

# Convolutions in the world of sequences

Alexey Zaytsev

Skoltech

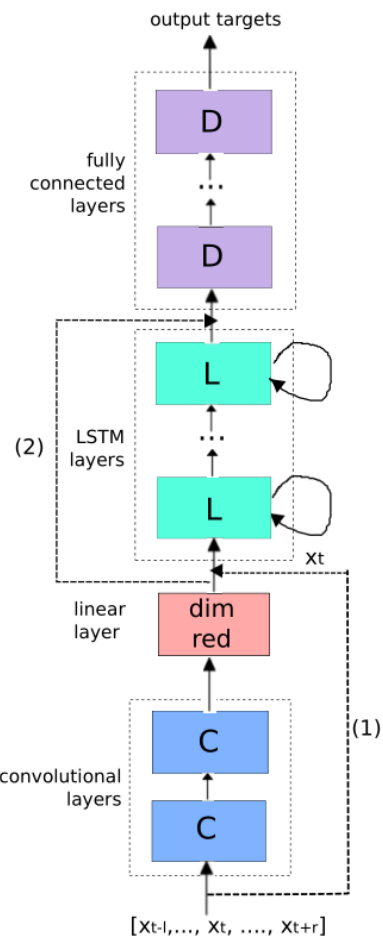


# Main types of models for sequential data

1. Recurrent Neural Networks
2. **One-dimensional  
convolutional neural  
networks**
3. Transformers

# Sequential data models for acoustic data processing

Model	WER ↓
FCNN	18.4
CNN	18.0
LSTM	18.0
CNN+LSTM	17.6
LSTM+FCNN	17.6
CNN+LSTM+FCNN	<b>17.3</b>



Sainath, Tara N., et al. "Convolutional, long short-term memory, fully connected deep neural networks." *IEEE ICASSP*. 2015.

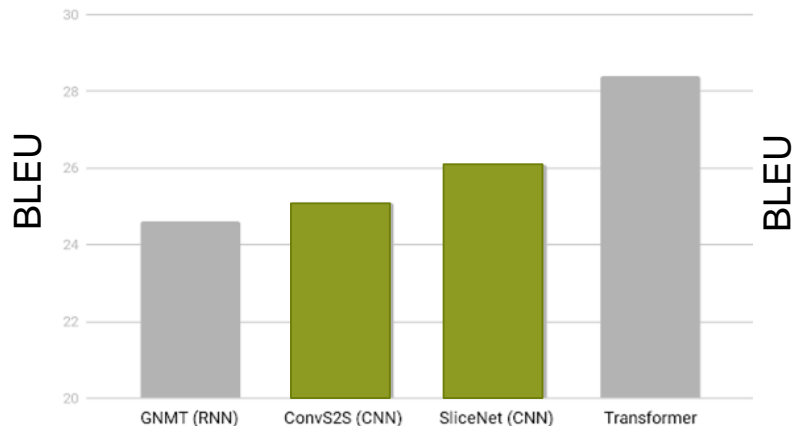
CNN+LSTM+FCNN architecture

# Motivation for CNNs

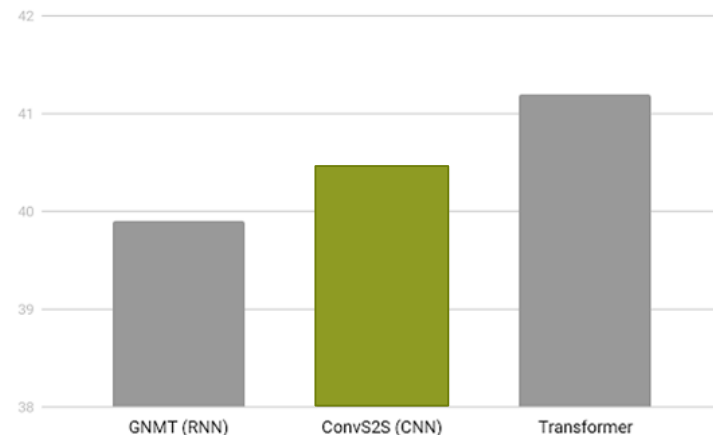
# CNNs are better than RNNs for Machine Translation

Skoltech

English German Translation quality



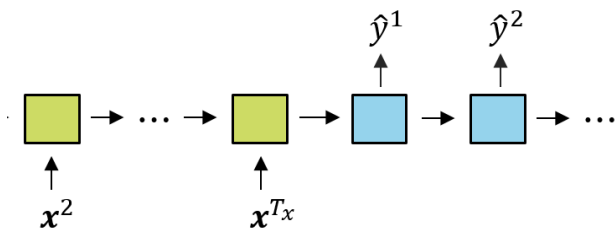
English French Translation Quality



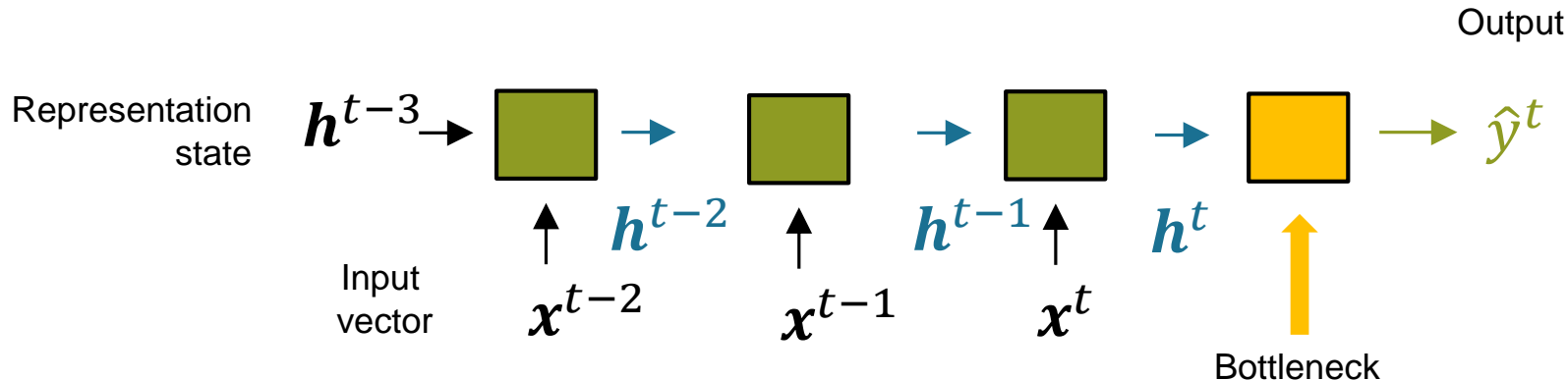
↑  
We maximize BLEU

# What is wrong with RNNs?

- A. The bottleneck for the last token
- B. Goes from past to future. *Bidirectional RNNs can help there*
- C. Long time to train with GPU



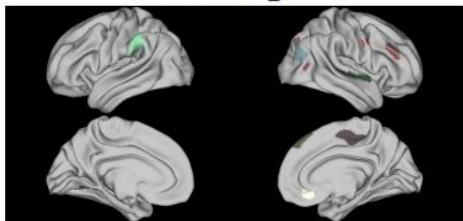
The bottleneck is more evident for seq2seq problems



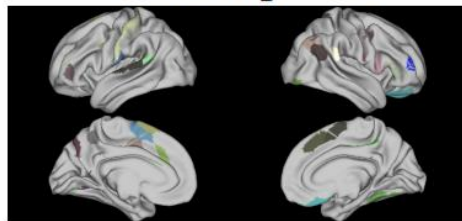
# What is wrong with RNNs?

A. Most important part is the end

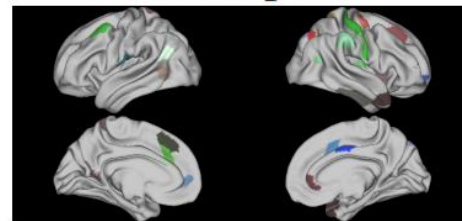
Time step = 0



Time step = 20



Time step = 40



# CNN idea

1. Compute embeddings for all subsequences of a certain length
2. Combine them in a meaningful way

Example: “Now I need a place to hide away” has vectors related to:  
“now I need”, “need a place”, “a place to”, “place to hide”, “to hide away”

Issues:

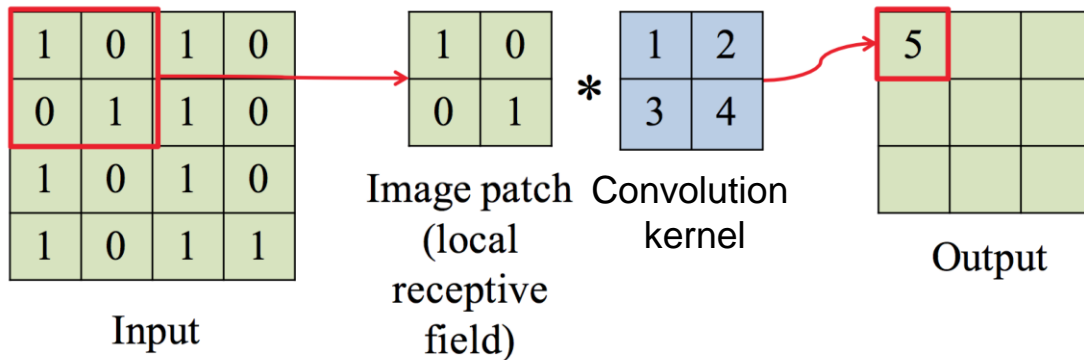
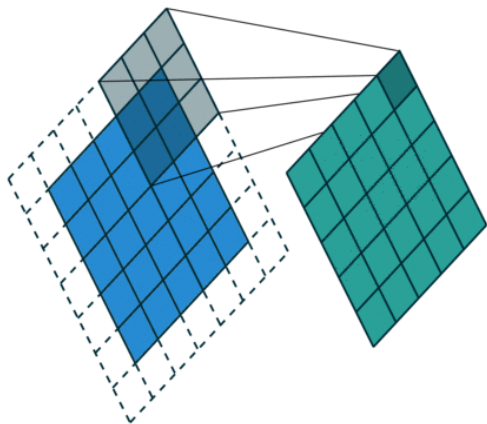
- No grammatical sense
- Approach not backed by linguistics



# 1D CNN architecture

# “Classic” 2D convolution for images

Input      Output



# Definition of 1D convolution

$$\tilde{x}_t = \sum_{i=-h}^h x_{t+i} w_i$$

convolution of size  $(2h + 1)$ ,  $t$  is from 1 to  $T$   
 $w$  is the convolution kernel,  $x$  is input data,  $\tilde{x}$  is input data

0.1
0.6
-0.2
0.9
0.4
0.7
0



-1
1
0

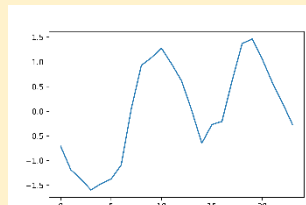


0.5
-0.4
-0.7
-1.3
-1.1

Apply a *filter* of size 3

# 1D CNN at home: ROCKET

Input: 1D  
time series



→ Convolved time series  
for 10K random kernels  
 $O(knTh)$

→ Max and percentage of  
positive values  
 $O(knTh)$

↓  $O(n(2k)^2)$  or  
 $O(n^2 2k)$

Ridge regression (or Logistic  
regression) as a classifier

$k$  – number of kernels  
 $n$  – number of examples  
 $T$  – length of the time series  
 $h$  – kernel width

# ROCKET and Mini ROCKET hyperparameters

	ROCKET	MINIROCKET
length	{7, 9, 11}	9
weights	$\mathcal{N}(0, 1)$	$\{-1, 2\}$
bias	$\mathcal{U}(-1, 1)$	from convolution output
dilation	random	fixed (rel. to input length)
padding	random	fixed
features	PPV + max	PPV
num. features	20K	10K

# 1D convolution for tokens

$$\tilde{x}_t = \sum_{i=-h}^h \sum_{j=1}^k x_{t-i,j} w_{i,j}$$

convolution of size  $(2h + 1) \times k$ ,  $j$  is from 1 to  $T$   
 $\mathbf{w}$  is the convolution kernel,  $\mathbf{x}$  is input data,  $\tilde{\mathbf{x}}$  is input data

now	0.1	-0.2	0.4
I	0.6	0.2	0.1
need	-0.2	0.1	0.0
a	0.9	0.1	0.2
place	0.4	0.2	0.6
to	0.7	-1.2	1
hide	0	0.5	0.5



3	2	-1
1	1	1
2	-1	0



n,l,n	-0.1
i,n,a	3.3
n,a,p	1.4
a,p,t	3.7
p,t,h	1.1

Apply a *filter* of size 3

# 1D convolution for tokens with padding

pad	0	0	0
now	0.1	-0.2	0.4
I	0.6	0.2	0.1
need	-0.2	0.1	0.0
a	0.9	0.1	0.2
place	0.4	0.2	0.6
to	0.7	-1.2	1
hide	0	0.5	0.5
pad	0	0	0

3	2	-1
1	1	1
2	-1	0

P,n,I	1.3
n,I,n	-0.1
i,n,a	3.3
n,a,p	1.4
a,p,t	3.7
p,t,h	1.1
t,h,P	-0.3

Apply a *filter* or *kernel* of size 3

# 1D convolution for tokens with max pooling

now	0.1	-0.2	0.4
I	0.6	0.2	0.1
need	-0.2	0.1	0.0
a	0.9	0.1	0.2
place	0.4	0.2	0.6
to	0.7	-1.2	1
hide	0	0.5	0.5

3	2	-1
1	1	1
2	-1	0

Apply a *filter* or *kernel* of size 3

n,l,n	-0.1
i,n,a	3.3
n,a,p	1.4
a,p,t	3.7
p,t,h	1.1

Max pool	3.7
-------------	-----



# 1D convolution for tokens with mean pooling

now	0.1	-0.2	0.4
I	0.6	0.2	0.1
need	-0.2	0.1	0.0
a	0.9	0.1	0.2
place	0.4	0.2	0.6
to	0.7	-1.2	1
hide	0	0.5	0.5

3	2	-1
1	1	1
2	-1	0

Apply a *filter* or *kernel* of size 3

n,I,n	-0.1
i,n,a	3.3
n,a,p	1.4
a,p,t	3.7
p,t,h	1.1

Mean pool	1.88
-----------	------

# 1D convolution for tokens with many kernels and mean pooling

now	0.1	-0.2	0.4
I	0.6	0.2	0.1
need	-0.2	0.1	0.0
a	0.9	0.1	0.2
place	0.4	0.2	0.6
to	0.7	-1.2	1
hide	0	0.5	0.5



3	2	-1
1	1	1
2	-1	0

2	2	-1
1	1	1
2	-1	0



Mean pool	1.88	1.52
-----------	------	------



n,l,n	-0.1	-0.2
i,n,a	3.3	2.7
n,a,p	1.4	1.6
a,p,t	3.7	2.8
p,t,h	1.1	0.7

Apply many *filters* or *kernel* of size 3

# PyTorch 1D convolution summary

```
batch_size = 16
word_embed_size = 4
seq_len = 7
input = torch.randn(batch_size, word_embed_size, seq_len)
conv1 = Conv1d(in_channels=word_embed_size, out_channels=3,
               kernel_size=3) # can add: padding=1
hidden1 = conv1(input)
hidden2 = torch.max(hidden1, dim=2) # max pool
```

# Other things to do: stride, pooling, dilation, multiple layers

- Dilation: increase receptive field
- Multiple convolutional layers: more layers
- Stride: stride=2 – skip odd triplets
- Max pooling of stride
- Top k-max pooling

now	0.1	-0.2	0.4
I	0.6	0.2	0.1
need	-0.2	0.1	0.0
a	0.9	0.1	0.2
place	0.4	0.2	0.6
to	0.7	-1.2	1
hide	0	0.5	0.5

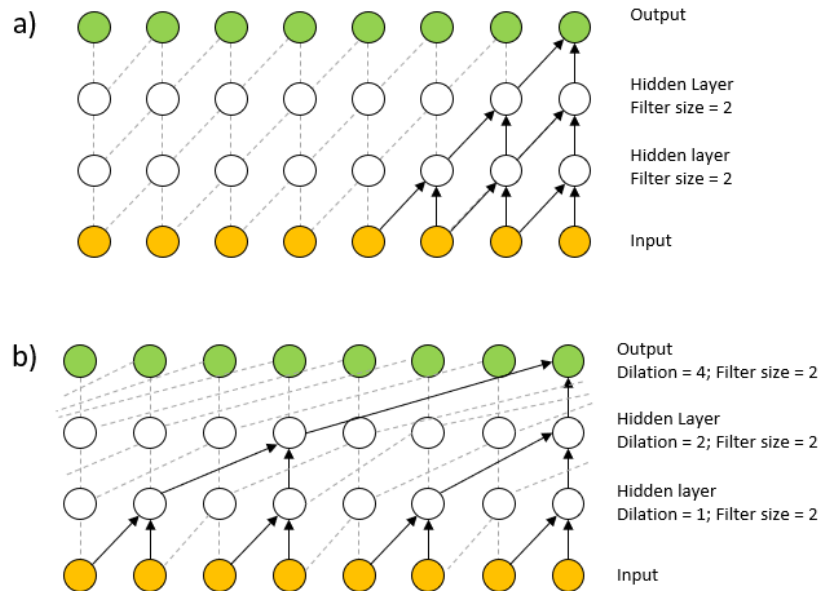


Figure source:  
Benson, B., et al. "Forecasting solar cycle 25  
using deep neural networks." *Solar Physics* 295  
(2020): 1-15.

# Single Layer CNN for sentence classification

- Simple: one convolutional layer + pooling
- Word vectors  $\mathbf{x}_i \in R^k$
- Sentence  $\mathbf{x}_{0:n} = [\mathbf{x}_0, \mathbf{x}_2, \dots, \mathbf{x}_{n-1}]$
- Concatenation of words  $\mathbf{x}_{i:i+h}$
- Convolutional filter  $\mathbf{w} \in R^{hk}$
- Filters can be of different size

# Single Feature Inference

- Convolutional filter  $\mathbf{w} \in R^{hk}$  applied to all possible windows
- Feature computation

$$c_i = f(\mathbf{w}^T \mathbf{x}_{i:i+h} + b)$$

- Result: a feature map

$$\mathbf{c} = [c_0, c_2, \dots, c_{n-h}] \in R^{n-h+1}$$

# Pooling and channels

- Idea: capture the most important information: maximum over time
- Feature map  $\mathbf{c} = [c_0, c_2, \dots, c_{n-h}]$
- Pooled single number  $\tilde{c} = \max(\mathbf{c})$
- Use multiple filters with different weights, different windows sizes
- For max pooling length is irrelevant

# Full network architecture

- One convolution followed by max-pooling
- Feature map  $\mathbf{z} = [\tilde{c}_1, \dots, \tilde{c}_m]$
- Final softmax layer
$$y = \text{softmax}(W\mathbf{z} + \mathbf{b})$$



# Multiple-channel input idea

- Initialize with pre-trained word vectors
- Start with two similar instances
- Change one instance, freeze other
- Both channels are added before max-pooling

# Regularization

Dropout!

- During *training*: Bernoulli random variable to drop with probability  $p$  features
- During *inference*: Multiply parameters matrix  $W$  by  $p$

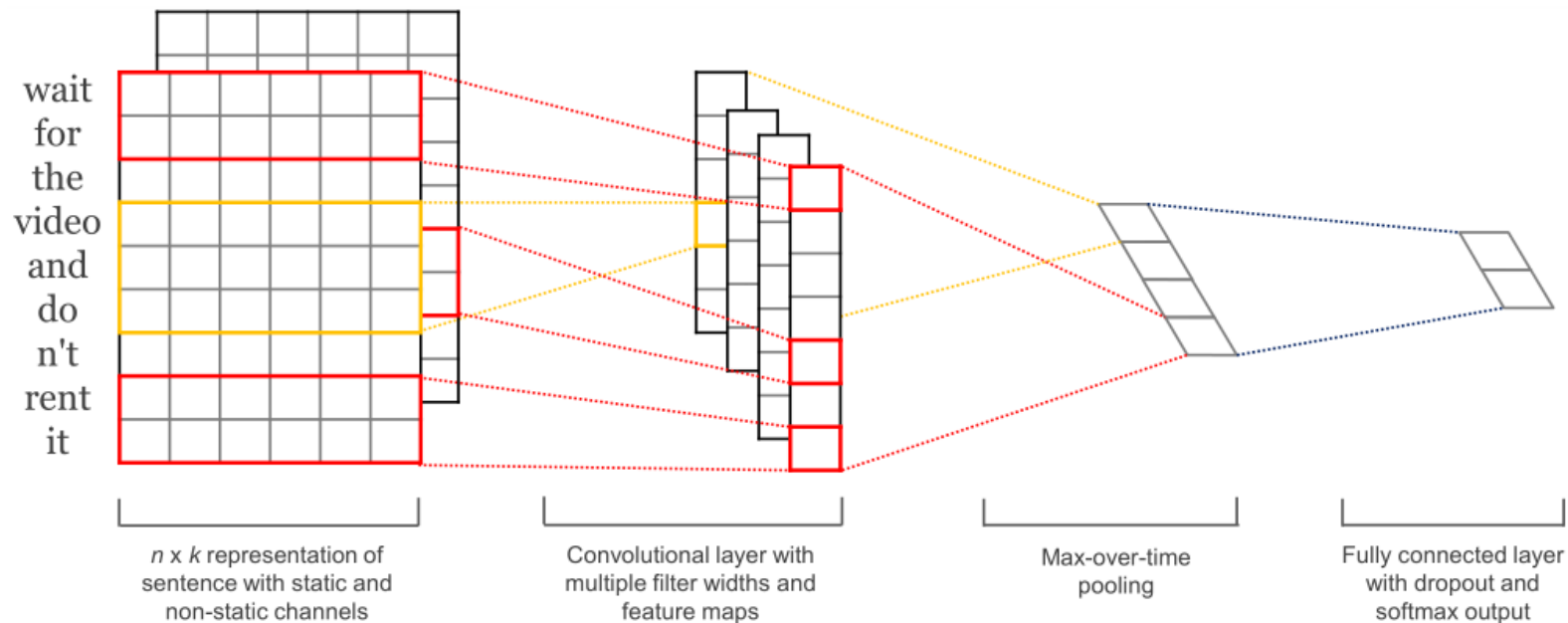
Constrain L2 norms for weights vector for each class:

- Row norm for each row is smaller than  $s$
- If bigger – rescale to meet the constraint
- Not very common (may be unnecessary)

# Hyperparameters

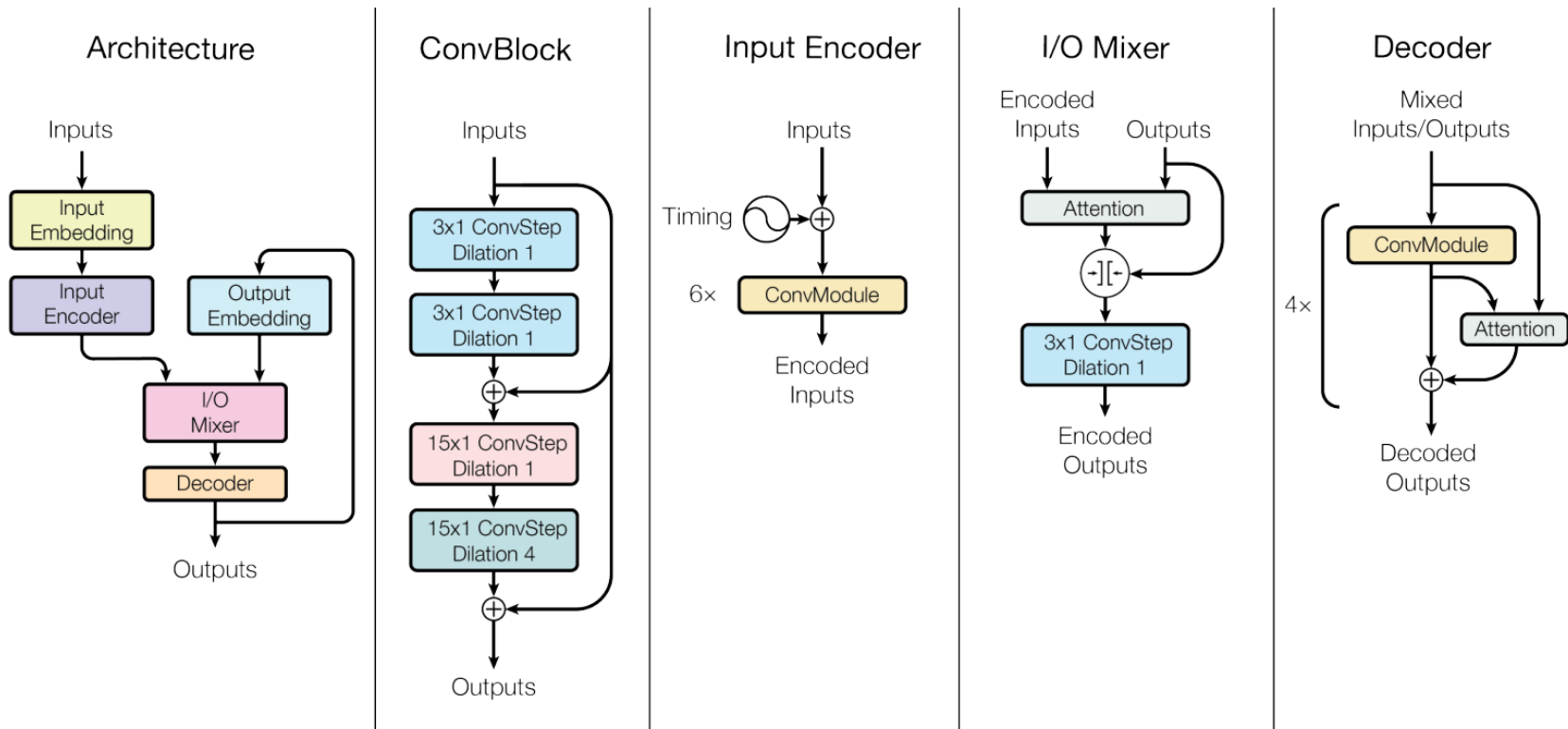
- ReLU activation
- Window filter sizes  $h = 3, 4, 5$
- Each filter size: 100 feature maps
- Dropout 2-4% accuracy improvement via dropout
- L2 constraint  $s$  for rows of softmax,  $s=3$
- Mini-batch size for SGD training: 5
- Word2vec with  $k=300$  features
- Use validation set error as a stopping criterion

# CNN overall architecture from Yoon Kim [2014]



Kim, Yoon. "Convolutional neural networks for sentence classification." *arXiv preprint arXiv:1408.5882* (2014).

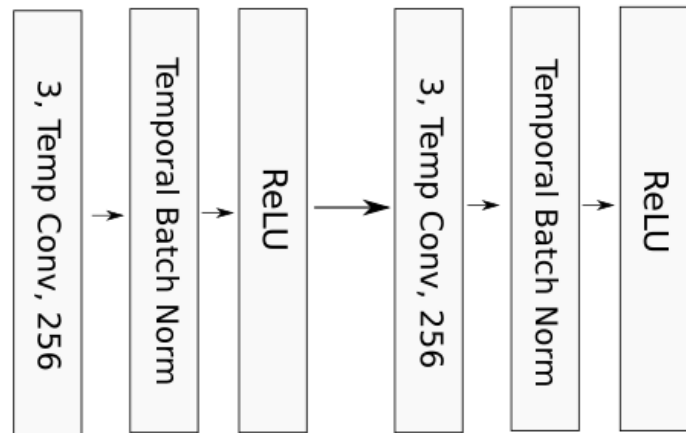
# CNNs with fewer parameters



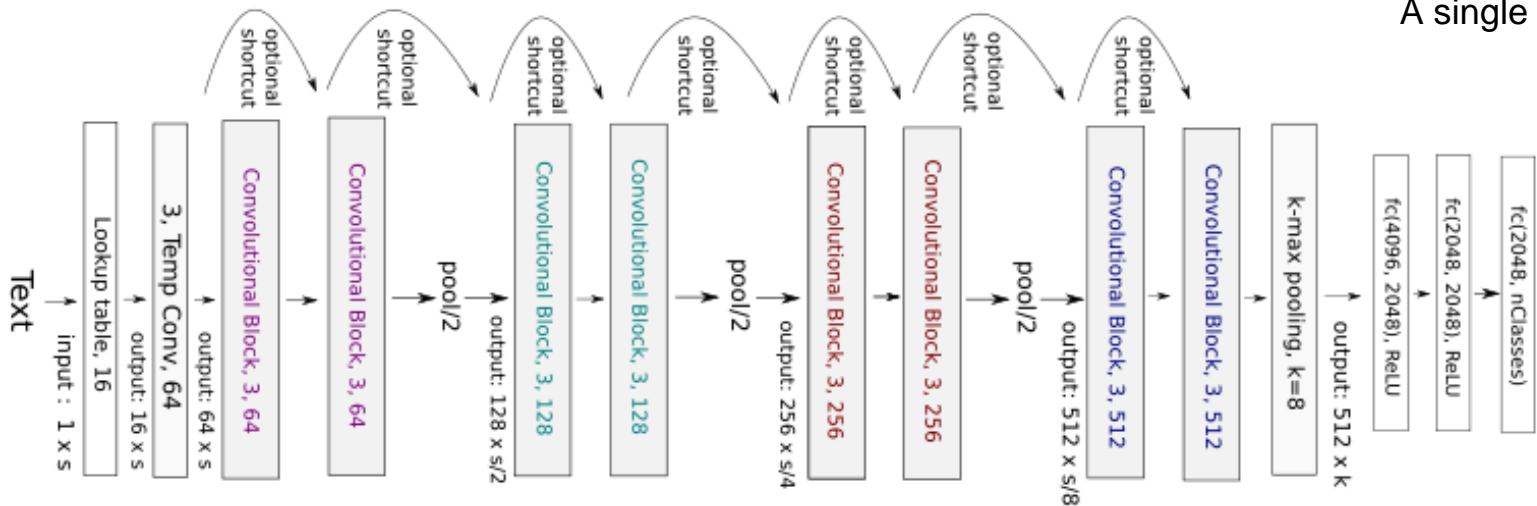
Kaiser, Lukasz, Aidan N. Gomez, and Francois Chollet. "Depthwise separable convolutions for neural machine translation." ICLR. 2018.

# Very deep CNNs for text classification

Brings ideas from ResNet + character level



A single convolutional block



Full network

Conneau, A., Schwenk, H., Barrault, L., & Le Cun, Y. Very deep convolutional networks for text classification. EACL. 2017

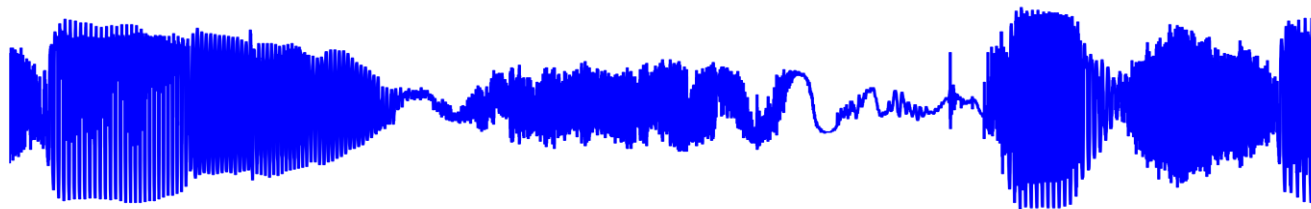
# WaveNet: a generative model for audio

van den Oord, Aäron, et al. "WaveNet: A Generative Model for Raw Audio." *9th ISCA Speech Synthesis Workshop*. 2016

# WaveNet key features

WaveNet, an audio generative model based on the PixelCNN architecture:

- SOTA in 2016 for the text-to-speech problem
- New architecture with dilated convolutions
- The main ingredient is causal convolutions



A second of generated speech

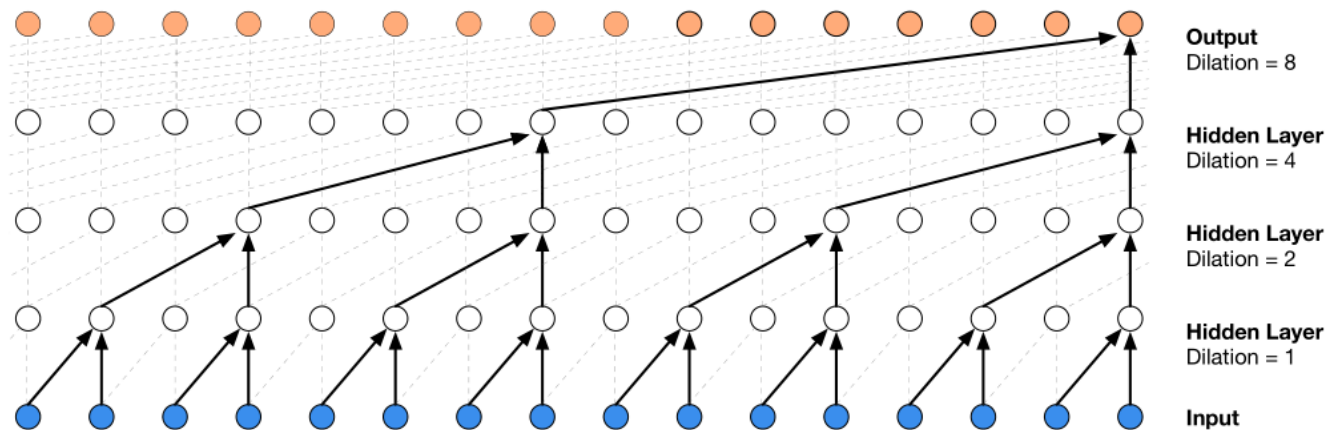
van den Oord, Aäron, et al. "WaveNet: A Generative Model for Raw Audio." *9th ISCA Speech Synthesis Workshop*. 2016



# WaveNet model

Causal convolutions looks only in the past and model

$$p(x^t | x^1, x^2, \dots, x^{t-1})$$



Dilated convolutions allow to make the receptive field larger

# WaveNet convolutions

$$f(x_t) = \text{sign}(x_t) \frac{\ln(1 + \mu |x_t|)}{\ln(1 + \mu)}$$

Convolutions with sigmoid gate

$$\mathbf{z} = \tanh(W_{f,k} * \mathbf{x}) \odot \sigma(W_{g,k} * \mathbf{x})$$

Convolutions with global context  $\mathbf{h}$

$$\mathbf{z} = \tanh(W_{f,k} * \mathbf{x} + V_{f,k}^T \mathbf{h}) \odot \sigma(W_{g,k} * \mathbf{x} + V_{g,k}^T \mathbf{h})$$

Convolutions with local context  $\mathbf{y}$

$$\mathbf{z} = \tanh(W_{f,k} * \mathbf{x} + V_{f,k} * \mathbf{y}) \odot \sigma(W_{g,k} * \mathbf{x} + V_{g,k} * \mathbf{y})$$

Prediction: softmax over quantized outputs

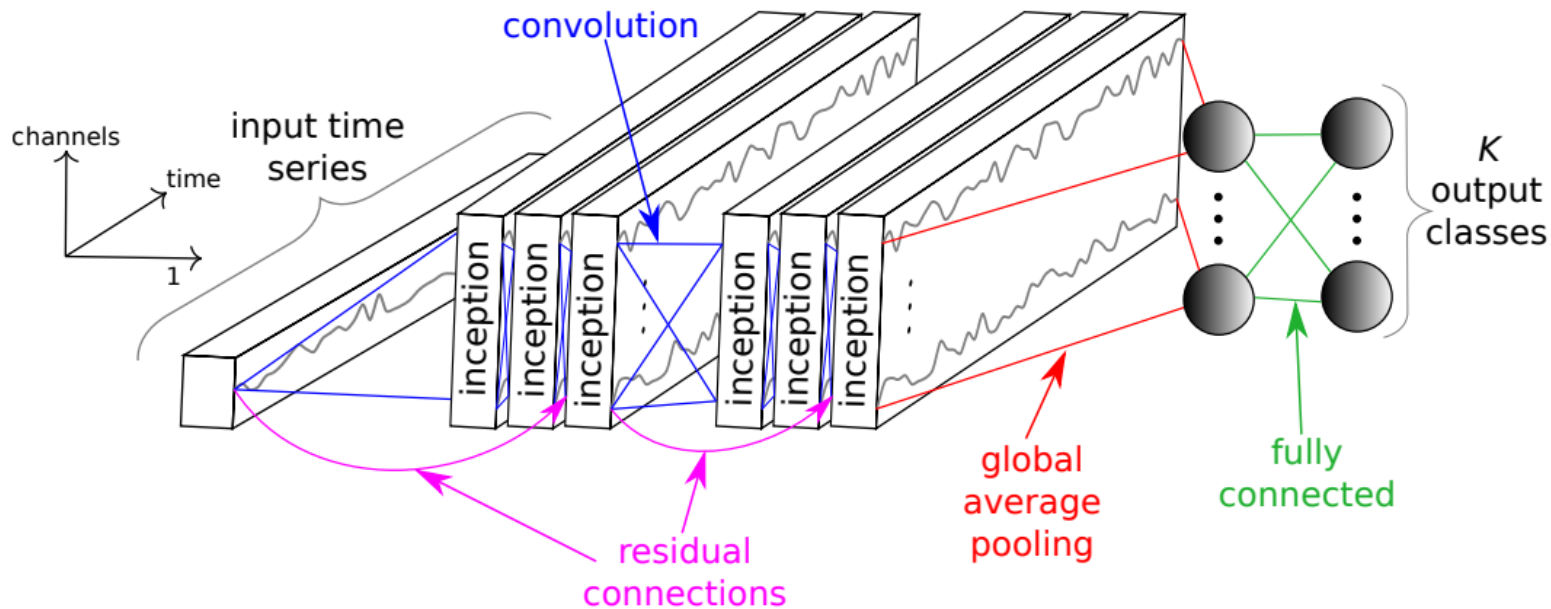
# InceptionTime: 1D CNN for time series classification

Ismail Fawaz, Hassan, et al. "InceptionTime: Finding AlexNet for time series classification." *Data Mining and Knowledge Discovery*. 2020.

# InceptionTime architecture

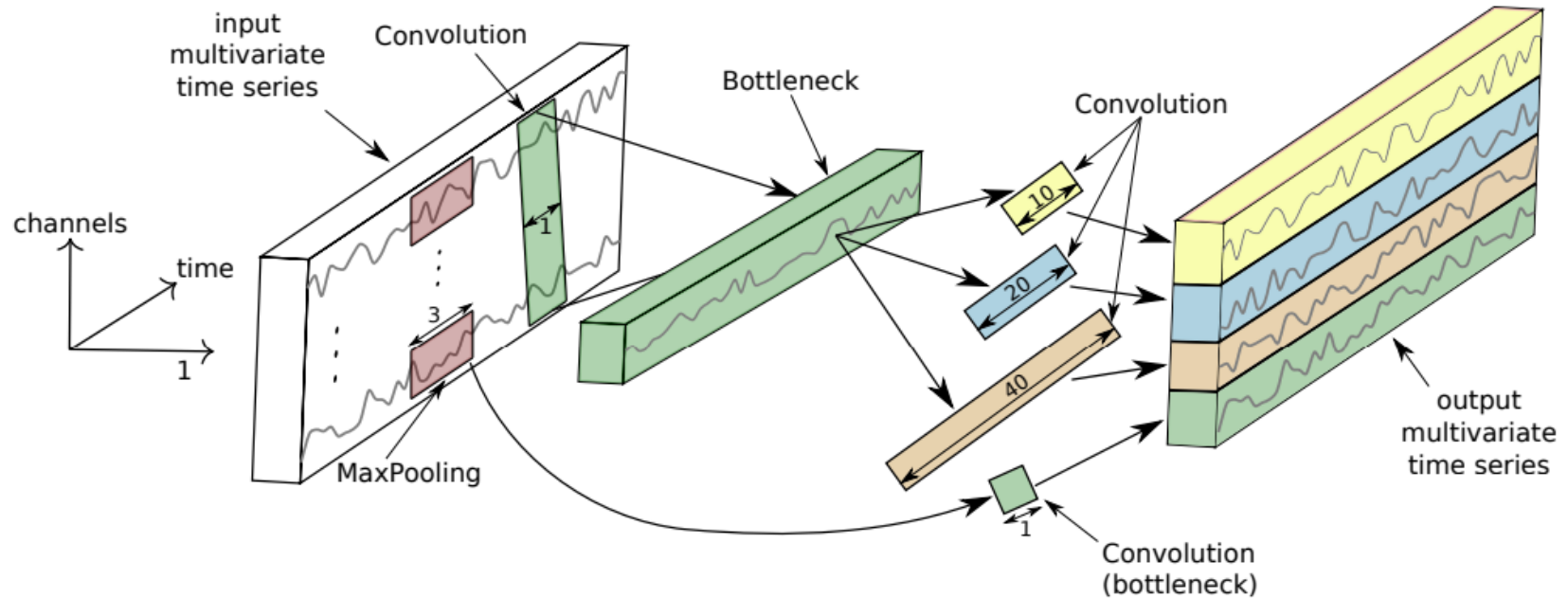
**Goal:** time series classification, one label for a series,

**Model:** ensemble of 5 independent neural networks



# InceptionTime block

...similar to ResNet

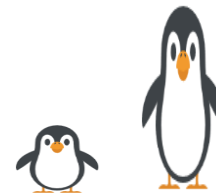


Ismail Fawaz, Hassan, et al. "InceptionTime: Finding AlexNet for time series classification." *Data Mining and Knowledge Discovery*. 2020.

# Conclusions

# Conclusions

- CNN is a powerful tool for sequential data processing
- With pooling we can overcome varying length of a sequence
- Simple CNN is a good point to start
- Many tricks from Computer Vision CNNs work
- Deeper CNNs for sequential data processing can be better



# References

1. van den Oord, Aäron, et al. "WaveNet: A Generative Model for Raw Audio." *9th ISCA Speech Synthesis Workshop*. 2016.
2. Ismail Fawaz, Hassan, et al. "InceptionTime: Finding AlexNet for time series classification." *Data Mining and Knowledge Discovery*. 2020.
3. Conneau, A., Schwenk, H., Barrault, L., & Le Cun, Y. Very deep convolutional networks for text classification. *EACL*. 2017.
4. Ismail, Aya Abdelsalam, et al. "Input-cell attention reduces vanishing saliency of recurrent neural networks." *NeurIPS*. 2019.