$  along with the weight matrices:
  - ❖  $\mathbf{U}$  – input-to-hidden
  - ❖  $\mathbf{V}$  – hidden-to-output
  - ❖  $\mathbf{W}$  – hidden-to-hidden



# RNN (CONT.)

❖ The computational graph and its unfolded version is shown in the following figure:



- ❖ Weight matrices:
- ❖  $U$  – input-to-hidden
- ❖  $V$  – hidden-to-output
- ❖  $W$  – hidden-to-hidden

- ❖ Computing the gradients involves performing a forward propagation pass through the unrolled graph followed by a backward propagation pass.
- ❖ The runtime is  $O(T)$  and cannot be reduced by parallelization because the forward propagation graph is inherently sequential, i.e., each time step may be computed only after the previous one.

# RNN (CONT.)

- ❖ Recurrent models construct very deep computational graphs by repeatedly applying the same operation at each time step of a long temporal sequence. This gives rise to the vanishing gradient problem and makes it notoriously difficult to train RNNs.
- ❖ To prevent these difficulties more elaborate recurrent architectures were developed, such as the long short-term memory (LSTM) (Hochreiter, 1997) and the gated recurrent unit (GRU) (Cho, 2014).

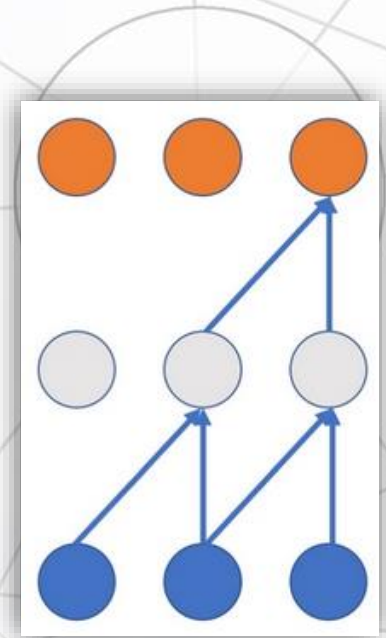
# TEMPORAL CONVOLUTIONAL NETWORK (TCN)

- ❖ TCN is inspired by recent convolutional architectures for sequential data and combines simplicity , autoregressive prediction, and very long memory.
- ❖ The distinguishing characteristics of TCNs are:
  - 1) The **convolutions in the architecture are causal**, meaning that there is no information “leakage” from future to past.
  - 2) The architecture can take a sequence of any length and map it to an output sequence of the same length, just as with an RNN.



# TCN — CAUSAL CONVOLUTIONS

- 1) TCN uses **causal convolutions**, i.e., convolutions where an output at time  $t$  is convolved only with elements from time  $t$  and earlier in the previous layer.
- 2) TCN uses a **1D fully-convolutional network (FCN)** architecture, where each hidden layer is the same length as the input layer.



**TCN = 1D FCN + causal convolutions**

## TCN (CONT.)

A major **disadvantage** of this basic design is that in order to achieve a long effective history size, we need an extremely deep network or very large filters.

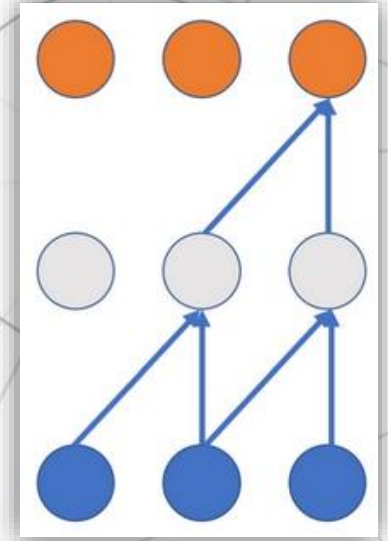
# TCN — DILATED CONVOLUTIONS

- ❖ A simple causal convolution is only able to look back at a history with size linear in the depth of the network.



Solution

To employ dilated convolutions that enable an exponentially large receptive field.





# TCN – DILATED CONVOLUTIONS

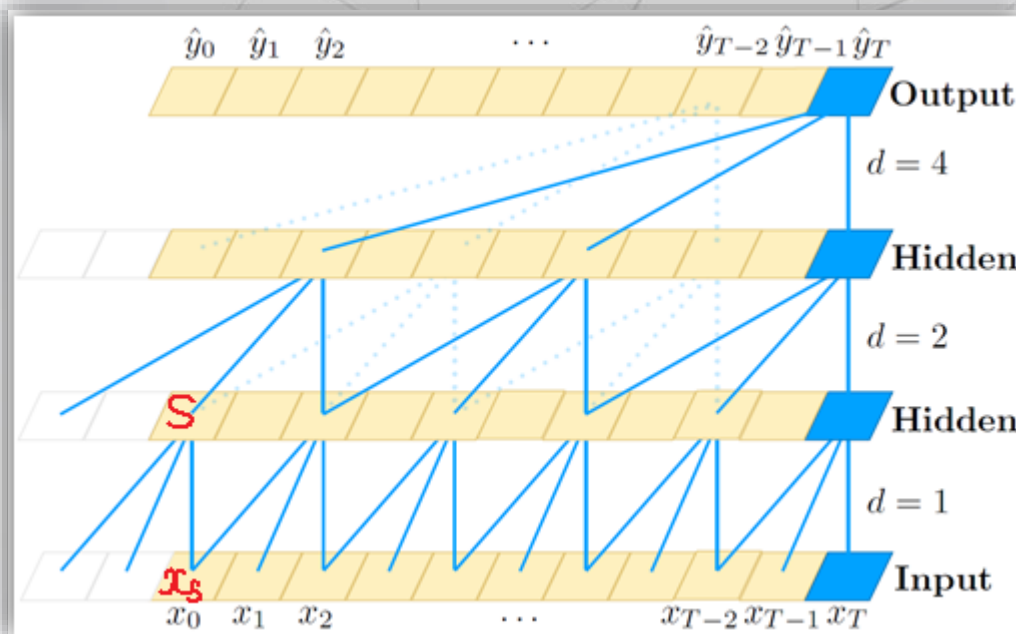
Given a 1-D sequence input  $\mathbf{x} \in \mathbb{R}^n$  and a filter  $f : \{f_0, \dots, f_{k-1}\} \rightarrow \mathbb{R}$ , the dilated convolution operation  $F$  on element  $s$  of the sequence is defined as:

$$F(s) = (\mathbf{x} *_d f)(s) = \sum_{i=0}^{k-1} f(i) \cdot \mathbf{x}_{s-d \cdot i}$$

$d = 2^v$  – dilation factor, with  $v$  the level of the network

$k$  – filter size

$(s - d \cdot i)$  – direction of the past



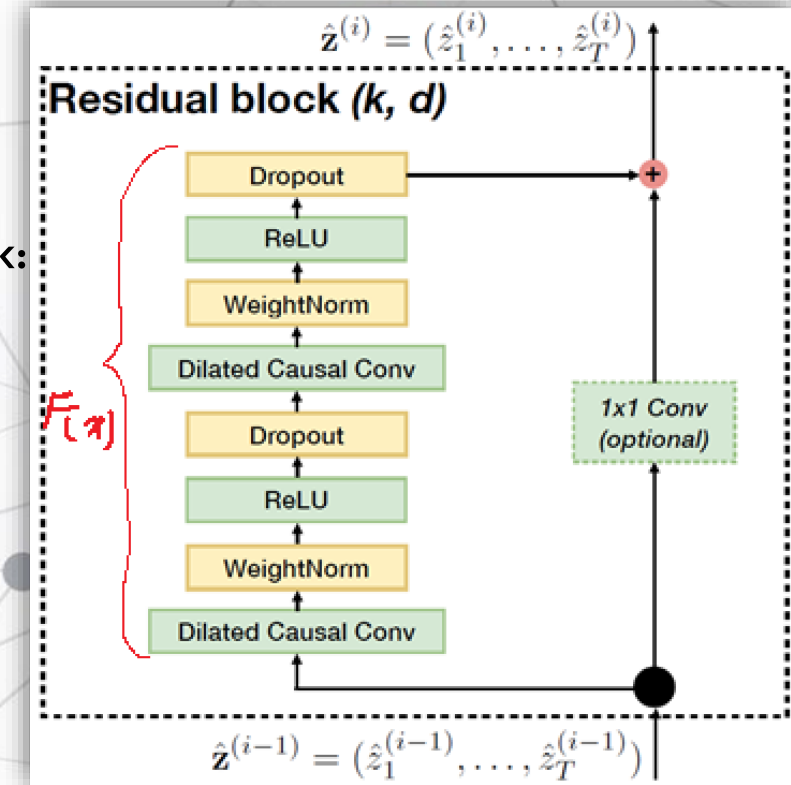
$k$  (filter size) = 3

# TCN – RESIDUAL (BLOCKS) CONNECTIONS

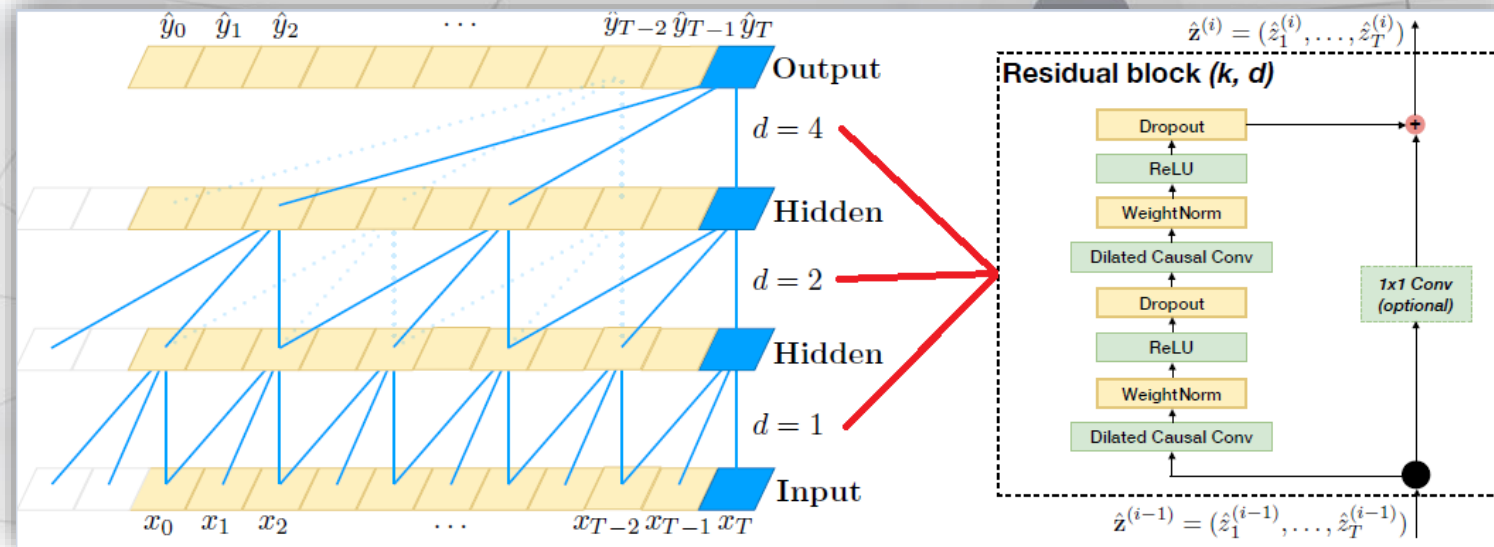
- ❖ A residual block contains a branch leading out to a series of transformations  $F$ , whose outputs are added to the input  $x$  of the block:

$$o = \text{Activation}(x + \mathcal{F}(x))$$

- ❖ This effectively allows layers to learn modifications to the identity mapping  $x$  rather than the entire transformation  $F(x)$ , which has repeatedly been shown to benefit very deep networks.



# TCN – RESIDUAL (BLOCKS) CONNECTIONS (CONT.)



- ❖ **Normalization:** applied weight normalization to the convolutional filters.
- ❖ **Regularization:** a spatial dropout was added after each dilated convolution and at each training step, a whole channel is zeroed out.
- ❖ An **1x1 convolution** is added when residual input and output have different dimensions.



# TCN — SUMMARY

A simple **temporal convolutional network (TCN)** that combines best practices such as **dilations** and **residual connections** with the **causal convolutions** needed for autoregressive prediction.

# TCN HYPERPARAMETERS SETTING

TCN SETTINGS							
Dataset/Task	Subtask	$k$	$n$	Hidden	Dropout	Grad Clip	Note
The Adding Problem	$T = 200$	6	7	27	0.0	N/A	
	$T = 400$	7	7	27			
	$T = 600$	8	8	24			
Seq. MNIST	-	7	8	25	0.0	N/A	
		6	8	20			
Permuted MNIST	-	7	8	25	0.0	N/A	
		6	8	20			
Copy Memory Task	$T = 500$	6	9	10	0.05	1.0	RMSprop 5e-4
	$T = 1000$	8	8	10			
	$T = 2000$	8	9	10			
Music JSB Chorales	-	3	2	150	0.5	0.4	
Music Nottingham	-	6	4	150	0.2	0.4	
Word-level LM	PTB	3	4	600	0.5	0.4	Embed. size 600
	Wiki-103	3	5	1000	0.4		Embed. size 400
	LAMBADA	4	5	500			Embed. size 500
Char-level LM	PTB	3	3	450	0.1	0.15	Embed. size 100
	text8	2	5	520			

- ❖ The most important factor for picking parameters is to make sure that the TCN has a sufficiently large receptive field by choosing  $k$  and  $d$  that can cover the amount of context needed for the task.
- ❖ There are two ways to increase the receptive field of a TCN: choosing larger filter sizes  $k$  and increasing the dilation factor  $d$ , since the effective history of one layer is  $(k-1)*d$ .

# RESULTS

$^h$  means that higher is better.

$^\ell$  means that lower is better.

Sequence Modeling Task	Model Size ( $\approx$ )	Models			
		LSTM	GRU	RNN	TCN
Seq. MNIST (accuracy $^h$ )	70K	87.2	96.2	21.5	<b>99.0</b>
Permuted MNIST (accuracy)	70K	85.7	87.3	25.3	<b>97.2</b>
Adding problem $T=600$ (loss $^\ell$ )	70K	0.164	<b>5.3e-5</b>	0.177	<b>5.8e-5</b>
Copy memory $T=1000$ (loss)	16K	0.0204	0.0197	0.0202	<b>3.5e-5</b>
Music JSB Chorales (loss)	300K	8.45	8.43	8.91	<b>8.10</b>
Music Nottingham (loss)	1M	3.29	3.46	4.05	<b>3.07</b>
Word-level PTB (perplexity $^\ell$ )	13M	<b>78.93</b>	92.48	114.50	88.68
Word-level Wiki-103 (perplexity)	-	48.4	-	-	<b>45.19</b>
Word-level LAMBADA (perplexity)	-	4186	-	14725	<b>1279</b>
Char-level PTB (bpc $^\ell$ )	3M	1.36	1.37	1.48	<b>1.31</b>
Char-level text8 (bpc)	5M	1.50	1.53	1.69	<b>1.45</b>



**NOTE:** the **number of hidden units** was chosen so that the model size is approximately at the same level as the recurrent models with which we are comparing.

TCN SETTINGS							
Dataset/Task	Subtask	$k$	$n$	Hidden	Dropout	Grad Clip	Note
The Adding Problem	$T = 200$	6	7	27	0.0	N/A	
	$T = 400$	7	7	27			
	$T = 600$	8	8	24			
Seq. MNIST	-	7	8	25	0.0	N/A	
		6	8	20			
Permuted MNIST	-	7	8	25	0.0	N/A	
		6	8	20			
Copy Memory Task	$T = 500$	6	9	10	0.05	1.0	RMSprop 5e-4
	$T = 1000$	8	8	10			
	$T = 2000$	8	9	10			
Music JSB Chorales	-	3	2	150	0.5	0.4	
Music Nottingham	-	6	4	150	0.2	0.4	
Word-level LM	PTB	3	4	600	0.5	0.4	Embed. size 600
	Wiki-103	3	5	1000	0.4		Embed. size 400
	LAMBADA	4	5	500			Embed. size 500
Char-level LM	PTB	3	3	450	0.1	0.15	Embed. size 100
	text8	2	5	520			

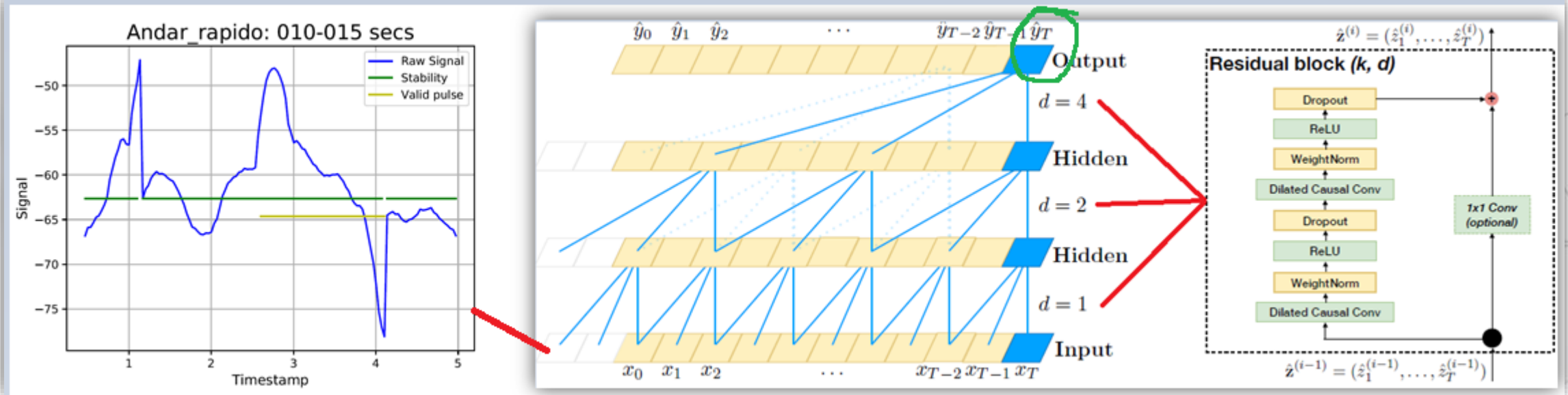
Sequence Modeling Task	Model Size ( $\approx$ )	Models			
		LSTM	GRU	RNN	TCN
Seq. MNIST (accuracy <sup><i>h</i></sup> )	70K	87.2	96.2	21.5	<b>99.0</b>
Permuted MNIST (accuracy)	70K	85.7	87.3	25.3	<b>97.2</b>
Adding problem $T=600$ (loss <sup><math>\ell</math></sup> )	70K	0.164	<b>5.3e-5</b>	0.177	<b>5.8e-5</b>
Copy memory $T=1000$ (loss)	16K	0.0204	0.0197	0.0202	<b>3.5e-5</b>
Music JSB Chorales (loss)	300K	8.45	8.43	8.91	<b>8.10</b>
Music Nottingham (loss)	1M	3.29	3.46	4.05	<b>3.07</b>
Word-level PTB (perplexity <sup><math>\ell</math></sup> )	13M	<b>78.93</b>	92.48	114.50	88.68
Word-level Wiki-103 (perplexity)	-	48.4	-	-	<b>45.19</b>
Word-level LAMBADA (perplexity)	-	4186	-	14725	<b>1279</b>
Char-level PTB (bpc <sup><math>\ell</math></sup> )	3M	1.36	1.37	1.48	<b>1.31</b>
Char-level text8 (bpc)	5M	1.50	1.53	1.69	<b>1.45</b>

- 1) “TCNs can be build to have very long effective history sizes, which means they have the ability to look very far into the past to make a prediction. To this end, a combination of very deep networks augmented with residual layers and dilated convolutions are deployed.”
- 2) “Generic TCN architecture outperform canonical recurrent architectures across a broad variety of sequence modelling tasks that are commonly used to benchmark the performance of recurrent architectures.”



Ruslan Padnevych n°47222  
r.padnevych@campus.fct.unl.pt

# MY THESIS





# Temporal Convolution Network - Liveness Detection

## Hyperparameters

Kernel size  $k = 3$

Dilation factor  $d = 2^v$ , with  $v$  the level of the network

Padding  $p = (k-1) * d$

