# Models of sequential data

20th February, 2024

Alexey Zaytsev,
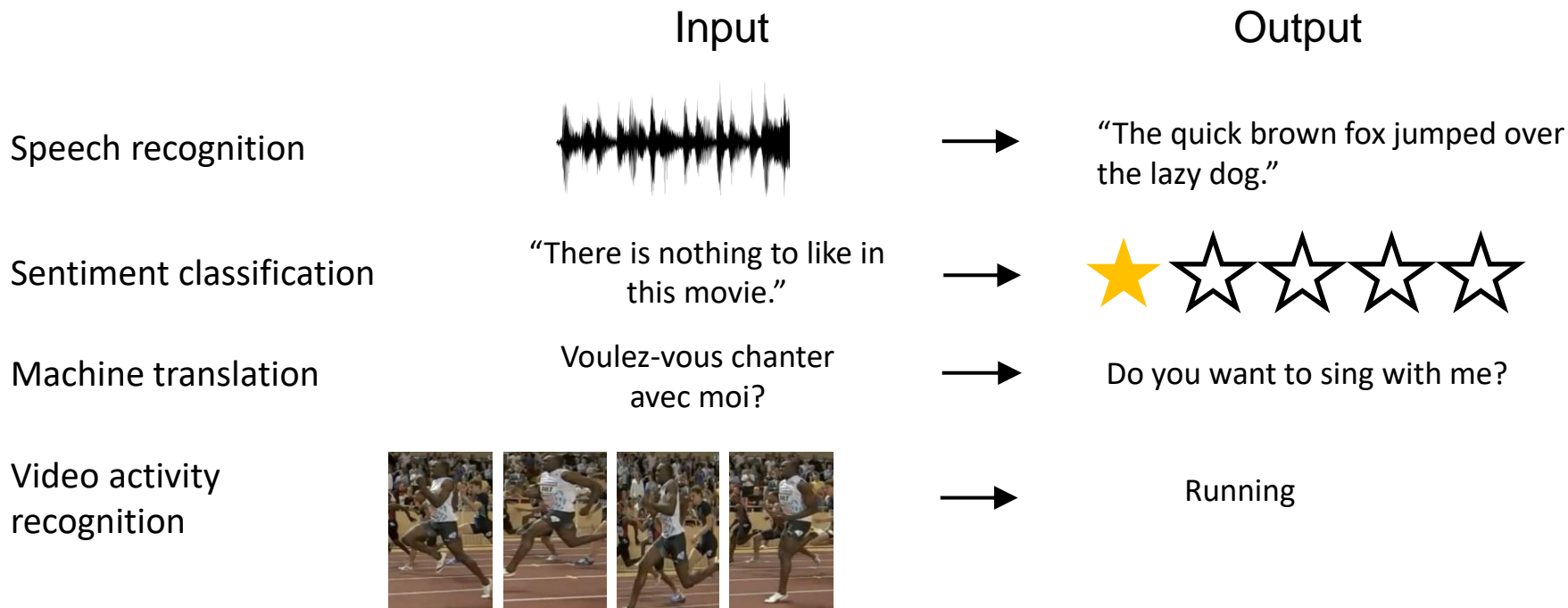
Assistant professor, Skoltech

Skoltech

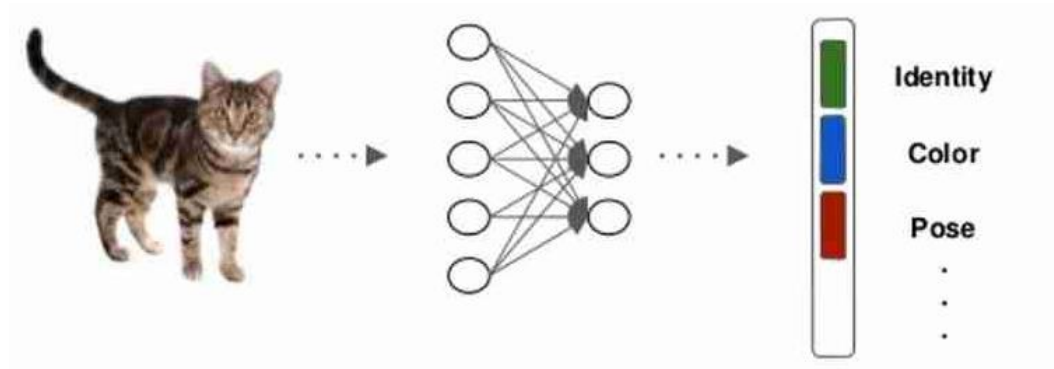# Common ways for classic ML application for time series data

1. Take input data including history for the target variable
2. Add differences, combinations, rolling means, medians, etc.
3. Add one-hot-encoding for important categorical features (day of week, holiday or not)

- Now we have input features for all points

- Let's apply our favorite ML regression algorithm

*dmlc*
**XGBoost**

Skoltech

# Deep Learning problems with sequential data: we need representations

Input | Output

Speech recognition → "The quick brown fox jumped over the lazy dog."

Sentiment classification → "There is nothing to like in this movie." → ★☆☆☆☆

Machine translation → Voulez-vous chanter avec moi? → Do you want to sing with me?

Video activity recognition → Running

Skoltech

Adopted from Andrew Ng

# Semi-structured data processing: why do we need NNs?

Skoltech

# Textbook example: next word prediction

The most complicated and difficult part of it was only just beginning.

Skoltech

# Textbook example: next word prediction
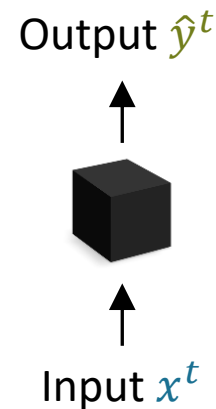
<u>Idea 1:</u> use previous word(s)

<u>Problem 1:</u> long-term dependencies

"**France** is where I grew up, but I now live in Boston. I speak fluent ____."

The most complicated and difficult part of it was only just beginning.

Feature representation: [0, 0, 0, 1, 0, 0]

Output $\hat{y}^t$

Input $x^t$

Skoltech

# Textbook example: next word prediction

Idea 1: use previous word(s)
Idea 2: use bag of words model

Problem 1: long-term dependencies

The most complicated and difficult part of it was only just beginning.

⬇

Feature representation:    [0, 3, 0, 2, 0, 0]

Bag of words: number of occurrences of each word

Adopted from Ava Soleimany, MIT

# Textbook example: next word prediction

Idea 1: use previous word(s)
Idea 2: use bag of words model

Problem 1: long-term dependencies
Problem 2: order preservation

The food was good, not bad at all.

vs.

The food was bad, not good at all.

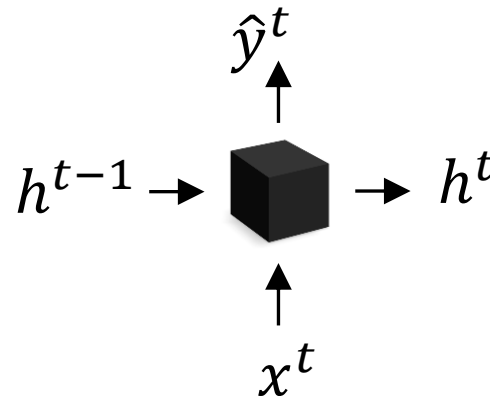The most complicated and difficult part of it was only just beginning.

Feature representation:   [0, 3, 0, 2, 0, 0]

Bag of words: number of occurrences of each word
(see also TF-IDF features)

# Model Design Criteria

1. Variable-length sequences processing

2. Long-term memory

3. Maintain order information

4. Natural preprocessing

$$\hat{y}^t$$

$$h^{t-1} \rightarrow \blacksquare \rightarrow h^t$$

$$x^t$$

**Recurrent Neural Networks** are the solution!

# Sequence processing with classic ML models

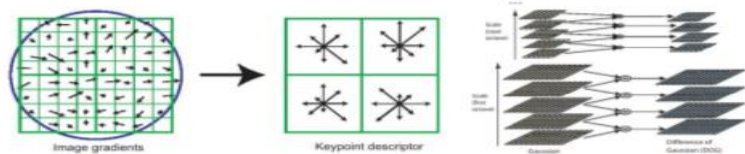1. Variable-length sequences processing
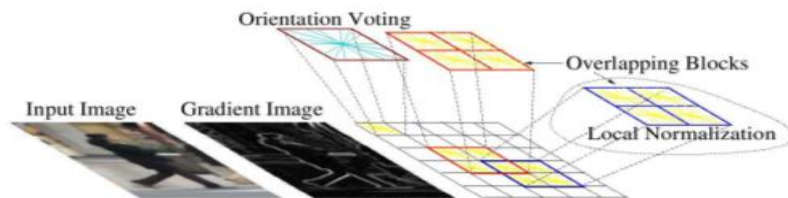
2. Long-term memory

3. Maintain order information

4. Natural preprocessing

**YES (if one to one)**

**NO** 🙁

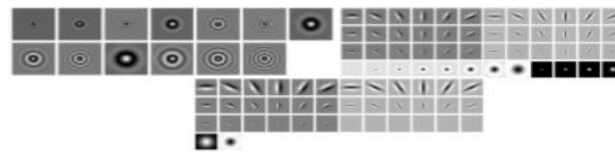**NO** 🙁

**a kind of**

Skoltech

# An art of feature construction



SIFT

Textons

HoG

RIFT

GIST

**No art now, just engineering**



SIFT

Textons

Orientation

Input Image    Gradient Image

HoG

RIFT

GIST

Deep Neural Networks:
a single concept, that
works better

Skoltech

# Sequence data examples

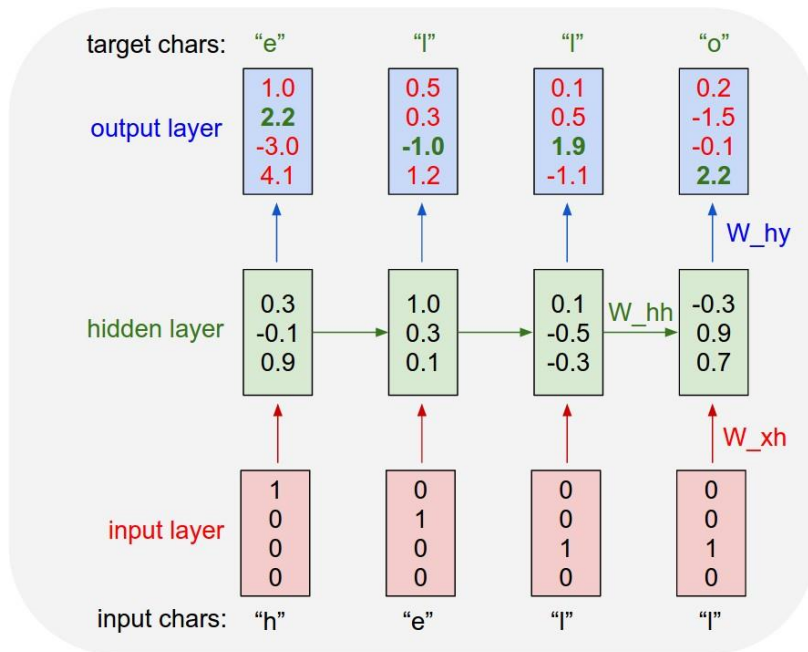| | Large data set | Semi-structured data |
|---|---|---|
| Speech recognition | + | + |
| Sentiment classification<br>"There is nothing to like in this movie." | + | + |
| DNA sequence analysis<br>AGCCCCTGTGAGGAACTAG | + | + |
| Machine translation<br>Voulez-vous chanter avec moi? | + | + |
| Video activity recognition | + | + |

Skoltech

# Examples of texts generated by LSTM (Long Short Term Memory NN)

- Shakespeare
- Wiki
- Algebraic geometrics articles
- Linux Source Code
- Dinosaurs names

# Conclusions

- Classic ML can't handle semi-structured data common in sequential data processing

- We can *learn representations* via Neural Networks

- Results are nice even for relatively simple models

Skoltech

# Machine translation: application example



Dog Translation Machine

# Machine translation, 50-s

**Cold war child: translator from Russian to English IBM 701 Translator**

Doctor Dostert predicted that "five, perhaps three years hence, interlingual meaning conversion by electronic process in important functional areas of several languages may well be an accomplished fact." (1954)

Rule-based approach that uses English-Russian dictionary

Project shut-down in 5 years: no significant progress



https://www.ibm.com/ibm/history/exhibits/701/701_translator.html
https://youtu.be/8ZtdVUB007A

# Statistical approach – the leading one before 2014

- Complicated and heavy model
- Many separate components
- Complex generation of inputs
- Support of a sophisticated system
- Quality is not great

Skoltech

# Neural network for machine translation

Input
sequence

Encoder

Embedding

Decoder

Output
sequence

**seq2seq** (sequence to sequence) architecture

Skoltech

# Vanilla Recurrent Neural Network



NO (LONG) MEMORY

# Simple RNN model

At each step:
$\hat{y}^t$ - output / model prediction
$x^t$ - input vector / new information
$h^t$ - cell / hidden state

Output $\qquad \hat{y}^{t-2} \qquad\qquad \hat{y}^{t-1} \qquad\qquad \hat{y}^t$

Cell state $\qquad h^{t-3} \rightarrow$

$h^{t-2} \qquad\qquad h^{t-1}$

Input vector $\qquad x^{t-2} \qquad\qquad x^{t-1} \qquad\qquad x^t$

# Simple RNN model block

Output $\hat{y}^t$

Cell
state $\boldsymbol{h}^{t-1} \rightarrow$ ■ $\rightarrow \boldsymbol{h}^t$

Input
vector $\boldsymbol{x}^t$

$$\boldsymbol{h}^t = f_h(\boldsymbol{x}^t, \boldsymbol{h}^{t-1})$$

$$\boldsymbol{h}^t = \tanh(V\boldsymbol{x}^t + W\boldsymbol{h}^{t-1} + b_h)$$

$$\hat{y}^t = f_y(\boldsymbol{h}^t)$$

$$\hat{y}^t = \mathrm{softmax}(U\boldsymbol{h}^t + b_y)$$

Skoltech

22

# Sequence processing with Vanilla RNN

- Long-term memory

- Maintain order information

- Natural preprocessing

- Variable-length sequences processing

**NO** 🙁

**YES**

**a kind of**

**YES (if one to one)**

Skoltech

# Another semi-structured data problems

Credit scoring:
default prediction

Money
Transactions data → ★☆☆☆☆

Fraud detection
in healthcare insurance

A history of visits → Is there a fraud?

Skoltech

# Problem statements for sequential data

$\hat{y}^1$

$x^1$

No recurrent structure

Sequence to a single output

$\hat{y}^T$

$x^1$     $x^2$     $\cdots$     $x^T$

Sequence to a sequence of outputs

$\hat{y}^1$     $\hat{y}^2$     $\hat{y}^T$

$x^1$     $x^2$     $\cdots$     $x^T$

Skoltech

# A sequence to sequence problem

Input sequence

$x^1$
$x^2$
$\ldots$
$x^{T_x}$

Encoder

Embedding

Decoder

$y^1$
$y^2$
$\ldots$
$y^{T_y}$

Output sequence

Skoltech

$\hat{y}^1$  $\hat{y}^2$  $\hat{y}^{T_y}$

$x^1$  $x^2$  $x^{T_x}$

# More on Recurrent Neural Networks

**Forward propagation through Recurrent Neural Network**



Outputs: $\hat{y}^1$ $\hat{y}^2$ $\hat{y}^3$ $\hat{y}^{T_y}$

$h^0$ $h^1$ $h^2$ $h^{T_x-1}$ $\cdots$

Inputs: $x^1$ $x^2$ $x^3$ $x^{T_x}$

# Forward propagation through Recurrent Neural Network

Outputs:

$\hat{y}^1$    $\hat{y}^2$    $\hat{y}^3$    $\hat{y}^{T_y}$

$h^0$ $\quad$ $h^1$ $\quad$ $h^2$ $\quad$ $h^{T_x-1}$ $\quad$ ...

Inputs: $\quad$ $x^1$ $\quad$ $x^2$ $\quad$ $x^3$ $\quad$ $x^{T_x}$

Processing unit is the same for all time moments.
Units has parameters we want to learn!

Skoltech

Adopted from Andrew Ng

# Forward propagation through Recurrent Neural Network

Outputs:

$\hat{y}^1$  $\hat{y}^2$  $\hat{y}^3$  $\hat{y}^{T_y}$

$h^0$  $h^1$  $h^2$  $h^{T_x-1}$  $\cdots$

Information flows between states

Inputs:  $x^1$  $x^2$  $x^3$  $x^{T_x}$

Processing unit is the same for all time moments.
Units has parameters we want to learn!

Adopted from Andrew Ng

**Forward propagation through Recurrent Neural Network**

Outputs:

$$\hat{y}^1 \qquad \hat{y}^2 \qquad \hat{y}^3 \qquad\qquad\qquad \hat{y}^{T_y}$$

$$h^0 \longrightarrow \boxed{\begin{matrix}\bigcirc\\\bigcirc\\\bigcirc\\\bigcirc\end{matrix}} \overset{h^1}{\longrightarrow} \boxed{\begin{matrix}\bigcirc\\\bigcirc\\\bigcirc\\\bigcirc\end{matrix}} \overset{h^2}{\longrightarrow} \boxed{\begin{matrix}\bigcirc\\\bigcirc\\\bigcirc\\\bigcirc\end{matrix}} \longrightarrow \quad \dots \quad \overset{h^{T_x-1}}{\longrightarrow} \boxed{\begin{matrix}\bigcirc\\\bigcirc\\\bigcirc\\\bigcirc\end{matrix}}$$

Inputs:

$$x^1 \qquad\quad x^2 \qquad\quad x^3 \qquad\qquad\qquad\quad x^{T_x}$$

Input
For NLP
problems

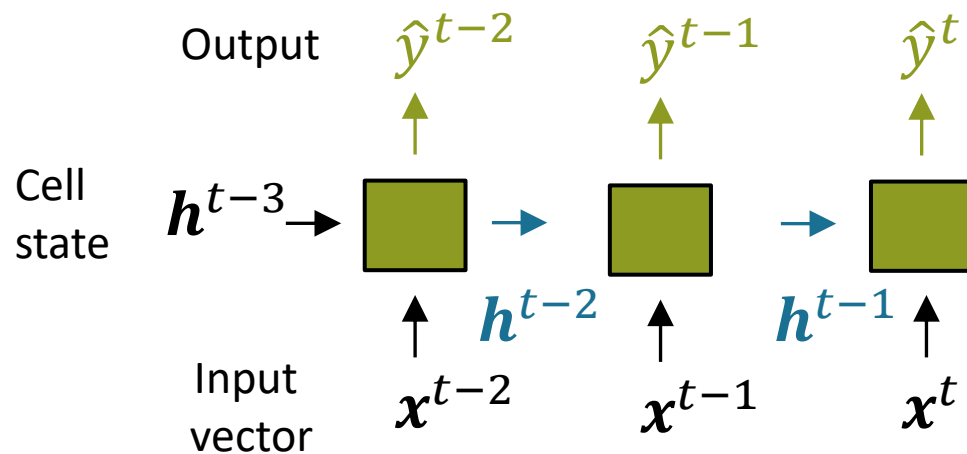It          is          truth          universally …
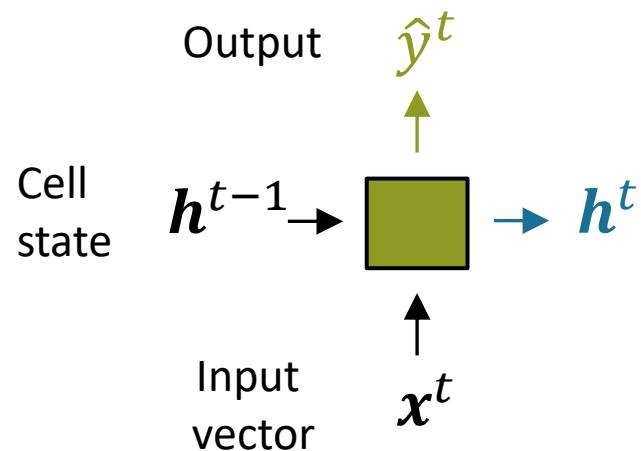
Skoltech

Adopted from Andrew Ng

# Simple RNN model

At each step:
$\hat{y}^t$ - output / model prediction
$x^t$ - input vector / new information
$h^t$ - cell / hidden state

Output

$\hat{y}^{t-2}$     $\hat{y}^{t-1}$     $\hat{y}^t$

Cell state

$h^{t-3} \rightarrow$ [ ] $\rightarrow$ [ ] $\rightarrow$ [ ]

$h^{t-2}$     $h^{t-1}$

Input vector

$x^{t-2}$     $x^{t-1}$     $x^t$

# Simple RNN model block

Output $\hat{y}^t$

Cell state $\boldsymbol{h}^{t-1}$ → ☐ → $\boldsymbol{h}^t$

Input vector $\boldsymbol{x}^t$

$$\boldsymbol{h}^t = f_h(\boldsymbol{x}^t, \boldsymbol{h}^{t-1})$$

$$\boldsymbol{h}^t = \tanh(V\boldsymbol{x}^t + W\boldsymbol{h}^{t-1} + b_h)$$

$$\hat{y}^t = f_y(\boldsymbol{h}^t)$$

$$\hat{y}^t = \text{softmax}(U\boldsymbol{h}^t + b_y)$$

**Skoltech**

# Backward propagation through Recurrent Neural Network



Loss function: $\quad L^1(\hat{y}^1, y^1) \qquad L^2(\hat{y}^2, y^2) \qquad L^3(\hat{y}^3, y^3) \qquad\qquad\qquad L^{T_y}(\hat{y}^{T_y}, y^{T_y})$

Outputs: $\qquad\qquad \hat{y}^1 \qquad\qquad\quad \hat{y}^2 \qquad\qquad\quad \hat{y}^3 \qquad\qquad\qquad\qquad \hat{y}^{T_y}$

$h^0 \qquad\qquad\quad h^1 \qquad\qquad h^2 \qquad\qquad \dots \qquad h^{T_x-1}$

Inputs: $\qquad\qquad x^1 \qquad\qquad\quad x^2 \qquad\qquad\quad x^3 \qquad\qquad\qquad\qquad x^{T_x}$

Skoltech

Adopted from Andrew Ng

# Backward propagation through Recurrent Neural Network



Loss function: $L^1(\hat{y}^1, y^1)$    $L^2(\hat{y}^2, y^2)$    $L^3(\hat{y}^3, y^3)$        $L^{T_y}(\hat{y}^{T_y}, y^{T_y})$

Outputs: $\hat{y}^1$    $\hat{y}^2$    $\hat{y}^3$        $\hat{y}^{T_y}$

$h^0$    $h^1$    $h^2$    $h^{T_x - 1}$

Inputs: $x^1$    $x^2$    $x^3$        $x^{T_x}$

Skoltech

Adopted from Andrew Ng

# Backpropagation w.r.t. U

$$L = \sum_{i=1}^{T_y} L^i\left(\hat{y}^i, y^i\right)$$

$$\frac{\partial L}{\partial U} = \sum_{i=1}^{T_y} \frac{\partial L_i}{\partial U} = \sum_{i=1}^{T_y} \frac{\partial L_i}{\partial \hat{y}^i} \frac{\partial \hat{y}^i}{\partial U}$$

$$\hat{y}^t = \text{softmax}(U\boldsymbol{h}^t + b_y)$$

$$L^t(\hat{y}^t, y^t)$$

$$\hat{y}^t$$

$$\boldsymbol{h}^{t-1} \rightarrow \boxed{\phantom{xx}} \rightarrow \boldsymbol{h}^t$$

$$\boldsymbol{x}^t$$

Skoltech

# Backpropagation w.r.t. W

$$\hat{y}^t = \text{softmax}(U\boldsymbol{h}^t + b_y)$$

$$L = \sum_{i=1}^{T_y} L^i(\hat{y}^i, y^i)$$

$$\frac{\partial L}{\partial W} = \sum_{i=1}^{T_y} \frac{\partial L_i}{\partial W} = \sum_{i=1}^{T_y} \frac{\partial L_i}{\partial \hat{y}^i} \frac{\partial \hat{y}^i}{\partial W}$$

$$\boldsymbol{h}^t = \tanh(V\boldsymbol{x}^t + W\boldsymbol{h}^{t-1} + b_h)$$

$$\frac{\partial L_i}{\partial \hat{y}^i} \frac{\partial \hat{y}^i}{\partial W} = \frac{\partial L_i}{\partial \hat{y}^i} \frac{\partial \hat{y}^i}{\partial h_t} \left( \frac{\partial h_t}{\partial W} + \frac{\partial h_t}{\partial h_{t-1}} \frac{\partial h_{t-1}}{\partial W} + \right)$$
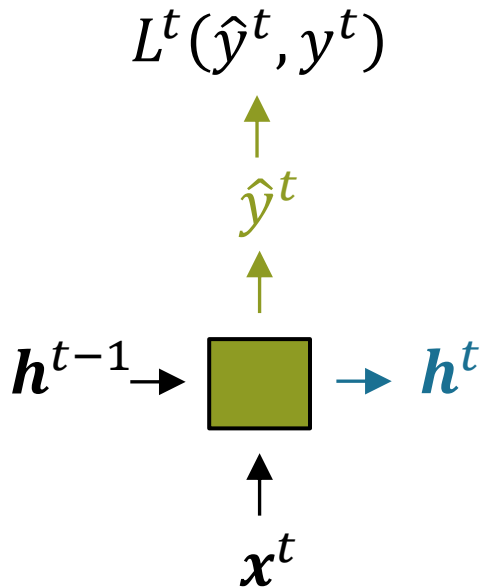
$$\frac{\partial L_i}{\partial \hat{y}^i} \frac{\partial \hat{y}^i}{\partial W} = \frac{\partial L_i}{\partial \hat{y}^i} \frac{\partial \hat{y}^i}{\partial h_t} \sum_{i=0}^{T_y} \left( \prod_{j=i+1}^{T_y} \frac{\partial h_j}{\partial h_{j-1}} \right) \frac{\partial h_i}{\partial W}$$

$$L^t(\hat{y}^t, y^t)$$

$$\hat{y}^t$$

$$\boldsymbol{h}^{t-1} \rightarrow \quad \rightarrow \boldsymbol{h}^t$$

$$\boldsymbol{x}^t$$

Skoltech

37

# Problems of classic RNN

**Solution:**

- **Gradient clipping to scale big gradients**

1. $g = \dfrac{\partial L}{\partial W}$

2. If $g > t$ for some threshold $t$:

$$g = \frac{t}{\|g\|}\, g$$

Threshold $t$ is selected given the dynamic of loss function over iterations



Train Cost (with Adam)

14-14x28-4 net
14-28x6-4 net
14-100-40-4 net

Skoltech

38

# Problems of classic RNN

Solution:

- Gradient clipping to scale big gradients

- Truncated backpropagation through time

Chunk

$$\hat{y}^1 \quad \hat{y}^2 \quad \hat{y}^3 \quad \hat{y}^4 \quad \hat{y}^5 \quad \hat{y}^6$$

$$h^0$$

$$x^1 \quad x^2 \quad x^3 \quad x^4 \quad x^5 \quad x^6$$

**Skoltech**

# Problems of classic RNN

**Gradients vanishing is a more**

   **serious problem**

**Many values < 1**

**Product << 1**

**Bias parameters to capture
long-term dependencies**

**Hard to detect!**

**Tricks:**

**Activation functions**

- **Use ReLU**

**Parameter initialization**

- **Initialize weights to identity matrix**

- **Initialize biases to zero**

Skoltech

# Another big problem of classic RNN

**Problem:** Neural networks forget fast, and it is hard to learn long-term dependencies

**Solution: Gated architectures**

More complex recurrent units with gates to control what information is passed through

- GRU (Gated Recurrent Unit)
- LSTM (Long-Short Term Memory)

Skoltech

# Selection of RNN architecture

Skoltech

# Better RNN units: LSTM and GRU

- LSTM: long short term memory [1]
- GRU: Gated recurrent unit [2]



Skoltech

1. Schmidhuber, J., & Hochreiter, S. (1997). Long short-term memory. *Neural Comput*, *9*(8), 1735-1780.
2. Cho, K., Van Merriënboer et al. (2014). Learning phrase representations using RNN encoder-decoder for statistical machine translation. *arXiv preprint arXiv:1406.1078*.

# Details on how LSTM works

Remembering information for long periods of time is practically the default behavior of LSTM

LSTM was proposed by J. Schmidhuber group in 1991

http://colah.github.io/posts/2015-08-Understanding-LSTMs/

# Long term memory part – Cell state



Forget part    Update part

Cell state before    $C_{t-1}$    $C_t$    Cell state after

http://colah.github.io/posts/2015-08-Understanding-LSTMs/

# Forget part

Identify how much should we forget: sigmoid
returns value between 0 and 1

hidden
state (a)

inputs

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

# Update part

Identify how much should we update:
sigmoid returns value between 0 and 1



hidden
state (a)

inputs        What update we do

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

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

Skoltech

# Long term memory part – Cell state

Forget part    Update part

Cell state
before

Cell state
after

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



Skoltech

# Update everything else



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

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

http://colah.github.io/posts/2015-08-Understanding-LSTMs/

# Details on how LSTM works

- There are cell and hidden (activation) states
- LSTM block forgets and updates cell state during processing at one block



Skoltech

## GRU – Gated Recurrent Unit

- Update gate – what to pay attention to
- Reset gate – what to forget

- Slightly worse than LSTM for NLP – but not in all problems
- Simpler and cheaper than LSTM

Forget part     Update part

Cho, Kyunghyun et al. Learning Phrase Representations using RNN Encoder-Decoder for Statistical Machine Translation. EMNLP.=-098]-[9o08i=-89-p9o0i8o0-9-9-980876-=-=-=-=-=-=[]jnhg76543209876543875r4e6666667y5yt2014.

Skoltech

# Multilayer architectures

# Other architectures: bidirectional LSTM

# Other architectures: bidirectional LSTM

Bidirectional LSTM are useful when we benefit from the future data or can use it:
- Handwriting Recognition
- Speech Recognition
- Protein Structure Prediction (Bioinformatics)

# Multiple layers RNN performance



ELO dataset (Boxplot with depths; ylim)

# Representation learning is the core feature of Neural Networks



3

224

224

Image

Feature Extractor

512

Representation

Skoltech

**Representation learning is still here for Recurrent Neural Networks**

Outputs:

$\hat{y}^1$     $\hat{y}^2$     $\hat{y}^3$     $\hat{y}^{T_y}$

$h^0$   $h^1$   $h^2$   $\boxed{h^{T_x-1}}$

...   We also have a sequence representation

Inputs:   $\boxed{x^1}$    $x^2$    $x^3$    $x^{T_x}$

Most of the time we also learn representations of objects in an end2end manner with backpropagation

Skoltech

# Attention mechanism

# New state of the art:
# **attention** is all we need



attention

Higher BLEU scores are better



English German Translation quality



English French Translation Quality

Skoltech

https://ai.googleblog.com/2017/08/transformer-novel-neural-network.html

# Bottleneck in seq2seq models



All information about the sequence is in this vector

# Attention

- Solution to the bottleneck problem

- Direction connection between parts of input and output sequence

Skoltech

# Sequence 2 sequence with attention

Attention
output

Attention
distribution

Attention
scores

**Skoltech**

$$h^0 \rightarrow \boxed{\phantom{x}} \rightarrow \cdots \rightarrow \boxed{\phantom{x}} \rightarrow \cdots \rightarrow \boxed{\phantom{x}} \rightarrow$$

$$x^1 \qquad x^{T_x} \qquad \text{<START>}$$

$s_1$

Encoder                          Decoder

# Sequence 2 sequence with attention

Attention output

Attention distribution

Attention scores

softmax

$\hat{y}^1$

$s_1$

$h^0 \rightarrow$

$x^1$

$x^{T_x}$

&lt;START&gt;

Encoder

Decoder

Skoltech

# Sequence 2 sequence with attention

# Attention: formulas

- First RNN produces encoder hidden states $\boldsymbol{h}_1, \dots, \boldsymbol{h}_{T_x} \in \mathbb{R}^h$
- Decoder hidden state $\boldsymbol{s}_t \in \mathbb{R}^h$ at time step $t$
- Attention scores for step $t$:
$$\boldsymbol{e^t} = [\boldsymbol{s}_t^T \boldsymbol{h}_1, \dots, \boldsymbol{s}_t^T \boldsymbol{h}_{T_x}] \in \mathbb{R}^{T_x}$$
- Softmax to get attention distribution: all values are positive, sum of all values is 1:
$$\boldsymbol{\alpha^t} = \text{softmax}(\boldsymbol{e^t}) \in \mathbb{R}^{T_x}$$
- Attention output $\boldsymbol{a}_t$ is a weighted sum of hidden states:
$$\boldsymbol{a}_t = \sum_{i=1}^{T_x} \alpha_i^t \boldsymbol{h}_i \in \mathbb{R}^h$$

- We concatenate the attention output $\boldsymbol{a}_t$ with the decoder hidden state $\boldsymbol{s}_t$ and proceed to the non-attention part of our seq2seq model
$$[\boldsymbol{a}_t, \boldsymbol{s}_t] \in \mathbb{R}^{2h}$$

# Attention is just great

- Significantly improves performance of NMT
- Solves the bottleneck problem
  - All encoder tokens are connected to all decoder tokens
- No more vanishing gradients
  - All to All connection
- Provides some interpretability
  - see alignment figure

- Similar to RNN seq2seq, but greater!



English French Translation Quality

BLEU scores (higher is better) of single models on the standard WMT newstest2014 English to French translation benchmark.

# Attention is a general deep learning idea

We can use attention in many architectures and many tasks
- Other NLP problems
- Sequential data processing
- Graph Neural Networks

Key value interpretation:

$s_i$ - query to a database               Hidden state of the decoder

$k_i$ - keys in the database               Hidden state of the encoder

$h_i$ - values in the database             Hidden state of the encoder

- Calculate correspondence $e(s_i, k_i)$
- Calculate weights on the base of correspondence values
- Extract information as weighted sum of values $\sum_{i=1}^{n} \alpha_i \boldsymbol{h}_i$

Skoltech

# Interpretability of attention

Remove a token $i^*$ with max attention…



**Remove random: Decision flip?**

| | Yahoo | | | | IMDB | |
| --- | --- | --- | --- | --- | --- | --- |
| | Yes | No | | | Yes | No |
| Yes | 0.5 | 8.7 | | Yes | 2.2 | 12.2 |
| No | 1.3 | 89.6 | | No | 1.4 | 84.2 |

| | Amazon | | | | Yelp | |
| --- | --- | --- | --- | --- | --- | --- |
| | Yes | No | | | Yes | No |
| Yes | 2.7 | 7.6 | | Yes | 1.5 | 8.9 |
| No | 2.7 | 87.1 | | No | 1.9 | 87.7 |

Remove $i^*$: Decision flip?

Figure source:
https://jalammar.github.io/illustrated-transformer/,
https://github.com/jessevig/bertviz

Serrano, Sofia, and Noah A. Smith. "Is attention interpretable?." arXiv preprint arXiv:1906.03731 (2019). ACL 2019.

# General attention idea

# Attention

We need a representation of an object.

A "coarse" description of an object is available.

For other objects we have more detailed stored information.

We use attention to extract these information

# Single-object key value interpretation

$\boldsymbol{q}_i$ - query to a database

$\boldsymbol{k}_j$ - keys in the database

$\boldsymbol{v}_j$ - values in the database

We calculate attention scores

$\alpha_j = e(\boldsymbol{q}_i, \boldsymbol{k}_j) = \boldsymbol{s}_i^T \boldsymbol{k}_j$ (other distances also possible)

$\boldsymbol{\alpha} = \text{softmax}(\boldsymbol{\alpha})$

Then we extract the information as weighted sum of values

$\mathbf{a}_i = \sum_{j=1}^{T_x} \alpha_j \boldsymbol{v}_j$

# Matrix key value interpretation

$q_i$ - query to a database

$k_j$ - keys in the database

$v_j$ - values in the database

We calculate correspondences

$$A(q, K, V) = \sum_i \frac{\exp(q_i^T k_j)}{\sum_l \exp(q_i^T k_l)} v_j$$

$$A(Q, K, V) = \text{softmax}(QK^T)V$$

# "Databases" for the attention

- Nodes in a graph

- Tokens in a sequence (we need to specify the position, as the order is important)

- Sets of objects

# Single-object key value interpretation

$\boldsymbol{q}_i$ - query to a database        Hidden state of the *decoder*

$\boldsymbol{k}_j$ - keys in the database        Hidden state of the *encoder*

$\boldsymbol{v}_j$ - values in the database       Hidden state of the *encoder*

We calculate attention scores

$\alpha_j = e(\boldsymbol{q}_i, \boldsymbol{k}_j) = \boldsymbol{s}_i^T \boldsymbol{k}_j$  (other distances also possible)

$\boldsymbol{\alpha} = \text{softmax}(\boldsymbol{\alpha})$

Then we extract the information as weighted sum of values

$\mathbf{a}_i = \sum_{j=1}^{T_x} \alpha_j \boldsymbol{v}_j$

# Scaled attention values

For large dimension of the space of keys $d_k$:

- Large variances dot products $\boldsymbol{q}_i^T \boldsymbol{k}_j$
- Softmax only pays attention to some keys
- Gradients are small, hard to learn

Old formula:
$$A(Q, K, V) = \mathrm{softmax}(QK^T)V$$

New scaled formula:
$$A(Q, K, V) = \mathrm{softmax}(QK^T/\sqrt{d_k})V$$
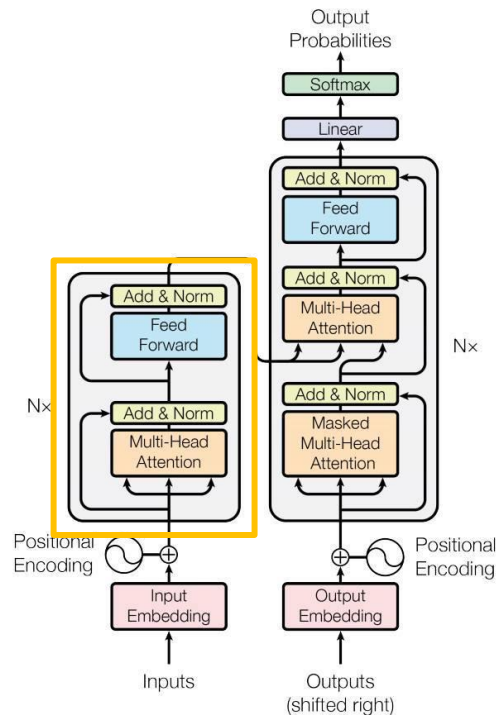
# Transformers with self-attention

# Transformer is based on the same idea

Now we completely drop RNN part

Also we repeat *self-attention* many times

Further we'll consider separate parts:
- Multi-head attention
- Feed Forward

Skoltech
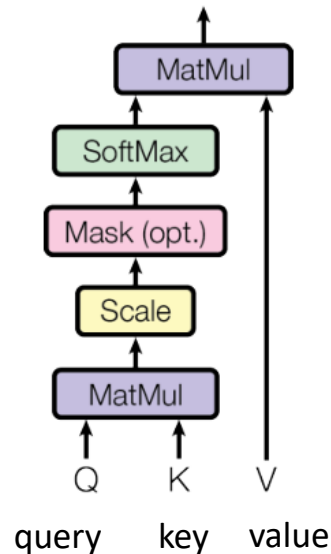
# Attention / Self-attention block

$$\text{Attention}(Q, K, V) = \text{softmax}(\frac{QK^T}{\sqrt{d_k}})V$$

$d_k$ is the dimension of query and key,
we scale to take control of large values of dot-product in high dimensions

A possible option is to replace scaled dot-product used here with additive attention: a single-hidden layer neural network.

Scaled Dot-Product Attention

# Self-attention block

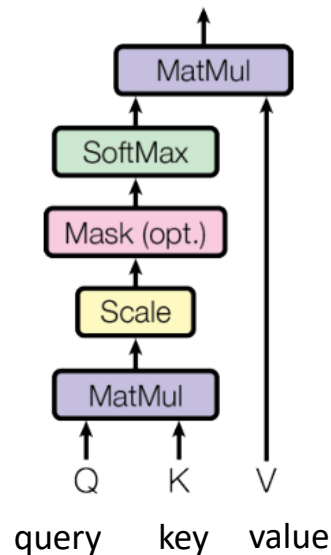$$\text{Attention}(Q, K, V) = \text{softmax}\left(\frac{QK^T}{\sqrt{d_k}}\right)V$$

We produce queries, keys, and values using initial word embeddings

$$Q = XW^Q, \dim(W^Q) = d_x \times d_q,$$

$$K = XW^K, \dim(W^Q) = d_x \times d_k,$$

$$V = XW^V, \dim(W^Q) = d_x \times d_v,$$

## Scaled Dot-Product Attention



query    key    value

Skoltech

# Multi-Head attention

$$\text{MultiHead}(Q, K, V) = \text{Concat}(\text{head}_1, ..., \text{head}_h)W^O$$
$$\text{where } \text{head}_i = \text{Attention}(QW_i^Q, KW_i^K, VW_i^V)$$

h heads in total

Scaled Dot-Product Attention



Multi-Head Attention



"Attention is all you need" paper

# Full block

Two linear transformation with ReLU activation in between
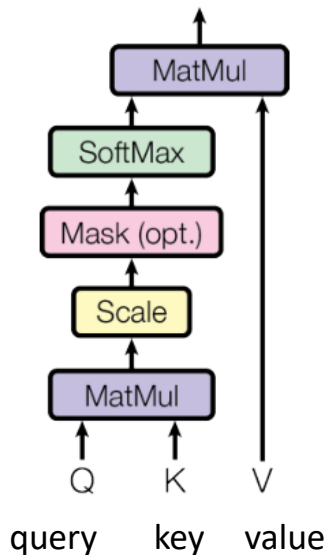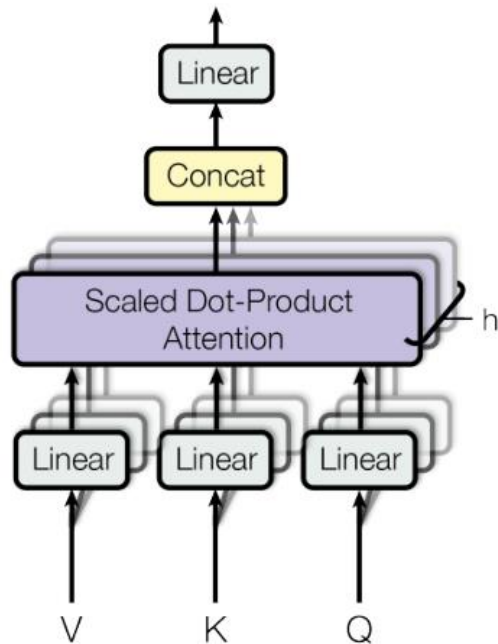
$$\text{FFN}(x) = \max(0, xW_1 + b_1)W_2 + b_2$$

Multi-head attention

$$\text{MultiHead}(Q, K, V) = \text{Concat}(\text{head}_1, ..., \text{head}_h)W^O$$
$$\text{where head}_i = \text{Attention}(QW_i^Q, KW_i^K, VW_i^V)$$

There are 6 consecutive *Full blocks* in the paper transformer architecture



LayerNorm$(x + \text{Sublayer}(x))$

Add & Norm

Feed Forward

N×

Add & Norm

Multi-Head Attention

**Skoltech**

"Attention is all you need" paper

# Full block: normalization and residual connection

$$\text{LayerNorm}(x + \text{Sublayer}(x))$$

Why? Speeds up training!

- Similar to Batch Normalization

- But can be used with batch size 1

- Can be used with RNNs and Transformers



$$\mu^l = \frac{1}{H} \sum_{i=1}^{H} a_i^l \qquad \sigma^l = \sqrt{\frac{1}{H} \sum_{i=1}^{H} \left(a_i^l - \mu^l\right)^2}$$

Layer Normalization, https://arxiv.org/abs/1607.06450

# Position encoding

In addition to usual embeddings of inputs we use position encoding to capture position

They are not one-hot vectors, as we want to handle various-length sequences
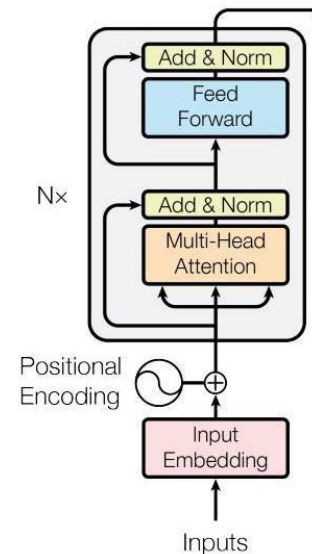
$$PE_{(pos,2i)} = sin(pos/10000^{2i/d_{model}})$$

$$PE_{(pos,2i+1)} = cos(pos/10000^{2i/d_{model}})$$

POSITIONAL ENCODING

| 0 | 0 | 1 | 1 |

| 0.84 | 0.0001 | 0.54 | 1 |

| 0.91 | 0.0002 | -0.42 | 1 |

EMBEDDINGS

$x_1$ $x_2$ $x_3$

INPUT  Je   suis   étudiant

Real example of positional encoding with a toy embedding size of 4

Add & Norm
Feed Forward
Nx
Add & Norm
Multi-Head Attention
Positional Encoding
Input Embedding
Inputs

Skoltech

# Transformer training

- Masked language model:

Replace random tokens with masks, try to reconstruct them using a Neural network
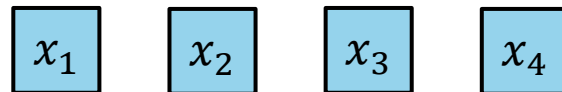
- Next token prediction

Predict next token

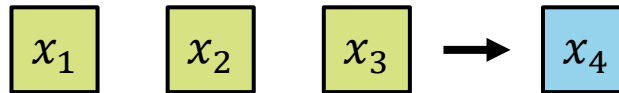**Self-supervised learning:** We don't need labeled examples, we just create them

BERT

$x_1$ $x_2$ $x_3$ $x_4$

↑ Transformer

$x_1$ MASK $x_3$ $x_4$

GPT

$x_1$ $x_2$ $x_3$ → $x_4$

Transformer

# Efficiency of transformers

| Model Name | $n_{params}$ | $n_{layers}$ | $d_{model}$ | $n_{heads}$ | $d_{head}$ | Batch Size | Learning Rate |
|---|---|---|---|---|---|---|---|
| GPT-3 Small | 125M | 12 | 768 | 12 | 64 | 0.5M | $6.0 \times 10^{-4}$ |
| GPT-3 Medium | 350M | 24 | 1024 | 16 | 64 | 0.5M | $3.0 \times 10^{-4}$ |
| GPT-3 Large | 760M | 24 | 1536 | 16 | 96 | 0.5M | $2.5 \times 10^{-4}$ |
| GPT-3 XL | 1.3B | 24 | 2048 | 24 | 128 | 1M | $2.0 \times 10^{-4}$ |
| GPT-3 2.7B | 2.7B | 32 | 2560 | 32 | 80 | 1M | $1.6 \times 10^{-4}$ |
| GPT-3 6.7B | 6.7B | 32 | 4096 | 32 | 128 | 2M | $1.2 \times 10^{-4}$ |
| GPT-3 13B | 13.0B | 40 | 5140 | 40 | 128 | 2M | $1.0 \times 10^{-4}$ |
| GPT-3 175B or "GPT-3" | 175.0B | 96 | 12288 | 96 | 128 | 3.2M | $0.6 \times 10^{-4}$ |

**Table 2.1:** Sizes, architectures, and learning hyper-parameters (batch size in tokens and learning rate) of the models which we trained. All models were trained for a total of 300 billion tokens.

| | Mean accuracy | 95% Confidence Interval (low, hi) | $t$ compared to control ($p$-value) | "I don't know" assignments |
|---|---|---|---|---|
| Control (deliberately bad model) | 86% | 83%–90% | - | 3.6 % |
| GPT-3 Small | 76% | 72%–80% | 3.9 (2e-4) | 4.9% |
| GPT-3 Medium | 61% | 58%–65% | 10.3 (7e-21) | 6.0% |
| GPT-3 Large | 68% | 64%–72% | 7.3 (3e-11) | 8.7% |
| GPT-3 XL | 62% | 59%–65% | 10.7 (1e-19) | 7.5% |
| GPT-3 2.7B | 62% | 58%–65% | 10.4 (5e-19) | 7.1% |
| GPT-3 6.7B | 60% | 56%–63% | 11.2 (3e-21) | 6.2% |
| GPT-3 13B | 55% | 52%–58% | 15.3 (1e-32) | 7.1% |
| GPT-3 175B | 52% | 49%–54% | 16.9 (1e-34) | 7.8% |

**Table 3.11: Human accuracy in identifying whether short ($\sim$200 word) news articles are model generated**. We find that human accuracy (measured by the ratio of correct assignments to non-neutral assignments) ranges from 86% on the control model to 52% on GPT-3 175B. This table compares mean accuracy between five different models, and shows the results of a two-sample T-Test for the difference in mean accuracy between each model and the control model (an unconditional GPT-3 Small model with increased output randomness).

Skoltech

# Transformer models are monstrous

2018 and 2019                                          2020          **175B**

**94M** **110M** 340M 465M 665M 330M **1.5B** 340M 355M 1.5B **8.3B** 1.5B **11B**  **2.6B** **17B** 66M **9.4B**

ELMo GPT BERT-Large Transformer XLM MT-DNN GPT-2 XLNET RoBERTa CTRL MegatronLM Grover-Mega T5-11B Meena Turing-NLG DistilBERT BST 9.4B GPT-3

Number of parameters

**Skoltech**

Training cost estimation: 10-50 MLN US$ for GPT-3

# Self-attention block complexity
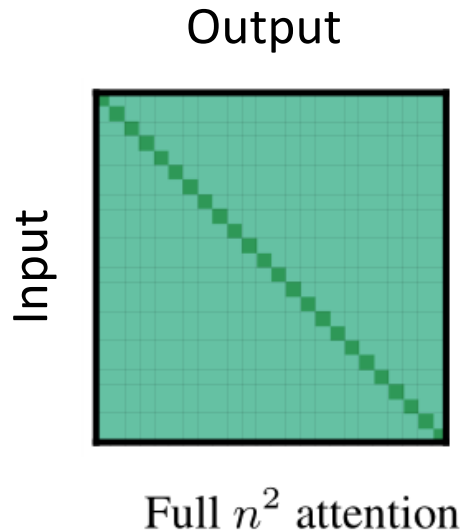
$$\text{Attention}(Q, K, V) = \text{softmax}(\frac{QK^T}{\sqrt{d_k}})V$$

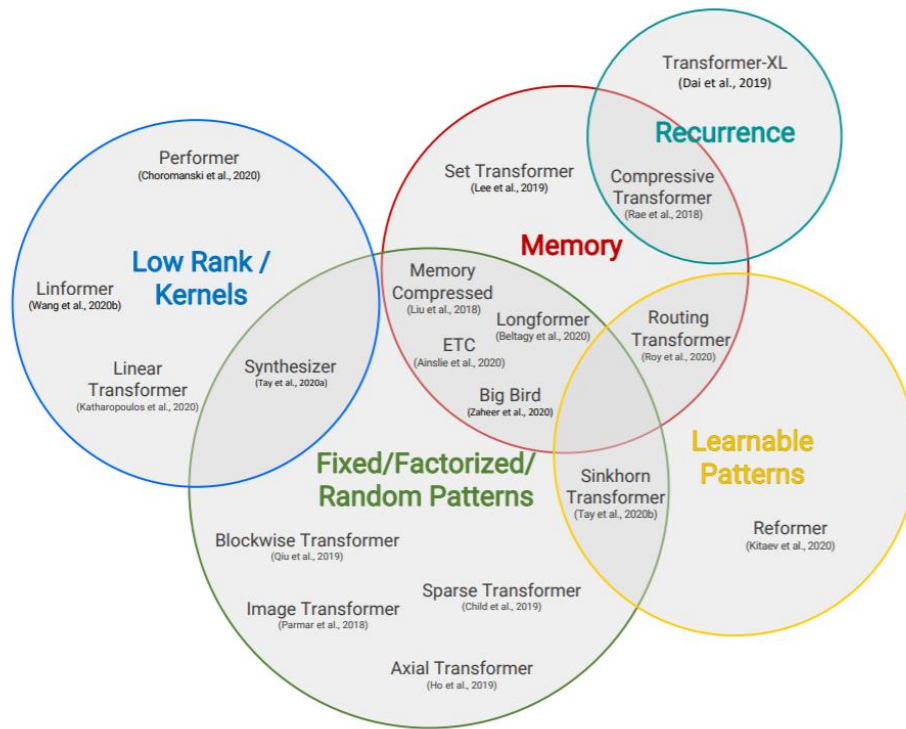Memory complexity is $O(d_x^2)$
Computational complexity is $O(d_x^2)$

Max sequence size in popular models (e.g. BERT) is only $n = d_x = 512$ tokens

In modern models tokens are parts of words

Output

Input

Full $n^2$ attention

# Transformer for long sequences

To work with sequences with significant length we should decrease memory consumption and computation complexity $O(n^2)$

Tay, Yi, et al. "Efficient transformers: A survey." *arXiv preprint arXiv:2009.06732* (2020).

Skoltech

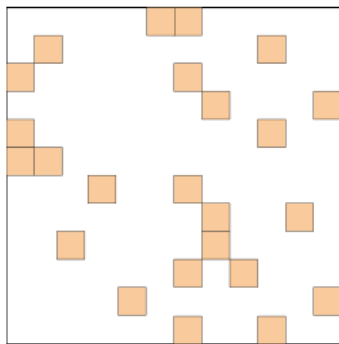# BigBird approach

Memory requirements, n is the sequence length:

Random attention requires $O(r \cdot n)$

Sliding window requires $O(h \cdot n)$,
h is the window size
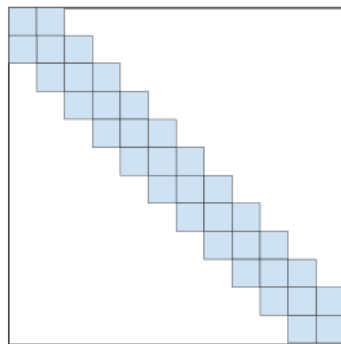
Global requires $O(g \cdot n)$, g is the global tokens number

BigBird combines <u>3 types of attention mechanism</u>. All of them have linear complexity.
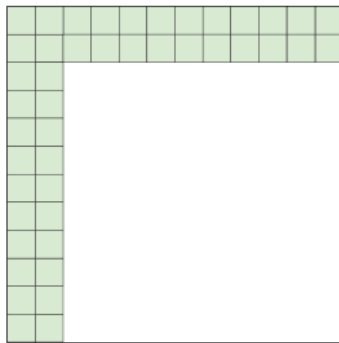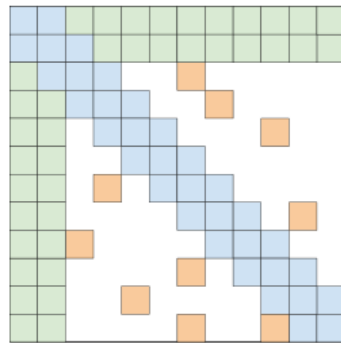
**Total: $\mathbf{O((r + h + \; g) \cdot n)}$**



(a) Random attention    (b) Window attention

(c) Global Attention    (d) BIGBIRD

Zaheer, Manzil, et al. "Big bird: Transformers for longer sequences." *Advances in Neural Information Processing Systems* 33 (2020): 17283-17297.

Skoltech

# Take-home messages

- Attention mechanism is for "information extraction"

- Transformer architecture is a multi-layer architecture based on the self attention layer