# Lecture 9:
# RNN, Seq2Seq & Attention Part2
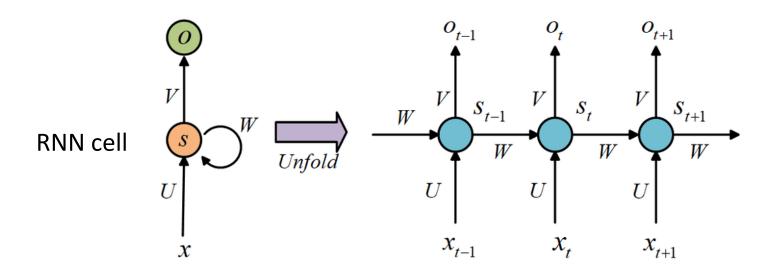
Olexandr Isayev

Department of Chemistry, CMU

olexandr@cmu.edu

Home Assignment #3
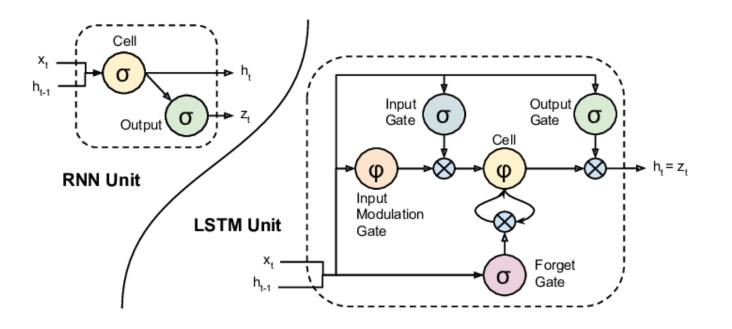

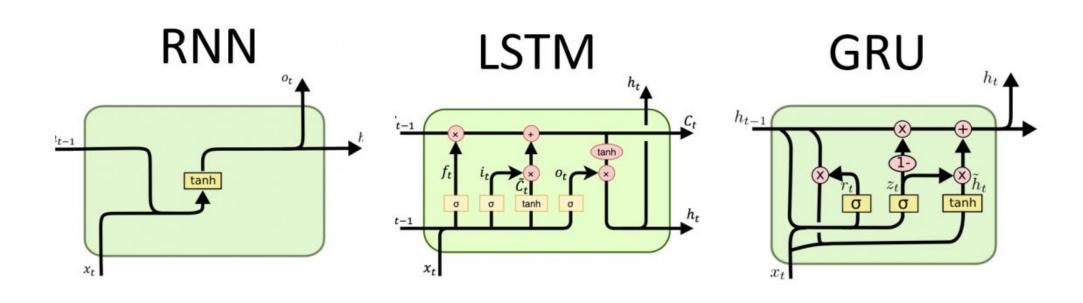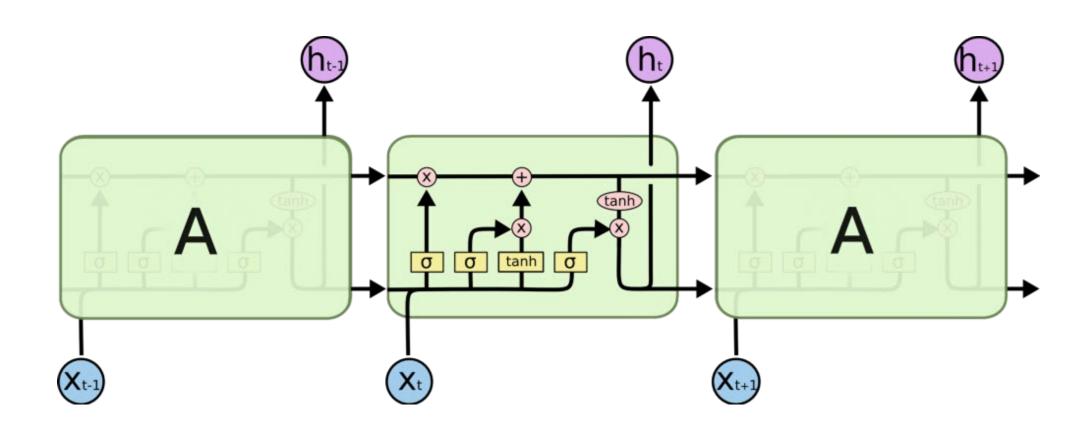In Class Project

# Main Concept of RNNs

RNN cell



Unfold

Time

# RNN Cells



RNN Unit

LSTM Unit

RNN

LSTM

GRU

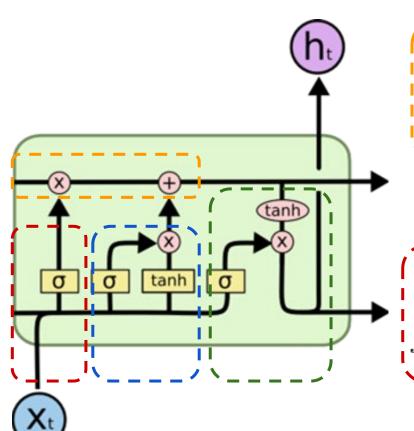# LSTM: Long short term memory

# LSTM big picture ...



**Cell State**

$$C_t = f_t * C_{t-1} + i_t * \tilde{C}_t$$

**Output Gate**

$$o_t = \sigma(W_o \cdot [h_{t-1}, x_t] + b_o)$$

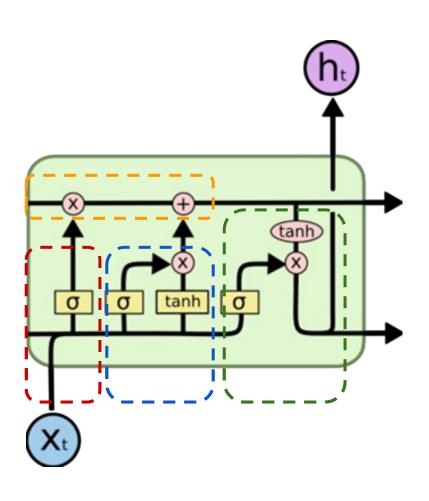$$h_t = o_t * \tanh(C_t)$$

**Forget Gate**

$$f_t = \sigma(W_f \cdot [h_{t-1}, x_{t-1}] + b_f)$$

**Input Gate**

$$i_t = \sigma(W_i \cdot [h_{t-1}, x_t] + b_i)$$

$$\tilde{C}_t = \tanh(W_c \cdot [h_{t-1}, x_t] + b_f)$$

# LSTM big picture ...



1. LSTM are recurrent neural network with a cell and a hidden state, boths of these are updated in each step and can be thought as memories.

2. Cell states work as a long term memory and the updates depends on the relation between the hidden state in t -1 and the input.

3. The hidden state of the next step is a transformation of the cell state and the output (which is the section that is in general used to calculate our loss, ie information that we want in a short memory).

# What makes Recurrent Networks so special?

# RNN Structures

$Y_t$

$X_t$

one to one

- The **one to one** structure is useless.
- It takes a single input and it produces a single output.
- Not useful because the RNN cell is making little use of its unique ability to remember things about its input sequence

# RNN Structures (cont)



The **many to one** structure reads in a sequence and gives us back a single value.

Example: Sentiment analysis, where the network is given a piece of text and then reports on some quality inherent in the writing. A common example is to look at a movie review and determine if it was positive or negative.

# RNN Structures (cont)



The **one to many** takes in a single piece of data and produces a sequence.
For example we give it the starting note for a song, and the network produces the rest of the melody for us.

Generative models

# RNN Structures (cont)



$Y_{t-2}$  $Y_{t-1}$  $Y_t$

$X_{t-2}$  $X_{t-1}$  $X_t$

many to many

The **many to many** structures are in some ways the most interesting. used for machine translation.

Seq2seq models

Language translation
…

# Bidirectional

- LSTM and RNN are designed to analyze sequence of values.

- For example: *Patrick said he needs a vacation.*
- *he* here means *Patrick* and we know this because *Patrick* was before the word *he.*

- However consider the following sentence:
- *He needs to work more, Peter said about Patrick.*

- Bidirectional RNN or BRNN or bidirectional LSTM or BLSTM when using LSTM units.

# Bidirectional (cond)

symbol for a BRNN

$Y_t$

$X_t$

$Y_{t-2}$   $Y_{t-1}$   $Y_t$

previous state

previous state

$X_{t-2}$   $X_{t-1}$   $X_t$

# Deep RNN

- LSTM units can be arranged in layers, so that each the output of each unit is the input to the other units. This is called **a deep RNN**, where the adjective "deep" refers to these multiple layers.

- Each layer feeds the LSTM on the next layer

- First time step of a feature is fed to the first LSTM, which processes that data and produces an output (and a new state for itself).

- That output is fed to the next LSTM, which does the same thing, and the next, and so on.

- Then the second time step arrives at the first LSTM, and the process repeats.

# Deep RNN

$Y$

$X$

$Y_{t-2}$   $Y_{t-1}$   $Y_t$   $Y_{t+1}$   $Y_{t+2}$

$X_{t-2}$   $X_{t-1}$   $X_t$   $X_{t+1}$   $X_{t+2}$

# Skip Connections

- Add additional connections between units $d$ time steps apart

- Creating paths through time where gradients neither vanish or explode

# Multi-layer RNNs

- We can of course design RNNs with multiple hidden layers

$y_1$  $y_2$  $y_3$  $y_4$  $y_5$  $y_6$

$x_1$  $x_2$  $x_3$  $x_4$  $x_5$  $x_6$

- Anything goes: skip connections across layers, across time, …

# Multi-layer RNNs

- We can of course design RNNs with multiple hidden layers



- Anything goes: skip connections across layers, across time, …

# Language modeling: Character RNN

Output symbol $y_i$

Output layer (linear + softmax)

Hidden state $h_i$

One-hot encoding $x_i$

Input symbol



$$p(y_1, y_2, \ldots, y_n)$$

$$= \prod_{i=1}^{n} p(y_i | y_1, \ldots, y_{i-1})$$

$$\approx \prod_{i=1}^{n} P_W(y_i | h_i)$$

# Language modeling: Character RNN

100th iteration

```
tyntd-iafhatawiaoihrdemot  lytdws  e ,tfti, astai f ogoh eoase rrranbyne 'nhthnee e
plia tklrgd t o idoe ns,smtt   h ne etie h,hregtrs nigtike,aoaenns lng
```

train more

300th iteration

```
"Tmont thithey" fomesscerliund
Keushey. Thom here
sheulke, anmerenith ol sivh I lalterthend Bleipile shuwy fil on aseterlome
coaniogennc Phe lism thond hon at. MeiDimorotion in ther thize."
```

train more

700th iteration

```
Aftair fall unsuch that the hall for Prince Velzonski's that me of
her hearly, and behs to so arwage fiving were to it beloge, pavu say falling misfort
how, and Gogition is so overelical and ofter.
```
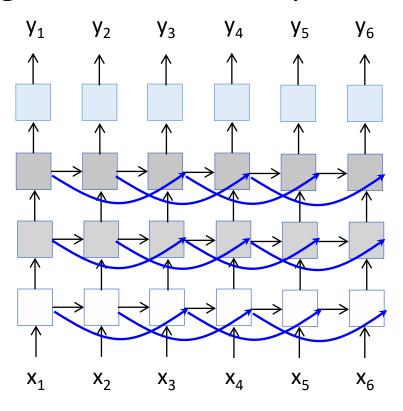
train more

2000th

iteration

```
"Why do what that day," replied Natasha, and wishing to himself the fact the
princess, Princess Mary was easier, fed in had oftened him.
Pierre aking his soul came to the packs and drove up his father-in-law women.
```

http://karpathy.github.io/2015/05/21/rnn-effectiveness/

# Searching for interpretable hidden units



quote detection cell

A. Karpathy, J. Johnson, and L. Fei-Fei, Visualizing and Understanding Recurrent Networks, ICLR Workshop 2016

# Searching for interpretable hidden units



line position tracking cell

A. Karpathy, J. Johnson, and L. Fei-Fei, Visualizing and Understanding Recurrent Networks, ICLR Workshop 2016

# Searching for interpretable hidden units



if statement cell

A. Karpathy, J. Johnson, and L. Fei-Fei, Visualizing and Understanding Recurrent Networks, ICLR Workshop 2016

# Searching for interpretable hidden units



quote/comment cell

A. Karpathy, J. Johnson, and L. Fei-Fei, Visualizing and Understanding Recurrent Networks, ICLR Workshop 2016

# Searching for interpretable hidden units



code depth cell

A. Karpathy, J. Johnson, and L. Fei-Fei, Visualizing and Understanding Recurrent Networks, ICLR Workshop 2016

# Searching for interpretable hidden units



```
/* Unpack a filter field's string representation from user-space
 * buffer. */
char *audit_unpack_string(void **bufp, size_t *remain, size_t len)
{
	char *str;
	if (!*bufp || (len == 0) || (len > *remain))
		return ERR_PTR(-EINVAL);
	/* Of the currently implemented string fields, PATH_MAX
	 * defines the longest valid length.
	 */
```

¯\\_(ツ)_/¯

A. Karpathy, J. Johnson, and L. Fei-Fei, Visualizing and Understanding Recurrent Networks, ICLR Workshop 2016

# Sequence classification

"The food is usually not so good" ➡️ 👍👎

RNN

*Word embedding*    $W_e\, x_1$

Embed

*One-hot* encoding    $x_1$

"The"

# Sequence classification

# Sequence classification

# Sequence classification



http://deeplearning.net/tutorial/lstm.html

# Sequence classification

# Image caption generation



Training time

- Maximize likelihood of reference captions



| | |
|---|---|
| $\log p_1(S_1)$  $\log p_2(S_2)$  $\log p_N(S_N)$ | Log-likelihood of next reference word |
| $p_1$  $p_2$  $p_N$ | Softmax probability over next word |
| $W_e S_0$  $W_e S_1$  $W_e S_{N-1}$ | Word embedding |
| $S_0$  $S_1$  $S_{N-1}$ | Words of reference caption (one-hot encoding) |

O. Vinyals, A. Toshev, S. Bengio, D. Erhan, Show and Tell: A Neural Image Caption Generator, CVPR 2015

# Image caption generation: Example outputs



A person riding a motorcycle on a dirt road.

Two dogs play in the grass.

A skateboarder does a trick on a ramp.

A dog is jumping to catch a frisbee.

A group of young people playing a game of frisbee.

Two hockey players are fighting over the puck.

A little girl in a pink hat is blowing bubbles.

A refrigerator filled with lots of food and drinks.

A herd of elephants walking across a dry grass field.

A close up of a cat laying on a couch.

A red motorcycle parked on the side of the road.

A yellow school bus parked in a parking lot.

| Describes without errors | Describes with minor errors | Somewhat related to the image | Unrelated to the image |

# Sequence-to-sequence models with attention



Many slides adapted from J. Johnson

"We are eating bread" ➡ "Estamos comiendo pan"

# Sequence-to-sequence with RNNs

**Decoder:** $s_t = g_U(y_{t-1}, s_{t-1}, c)$

**Encoder:** $h_t = f_W(x_t, h_{t-1})$

From final hidden state predict:
**Initial decoder state** $s_0$
**Context vector** $c$ (often $c = h_T$)

# Sequence-to-sequence with RNNs

**Decoder:** $s_t = g_U(y_{t-1}, s_{t-1}, c)$

estamos    comiendo    pan    [STOP]

From final hidden state predict:
**Initial decoder state** $s_0$
**Context vector** $c$ (often $c = h_T$)

**Encoder:** $h_t = f_W(x_t, h_{t-1})$

$h_1$    $h_2$    $h_3$    $h_4$    $s_0$    $y_1$    $y_2$    $y_3$    $y_4$

$s_1$    $s_2$    $s_3$    $s_4$

$x_1$    $x_2$    $x_3$    $x_4$    $c$    $y_0$    $y_1$    $y_2$    $y_3$

we    are    eating    bread    [START]    estamos    comiendo    pan

Problem: Input sequence bottlenecked through fixed-sized vector

Idea: use *new* context vector at each step of decoder!

A. Sutskever, O. Vinyals, Q. Le, Sequence to sequence learning with neural networks, NeurIPS 2014

# Sequence-to-sequence with RNNs *and attention*

- Intuition: translation requires *alignment*

# Sequence-to-sequence with RNNs and attention

- At each timestep of decoder, context vector "looks at" different parts of the input sequence

# Sequence-to-sequence with RNNs and attention



Normalize to get **attention weights** $a_{t,i}$

Compute **context vector** as
$$c_t = \sum_i a_{t,i} h_i$$

Compute scalar **alignment scores**
$$e_{t,i} = f_{\text{att}}(s_{t-1}, h_i)$$

Use context vector in decoder:
$$s_t = g_U(y_{t-1}, s_{t-1}, c_t)$$

**Intuition**: Context vector "attends" to the relevant part of the input sequence *"estamos"* = *"we are"* so maybe $a_{11} = a_{12} = 0.45$, $a_{13} = a_{14} = 0.05$

we    are    eating    bread

[START]

D. Bahdanau, K. Cho, Y. Bengio, Neural Machine Translation by Jointly Learning to Align and Translate, ICLR 2015

# Sequence-to-sequence with RNNs and attention

# Sequence-to-sequence with RNNs and attention

# Sequence-to-sequence with RNNs and attention

# Sequence-to-sequence with RNNs and attention

# Sequence-to-sequence with RNNs and attention

- Visualizing attention weights (English source, French target):



Same word order in source and target languages

Flipped word order

Verb conjugation is different

D. Bahdanau, K. Cho, Y. Bengio, Neural Machine Translation by Jointly Learning to Align and Translate, ICLR 2015

# Attention: Weights Visualization



**Decoder RNN (target language: French)**

l'   accord   sur   la   zone   économique   européenne   a   été   signé   en   août   1992   .   \<end\>

the   agreement   on   the   European   Economic   Area   was   signed   in   August   1992   .   \<end\>

**Encoder RNN (source language: English)**

Figure is from https://distill.pub/2016/augmented-rnns/

# Quantitative evaluation



With attention (trained with sentence length <= 50)

With attention (trained with sentence length <= 30)

No attention

D. Bahdanau, K. Cho, Y. Bengio, Neural Machine Translation by Jointly Learning to Align and Translate, ICLR 2015

# Generalizing attention



$$c_t = \sum_i a_{t,i} h_i$$

- The decoder doesn't use the fact that the $h_i$ form an ordered sequence – it just treats them as an unordered set

- Can use similar architecture given any set of input hidden vectors $\{h_i\}$!

# Image captioning with RNNs and attention

- Idea: pay attention to different parts of the image when generating different words

- Automatically learn this *grounding* of words to image regions without direct supervision

K. Xu et al., Show, Attend and Tell: Neural Image Caption Generation with Visual Attention, ICML 2015

# Image captioning with RNNs and attention

Alignment scores

| $e_{1,1,1}$ | $e_{1,1,2}$ | $e_{1,1,3}$ |
|---|---|---|
| $e_{1,2,1}$ | $e_{1,2,2}$ | $e_{1,2,3}$ |
| $e_{1,3,1}$ | $e_{1,3,2}$ | $e_{1,3,3}$ |

Attention weights

| $a_{1,1,1}$ | $a_{1,1,2}$ | $a_{1,1,3}$ |
|---|---|---|
| $a_{1,2,1}$ | $a_{1,2,2}$ | $a_{1,2,3}$ |
| $a_{1,3,1}$ | $a_{1,3,2}$ | $a_{1,3,3}$ |

softmax

$$e_{t,i,j} = f_{att}(s_{t-1}, h_{i,j})$$

CNN

| $h_{1,1}$ | $h_{1,2}$ | $h_{1,3}$ |
|---|---|---|
| $h_{2,1}$ | $h_{2,2}$ | $h_{2,3}$ |
| $h_{3,1}$ | $h_{3,2}$ | $h_{3,3}$ |

$s_0$

Use CNN to extract a grid of features

K. Xu et al., Show, Attend and Tell: Neural Image Caption Generation with Visual Attention, ICML 2015

# Image captioning with RNNs and attention

Alignment scores

| $e_{1,1,1}$ | $e_{1,1,2}$ | $e_{1,1,3}$ |
|---|---|---|
| $e_{1,2,1}$ | $e_{1,2,2}$ | $e_{1,2,3}$ |
| $e_{1,3,1}$ | $e_{1,3,2}$ | $e_{1,3,3}$ |

Attention weights

| $a_{1,1,1}$ | $a_{1,1,2}$ | $a_{1,1,3}$ |
|---|---|---|
| $a_{1,2,1}$ | $a_{1,2,2}$ | $a_{1,2,3}$ |
| $a_{1,3,1}$ | $a_{1,3,2}$ | $a_{1,3,3}$ |

softmax

$$e_{t,i,j} = f_{\text{att}}(s_{t-1}, h_{i,j})$$

cat

$y_1$

CNN

| $h_{1,1}$ | $h_{1,2}$ | $h_{1,3}$ |
|---|---|---|
| $h_{2,1}$ | $h_{2,2}$ | $h_{2,3}$ |
| $h_{3,1}$ | $h_{3,2}$ | $h_{3,3}$ |

$s_0$

$s_1$

$c_1$

$y_0$

[START]

$$c_t = \sum_i a_{t,i,j} h_i$$

# Image captioning with RNNs and attention

Alignment scores

| $e_{2,1,1}$ | $e_{2,1,2}$ | $e_{2,1,3}$ |
|---|---|---|
| $e_{2,2,1}$ | $e_{2,2,2}$ | $e_{2,2,3}$ |
| $e_{2,3,1}$ | $e_{2,3,2}$ | $e_{2,3,3}$ |

Attention weights

| $a_{2,1,1}$ | $a_{2,1,2}$ | $a_{2,1,3}$ |
|---|---|---|
| $a_{2,2,1}$ | $a_{2,2,2}$ | $a_{2,2,3}$ |
| $a_{2,3,1}$ | $a_{2,3,2}$ | $a_{2,3,3}$ |

softmax

$$e_{t,i,j} = f_{\text{att}}(s_{t-1}, h_{i,j})$$

cat    sitting

$y_1$    $y_2$

Use a CNN to compute a grid of features for an image

CNN

| $h_{1,1}$ | $h_{1,2}$ | $h_{1,3}$ |
|---|---|---|
| $h_{2,1}$ | $h_{2,2}$ | $h_{2,3}$ |
| $h_{3,1}$ | $h_{3,2}$ | $h_{3,3}$ |

$s_0$    $s_1$    $s_2$

$c_1$  $y_0$    $c_2$  $y_1$

[START]    cat

$$c_t = \sum_i a_{t,i,j} h_i$$

# Image captioning with RNNs and attention



$$e_{t,i,j} = f_{\text{att}}(s_{t-1}, h_{i,j})$$

Alignment scores

| $e_{t,1,1}$ | $e_{t,1,2}$ | $e_{t,1,3}$ |
| --- | --- | --- |
| $e_{t,2,1}$ | $e_{t,2,2}$ | $e_{t,2,3}$ |
| $e_{t,3,1}$ | $e_{t,3,2}$ | $e_{t,3,3}$ |

softmax

Attention weights

| $a_{t,1,1}$ | $a_{t,1,2}$ | $a_{t,1,3}$ |
| --- | --- | --- |
| $a_{t,2,1}$ | $a_{t,2,2}$ | $a_{t,2,3}$ |
| $a_{t,3,1}$ | $a_{t,3,2}$ | $a_{t,3,3}$ |

$$c_t = \sum_i a_{t,i,j} h_i$$

Each time step of decoder uses a different context vector that looks at different parts of the input image

# Example results

- Good captions



A woman is throwing a <u>frisbee</u> in a park.

A <u>dog</u> is standing on a hardwood floor.

A <u>stop</u> sign is on a road with a mountain in the background.

A little <u>girl</u> sitting on a bed with a teddy bear.

A group of <u>people</u> sitting on a boat in the water.

A giraffe standing in a forest with <u>trees</u> in the background.

# Example results

- Mistakes



A large white <u>bird</u> standing in a forest.

A woman holding a <u>clock</u> in her hand.

A man wearing a hat and a hat on a <u>skateboard</u>.

A person is standing on a beach with a <u>surfboard</u>.

A woman is sitting at a table with a large <u>pizza</u>.

A man is talking on his cell <u>phone</u> while another man watches.

# Recurrent vs. convolutional sequence models

- **Recurrent models:**
  - Treat input as ordered sequence (inherently sequential processing)
  - Build up context using the hidden vector

- **Convolutional models:**
  - Treat input as a grid indexed by time and feature dimension
  - Build up context using multiple layers of convolutions
  - Processing can be parallel at training time, but convolutions must be *causal*

# WaveNet

- Goal: generate raw audio
  - Represented as sequence of 16-bit integer values (can be quantized to 256 discrete levels), 16K samples per second
- Applications: text-to-speech, music generation
  - Also works for speech recognition
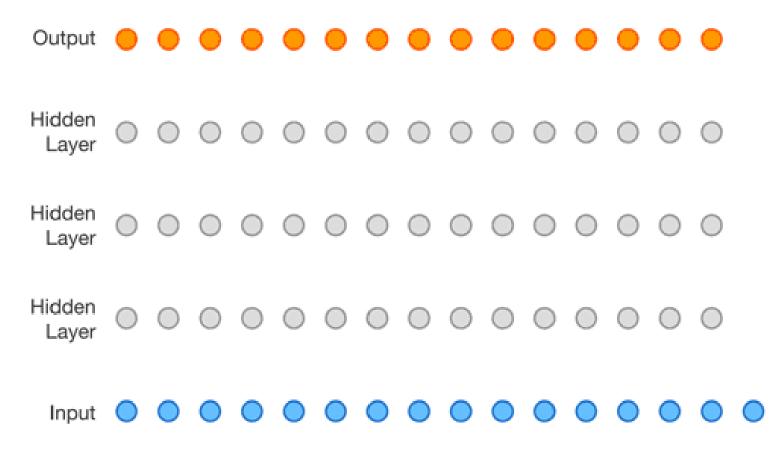


Figure 1: A second of generated speech.

A. van den Oord et al., WaveNet: A generative model for raw audio, arXiv 2016

# WaveNet

- Training time: compute predictions of all timesteps in parallel (conditioned on ground truth)

# WaveNet

- Test time: feed each predicted sample back into the model to make prediction at next timestep
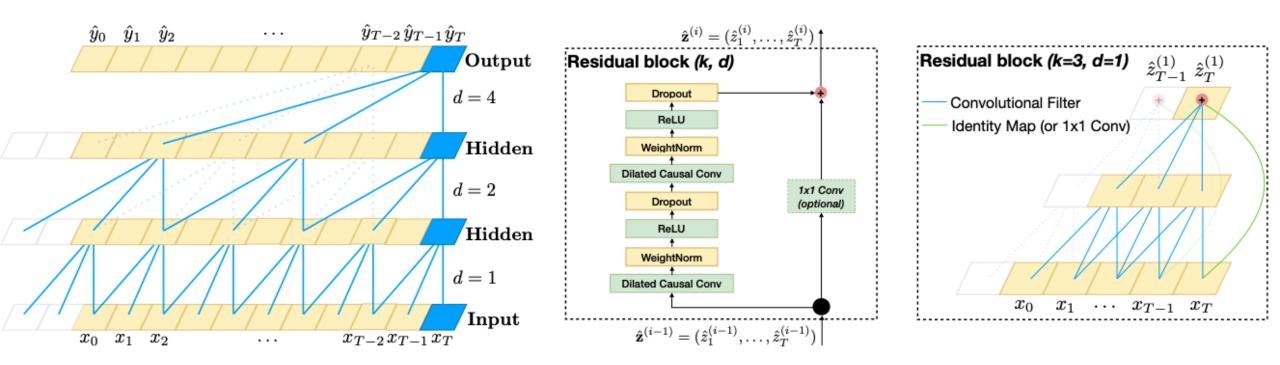
# WaveNet: Results

- Text-to-speech with different speaker identities:

- Generated sample of classical piano music:

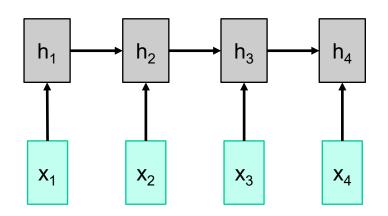# Temporal convolutional networks (TCNs)

- TCNs can be competitive with RNNs for a variety of sequence modeling tasks
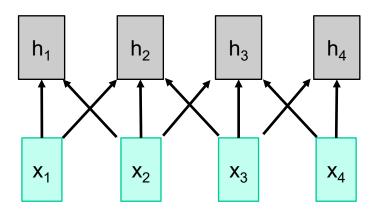


S. Bai, J. Kolter, and V. Koltun, An Empirical Evaluation of Generic Convolutional and Recurrent Networks for Sequence Modeling, arXiv 2018
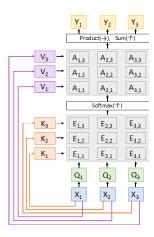
# Different ways of processing sequences

## RNN



Works on **ordered sequences**
- Pros: Good for long sequences: After one RNN layer, $h_T$ "sees" the whole sequence
- Con: Not parallelizable: need to compute hidden states sequentially
- Con: Hidden states have limited expressive capacity

## 1D convolutional network



Works on **multidimensional grids**
- Pro: Each output can be computed in parallel (at training time)
- Con: Bad at long sequences: Need to stack many conv layers for outputs to "see" the whole sequence

## Transformer



- Works on **sets of vectors**
- Pro: Good at long sequences: after one self-attention layer, each output "sees" all inputs!
- Pro: Each output can be computed in parallel (at training time)
- Con: Memory-intensive