# Attention!

22th November, 2022

Alexey Zaytsev, Assistant professor

Models of Sequential Data

Skoltech

# 3 generations of Machine translation models



Skoltech

# Machine translation: a sequence to sequence problem

Translate a sentence from one language to another

source language                            target language

$$x \quad \longrightarrow \quad y$$

*A la guerre comme a la guerre*          *На войне как на войне*

# First generation: Early machine translation, 50s

**Cold war child: 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, uses a dictionary to map Russian words to English

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

# Second generation: Statistical machine translation, 90s-2010s

Learn probabilistic model from data

$$p(\boldsymbol{y}|\boldsymbol{x})$$

To translate: for an input English sequence $\boldsymbol{x}$ find the most probable Russian sentence $\boldsymbol{y}$

$$p(\boldsymbol{y}|\boldsymbol{x}) \rightarrow \max_{\boldsymbol{y}}$$

Bayesian perspective:

$$p(\boldsymbol{y}|\boldsymbol{x}) \sim p(\boldsymbol{x}|\boldsymbol{y})p(\boldsymbol{y})$$

Translation model: learn from parallel corpus
Language model: learn from monolingual corpus

Skoltech

# Learning translation model

Learn translation model from data

$$p(x|y)$$
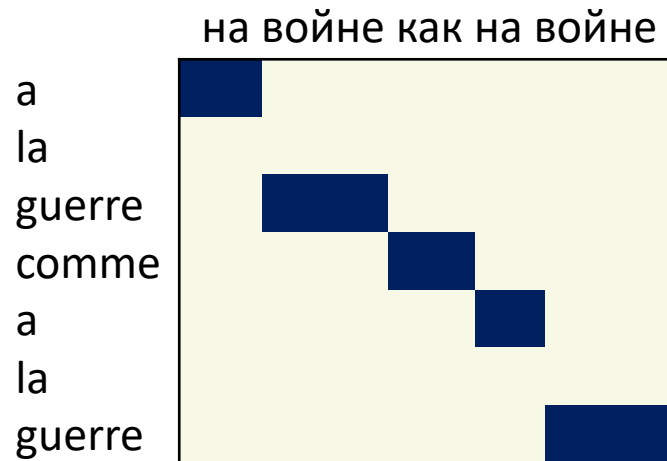
- Large amount of parallel data
- Alignment

$$p(x, a|y)$$

correspondence between words in different languages

Alignment types:
- one-to-one
- spurious words
- one-to-many
- Many-to-one

Many problems on the way

на войне как на войне

a
la
guerre
comme
a
la
guerre

Skoltech

# Decoding

The optimization problem is hard

$$p(\boldsymbol{y}|\boldsymbol{x}) \to \max_{\boldsymbol{y}}$$

- Full search is not possible
- Heuristic search algorithm to search for the best translation: look through a tree of possible options

# The best SMT systems are very complex

- Language itself is very complex
- Many details we don't even mention
- Separately designed subcomponents
- Tricky feature engineering
- Extra information
- The language changes – we need to maintain the system

All difficulties we saw about pre-Neural
approaches to sequence processing
multiplied x100!

Skoltech

Third generation: Neural Machine Translation!

# Neural machine translation

Our goal is to do machine translation with one Neural network
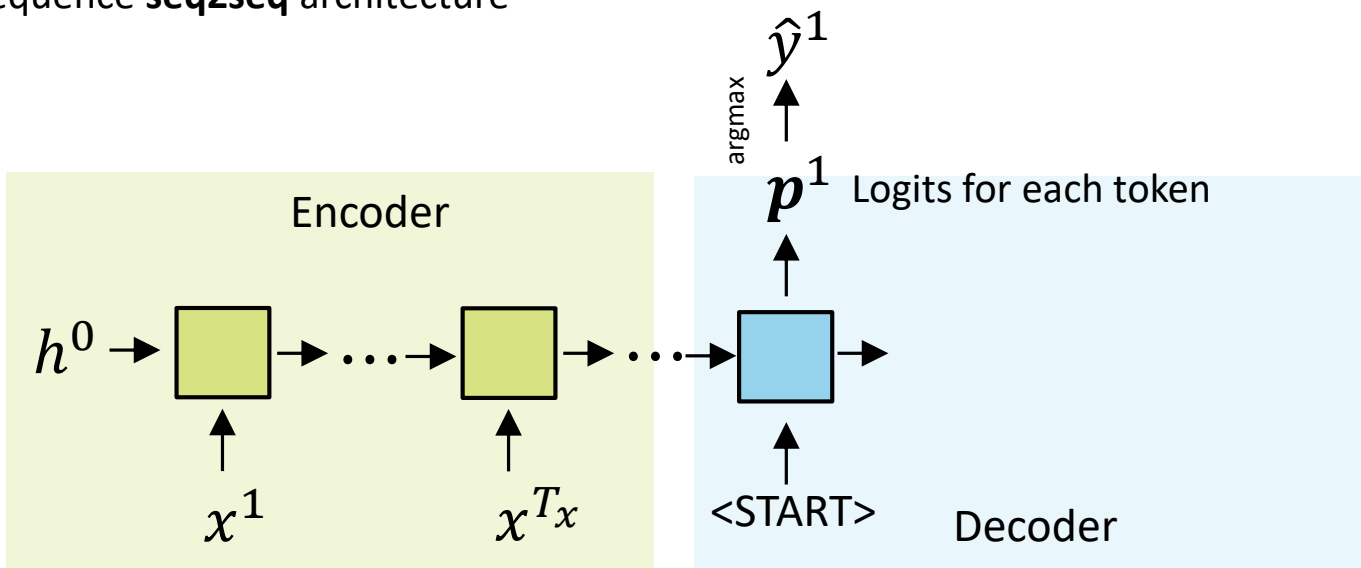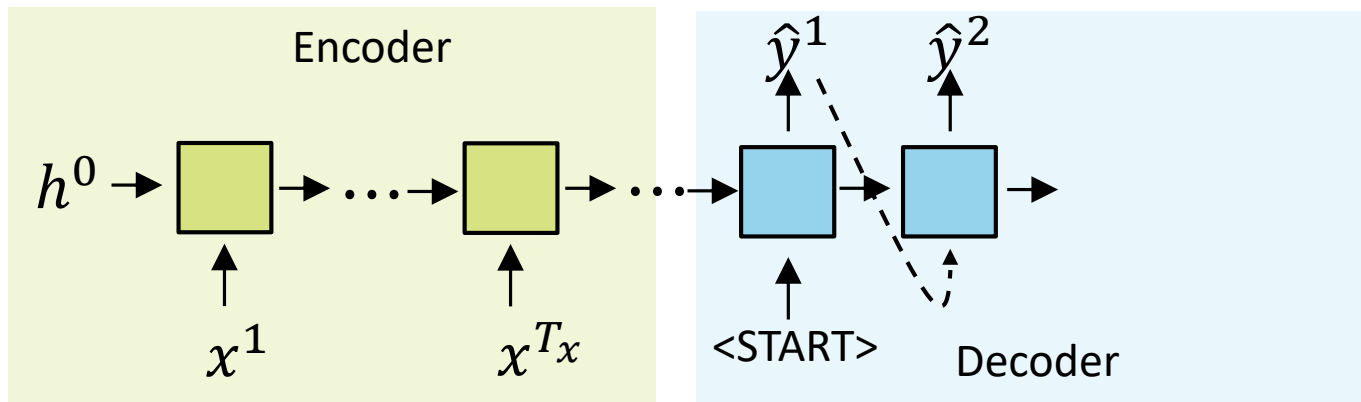
It works with two RNNs

Sequence to sequence **seq2seq** architecture

Skoltech

# Neural machine translation

Our goal is to do machine translation with one Neural network

It works with two RNNs

Sequence to sequence **seq2seq** architecture

$$\hat{y}^1$$

argmax

$$\boldsymbol{p}^1 \quad \text{Logits for each token}$$

Encoder

$$h^0 \rightarrow$$

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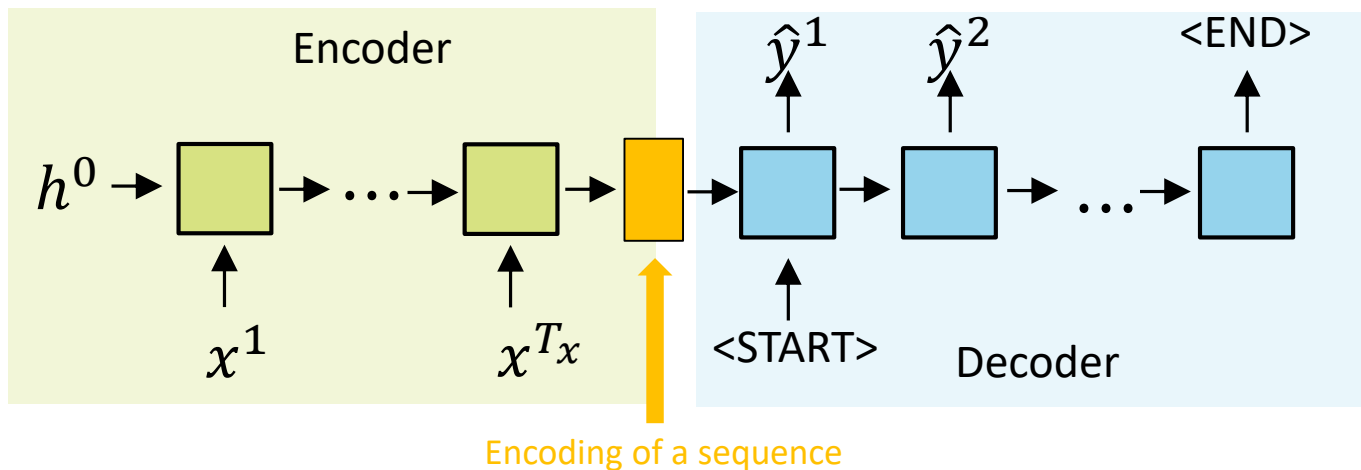
<START>

Decoder

Skoltech

# Neural machine translation

Our goal is to do machine translation with one Neural network

It works with two RNNs

Sequence to sequence **seq2seq** architecture
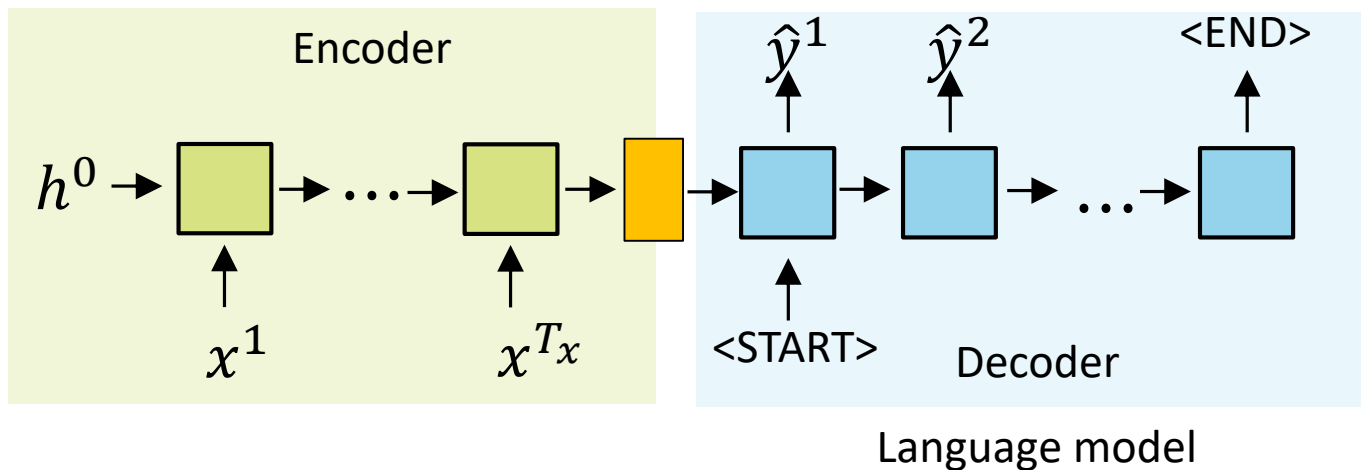
# Neural machine translation

Our goal is to do machine translation with one Neural network

It works with two RNNs

Sequence to sequence **seq2seq** architecture



Encoder

Decoder

Encoding of a sequence

$h^0$  $x^1$  $x^{T_x}$  $\hat{y}^1$  $\hat{y}^2$  <END>  <START>

# Neural machine translation

Our goal is to do machine translation with one Neural network

It works with two RNNs

Sequence to sequence **seq2seq** architecture
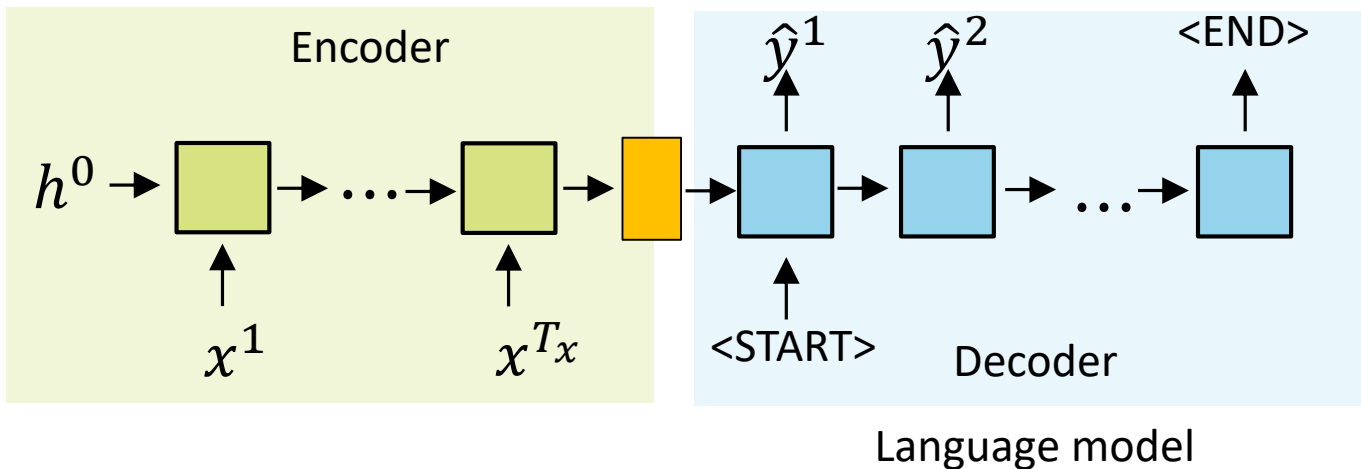


Skoltech

# New world of seq2seq models

Summarization: long text – text summary
Dialogue: one phrase – another phrase
Parsing: input – output parse as a sequence
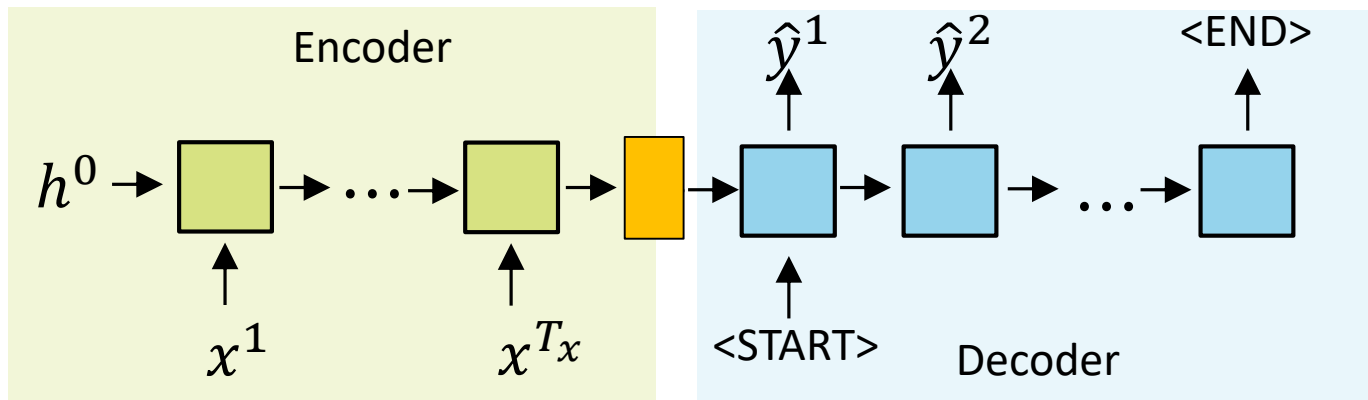Code generation: task description – python code



Encoder

Decoder

$\hat{y}^1$   $\hat{y}^2$   <END>

$h^0 \rightarrow$

$x^1$   $x^{T_x}$   <START>

Language model

Skoltech

# NMT directly models the conditional language model

Learn probabilistic model from data

$$p(\boldsymbol{y}|\boldsymbol{x}) = p(y_T|y_1, y_2, \ldots, y_{T-1}, \boldsymbol{x})\, p(y_{T-1}|y_1, y_2, \ldots, y_{T-2}, \boldsymbol{x})\ldots p(y_1|\boldsymbol{x})$$

Each term is an RNN block

Loss function: compare logits to true words, now we have a logloss
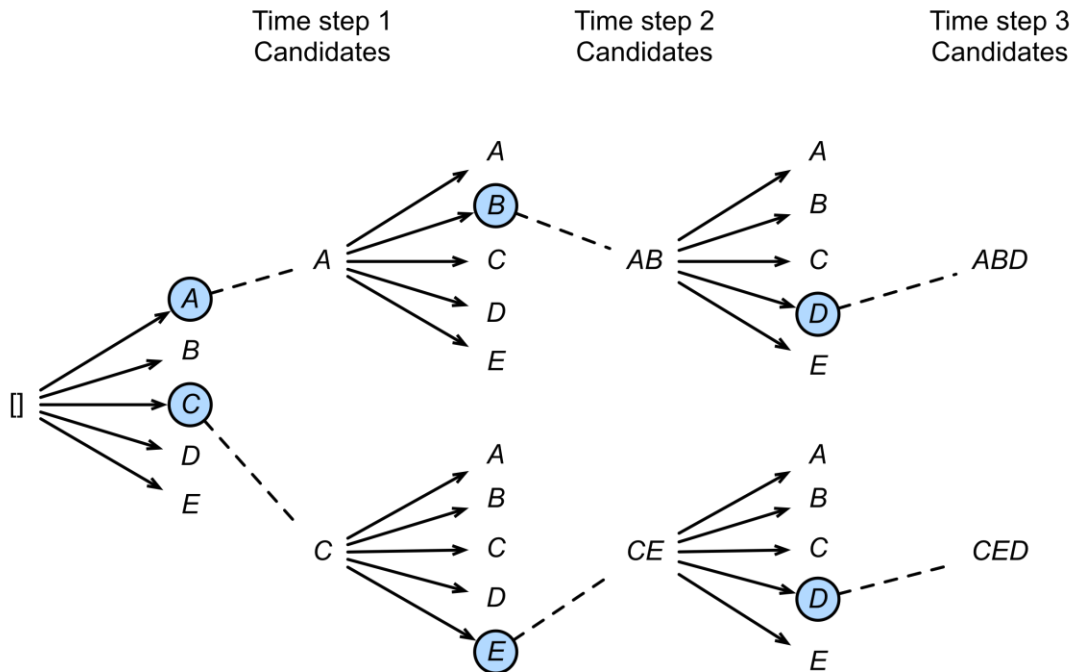if we have a parallel corpus

# Beam search decoding

**Problem:**

We generate tokens one by one from the very beginning Wrong decision at some step leads to wrong translation as a whole

**Solution:**

"Beam search": Keep a population of solutions and select the most probable sequences at each step

Skoltech

Time step 1
Candidates

Time step 2
Candidates

Time step 3
Candidates
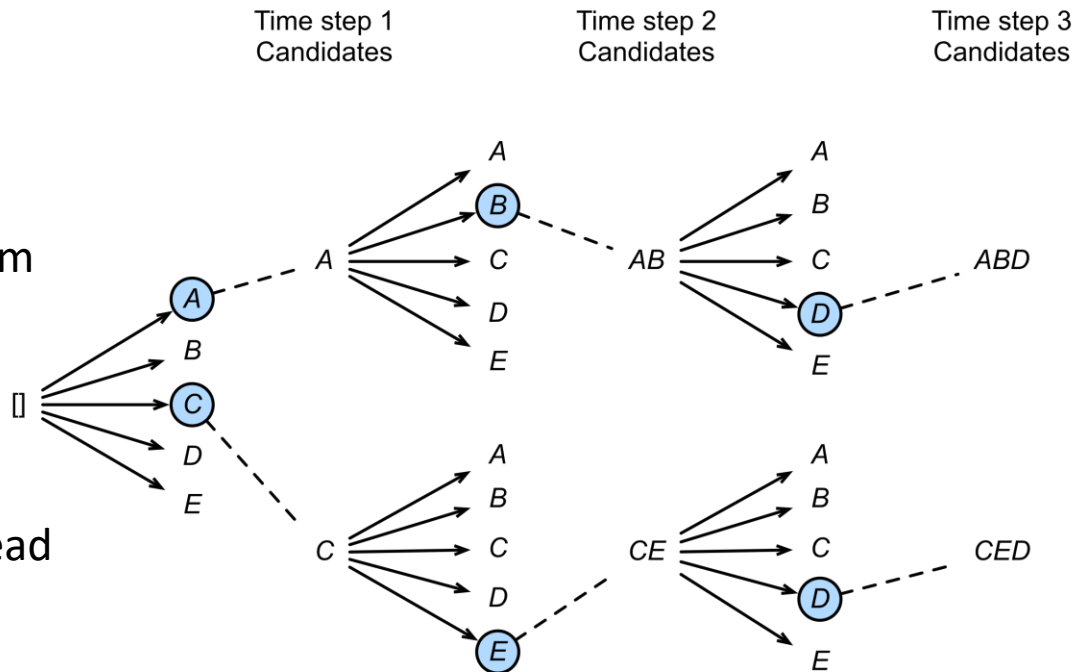
# Beam search decoding: practical implementation

**Stopping criteria:**
- Generated <END> token
- Reached maximum sequence length

We have $\{y_1, y_2, \ldots, y_m\}$ after beam search

**Selection criterion:**

Compare normalized $\tilde{p}(y_i|x)$ instead of $p(y_i|x)$

$$\tilde{p}(y|x) = \frac{1}{|y|} p(y|x)$$



Time step 1
Candidates

Time step 2
Candidates

Time step 3
Candidates

Skoltech

# Advantages of NMT

- Better performance compared to SLT (Statistical Learning Translation)

- Single end2end neural network

- Much less human engineering effort
  - No feature engineering
  - Same method for all language pairs

Skoltech

# Evaluations of Machine Translation: BLEU

"the closer a machine translation is to a professional human translation, the better it is"

**Precision**: share of words from $\hat{y}$ that appear in $y$

| True $y$ | the | cat | is | on | the | mat |
|---|---|---|---|---|---|---|
| Candidate $\hat{y}$ | the | the | the | the | the | the |

Unigram precision is 6/6 = 1

Limit with number of words from references: there are 2 "the" in $y$

BLEU modified unigram precision is 2/6 = 0.33

Skoltech

Papineni, K.; Roukos, S.; Ward, T.; Zhu, W. J. 2002. *BLEU: a method for automatic evaluation of machine translation*. ACL-2002: 40th Annual meeting of the Association for Computational Linguistics

# BLEU evaluation examples

| True $y$ | the | cat | is | on | the | mat |
|----------|-----|-----|-----|-----|-----|-----|
| Candidate $\hat{y}$ | the | the | cat | cat | cat | the |

Unigram precision is 6/6 = 1
Bigram precision is 1/5 = 0.2
BLEU unigram score 2/6
BLEU bigram score is 1/5

Good BLEU scores is about 0.7
Higher scores mean overfitting or other problems

Papineni, K.; Roukos, S.; Ward, T.; Zhu, W. J. 2002. *BLEU: a method for automatic evaluation of machine translation.* ACL-2002: 40th Annual meeting of the Association for Computational Linguistics
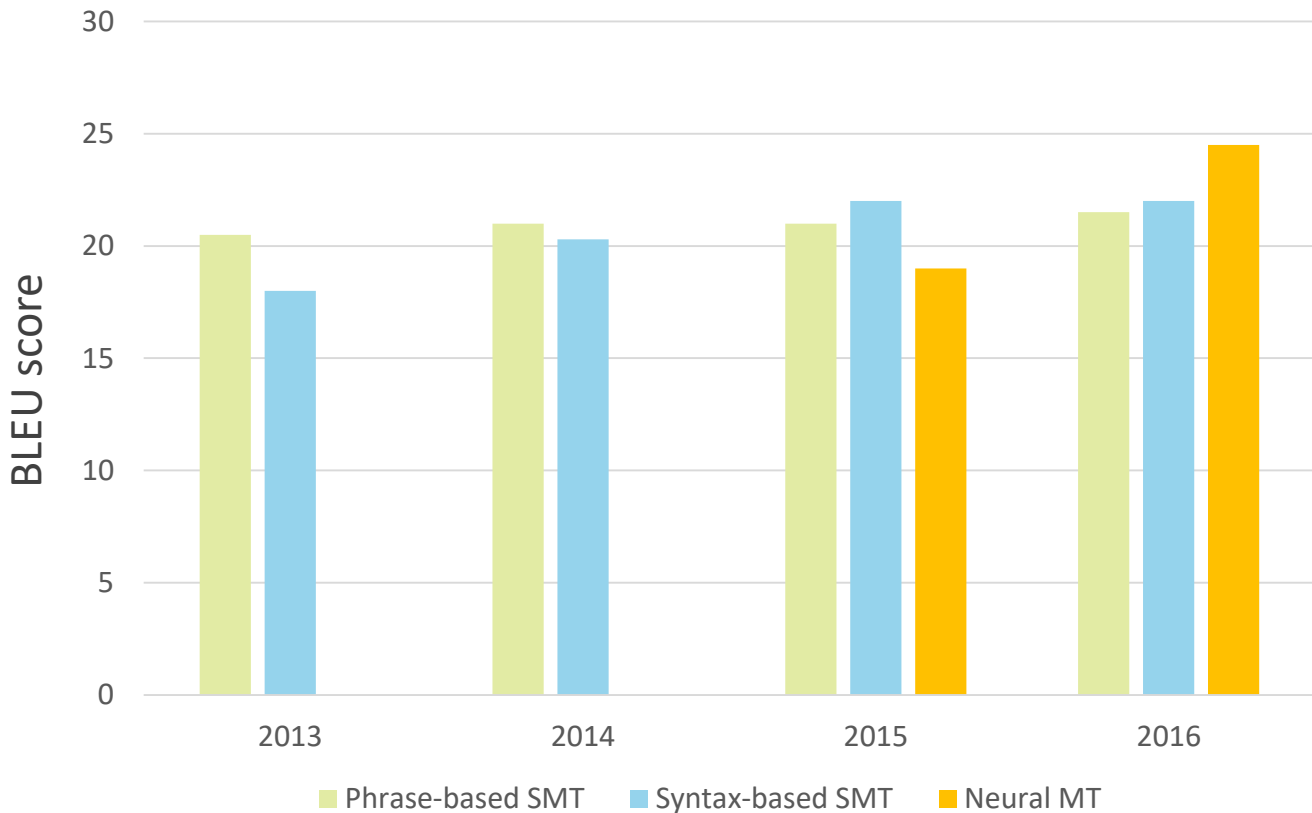
Skoltech

# General advices on BLEU

In practice we use **n-gram modified precision** with n = 4 that "best correlates to human judgement".

Good BLEU scores is about 0.7

Higher scores mean overfitting or other problems
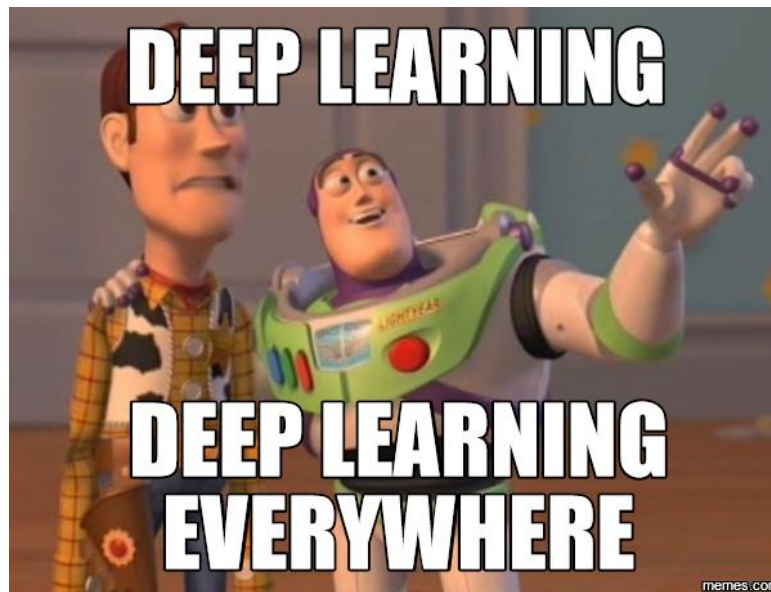
**Recall** is also important

**Skoltech**

Papineni, K.; Roukos, S.; Ward, T.; Zhu, W. J. 2002. *BLEU: a method for automatic evaluation of machine translation*. ACL-2002: 40th Annual meeting of the Association for Computational Linguistics

# Evaluations of Machine Translation

Skoltech

# NMT: success story for deep learning

**2014:** first seq2seq paper published
**2016:** Google translate switches from SMT to NMT

**SMT:** hundreds of engineers *for many years*
**NMT:** handful of engineers *in a few months*
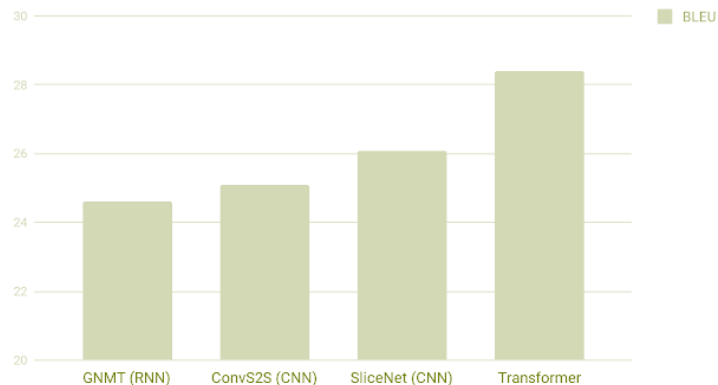
# NMT problems

- Out-of-vocabulary words
- Domain mismatch: training and test data
- Context over longer texts
- Low-resource language pairs

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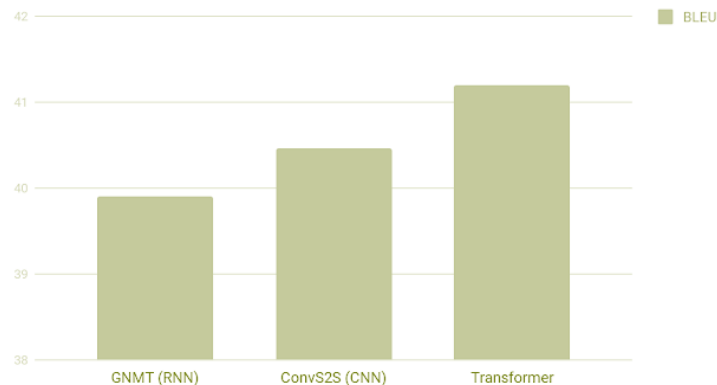
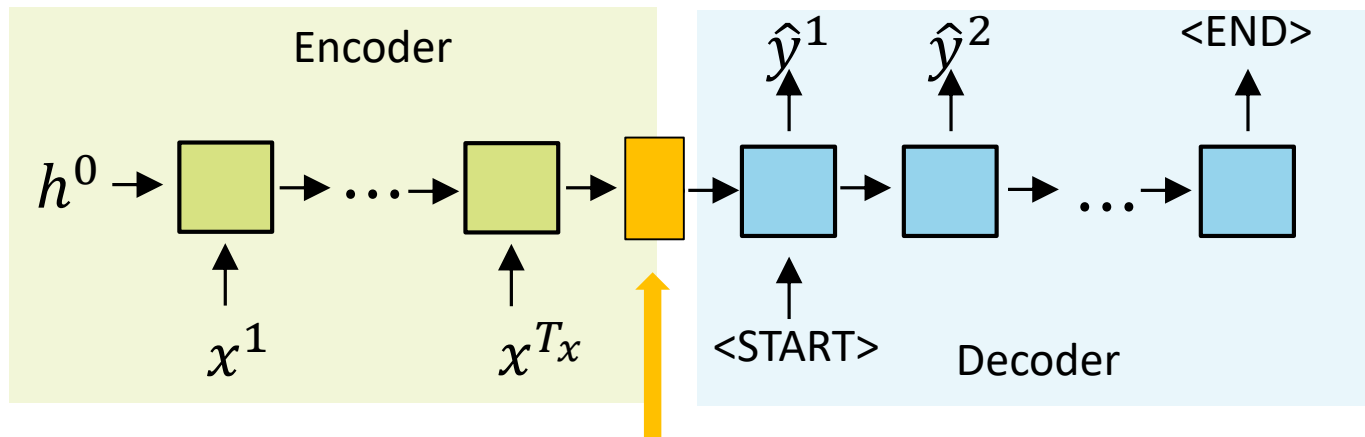# New state of the art: **attention** is all we need



attention



English German Translation quality



English French Translation Quality

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

# Bottleneck in seq2seq models



Encoder

Decoder

$h^0 \rightarrow$

$x^1$

$x^{T_x}$

$\hat{y}^1$

$\hat{y}^2$

<END>

<START>

All information about the sequence is in this vector

Skoltech

# Attention

- Solution to the bottleneck problem

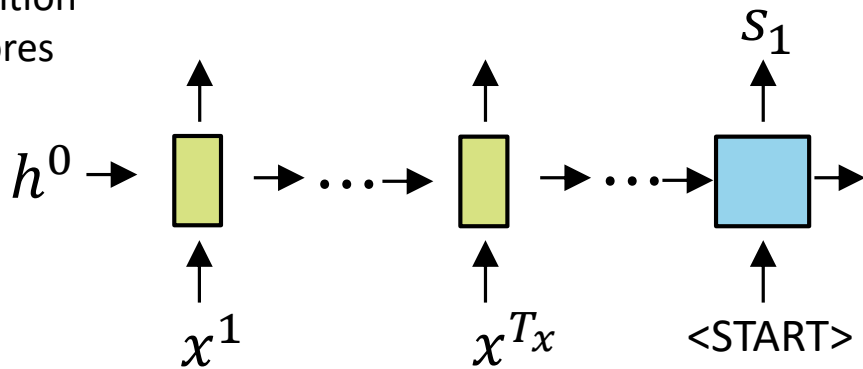- Direction connection between parts of input and output sequence

# Sequence 2 sequence with attention

Attention
output
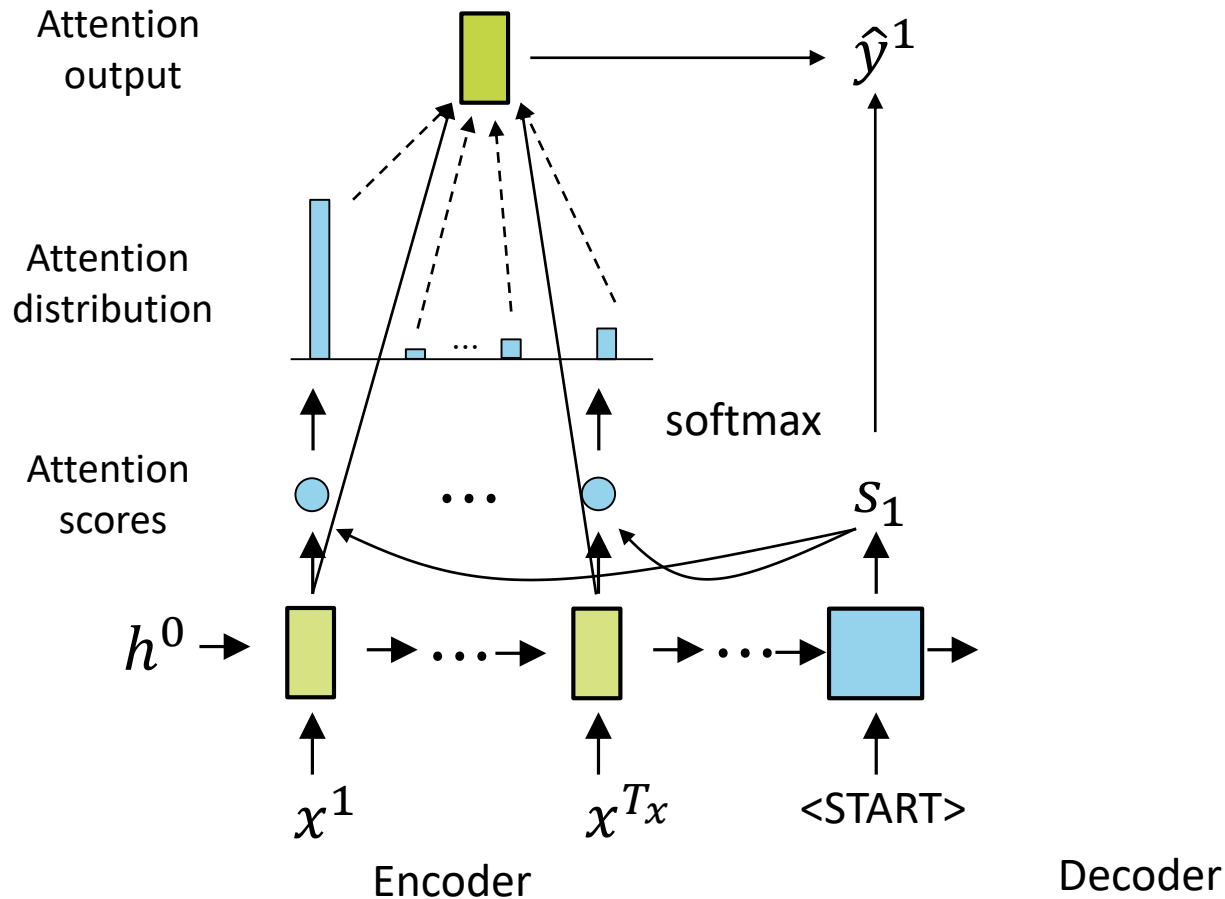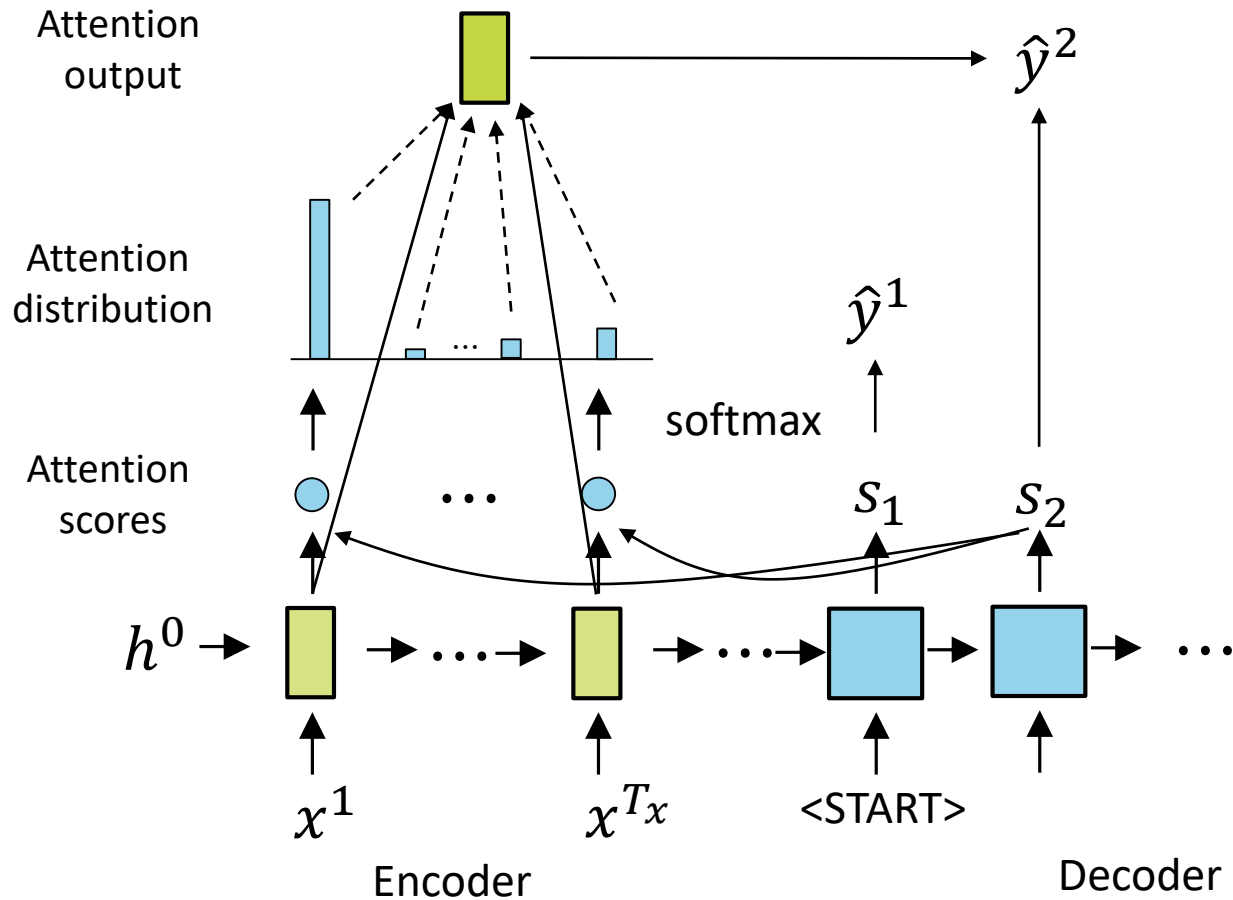
Attention
distribution

Attention
scores

$h^0$

$s_1$

$x^1$

$x^{T_x}$

<START>

Encoder

Decoder

Skoltech

# Sequence 2 sequence with attention

Attention output

Attention distribution

$\hat{y}^1$

softmax

Attention scores

$s_1$

$h^0 \rightarrow$

$x^1$

$x^{T_x}$

<START>

Encoder

Decoder

Skoltech

# Sequence 2 sequence with attention



Attention output

Attention distribution

Attention scores

$\hat{y}^2$

$\hat{y}^1$

softmax

$s_1$ $s_2$

$h^0$

$x^1$ $x^{T_x}$ <START>

Encoder Decoder

Skoltech

# Attention: formulas

- First RNN produces encoder hidden states $\boldsymbol{h}_1, \ldots, \boldsymbol{h}_{T_x} \in \mathbb{R}^h$
- Decoder hidden state $\boldsymbol{s}_t \in \mathbb{R}^h$ at time step $t$
- Attention scores for step $t$:

$$\boldsymbol{e^t} = [\boldsymbol{s}_t^T \boldsymbol{h}_1, \ldots, \boldsymbol{s}_t^T \boldsymbol{h}_{T_x}] \in \mathbb{R}^{T_x}$$

- 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:

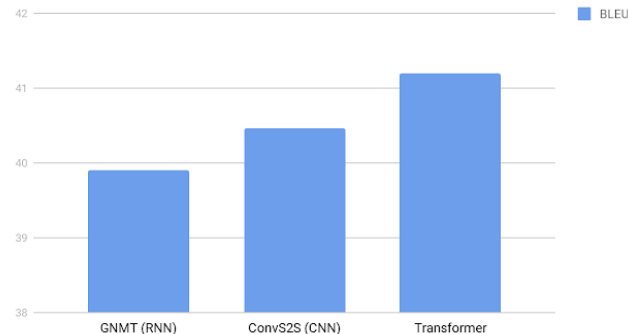$$\boldsymbol{a}_t = \sum_{i=1}^{T_x} \alpha_i^t \boldsymbol{h}_i \in \mathbb{R}^h$$

- 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

$$[\boldsymbol{a}_t, \boldsymbol{s}_t] \in \mathbb{R}^{2h}$$
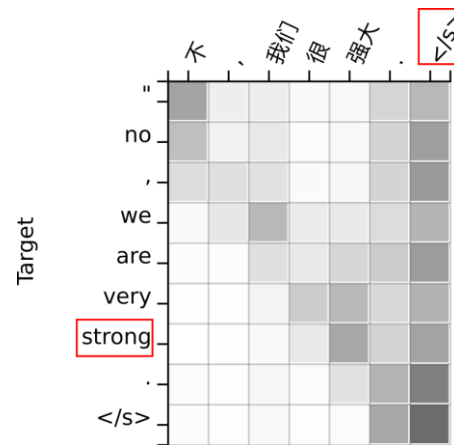
Skoltech

# 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.



Skoltech

# 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 $\sum_{i=1}^{n} \alpha_i \boldsymbol{h}_i$

# Transformers

# Key value interpretation

$q_i$ - query to a database        Hidden state of the *decoder*

$k_j$ - keys in the database        Hidden state of the *encoder*

$v_j$ - values in the database        Hidden state of the *encoder*

We **calculate** attention scores

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

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

Then we extract the information as weighted sum of values

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

Skoltech

# Matrix key value interpretation

$q_i$ - query to a database          Hidden state of the *decoder*

$k_j$ - keys in the database          Hidden state of the *encoder*

$v_j$ - values in the database          Hidden state of the *encoder*

We calculate correspondences

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

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

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

New scaled formula:
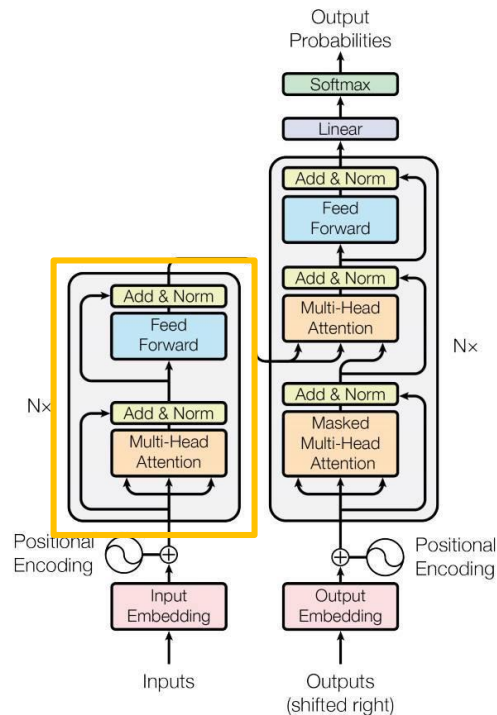$$A(Q, K, V) = \text{softmax}(QK^T/\sqrt{d_k})V$$

Skoltech

# 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:
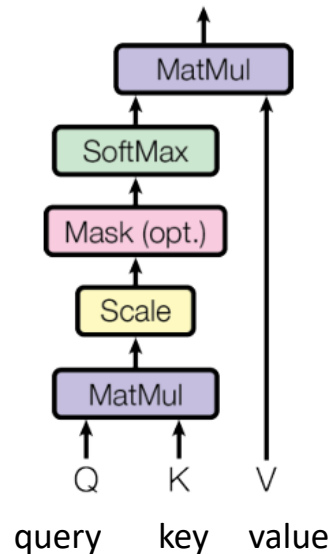- Multi-head attention
- Feed Forward



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



**Skoltech**

# Self-attention block

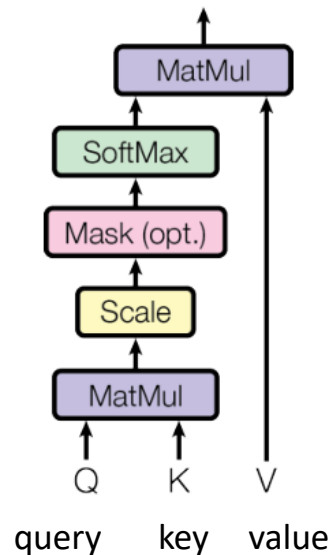$$\text{Attention}(Q, K, V) = \text{softmax}(\frac{QK^T}{\sqrt{d_k}})V$$

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

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

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

$$V = XW^V, \dim(W^Q) = d_x \times d_v,$$
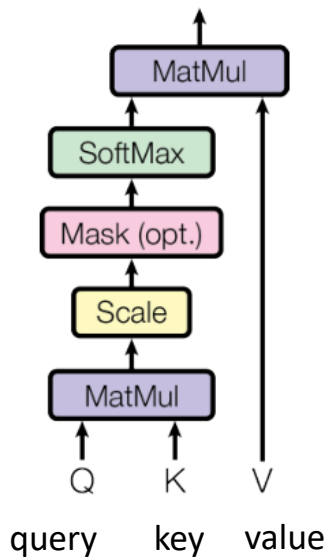
Scaled Dot-Product Attention



query    key    value

# Multi-Head attention

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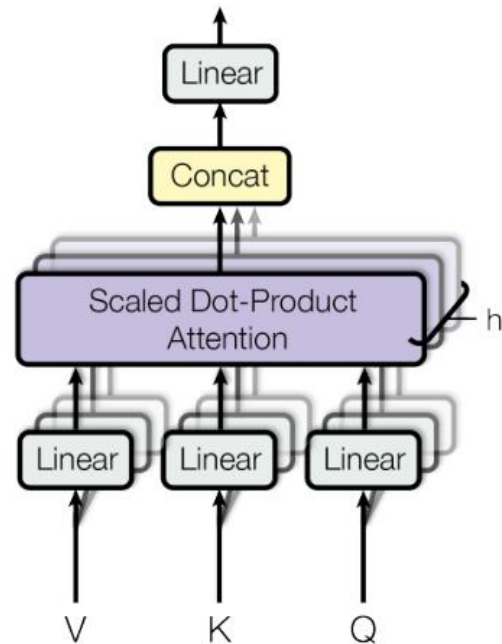
$$\text{where head}_i = \text{Attention}(QW_i^Q, KW_i^K, VW_i^V)$$

h heads in total



Scaled Dot-Product Attention

query    key    value

Multi-Head Attention

"Attention is all you need" paper

# Full block

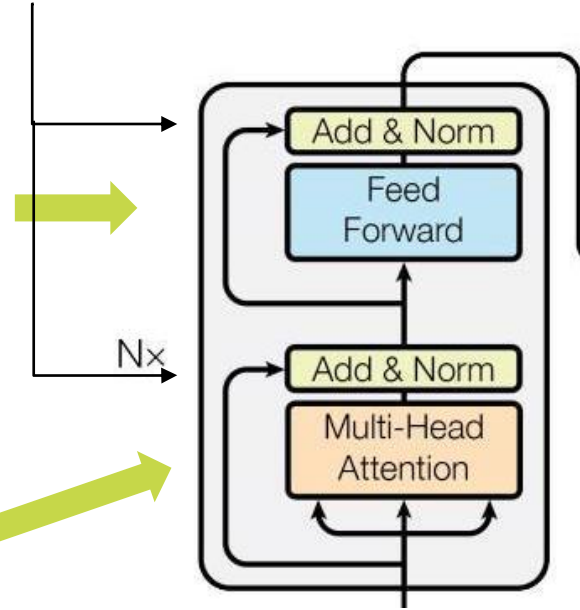Two linear transformation with ReLU activation in between

$$\text{FFN}(x) = \max(0, xW_1 + b_1)W_2 + b_2$$

Multi-head attention

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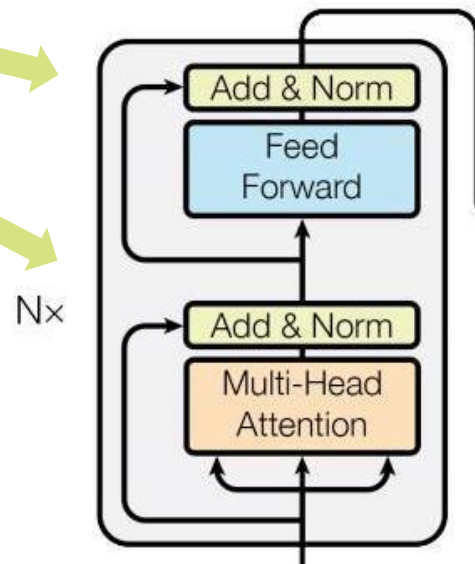


LaverNorm$(x + \text{Sublayer}(x))$

Add & Norm
Feed Forward
N×
Add & Norm
Multi-Head Attention

**Skoltech**

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

"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

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

Skoltech

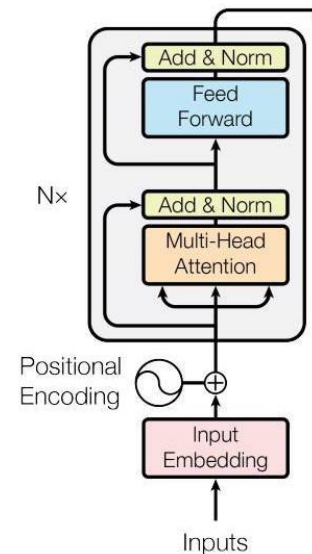Layer Normalization, https://arxiv.org/abs/1607.06450

# 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

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

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



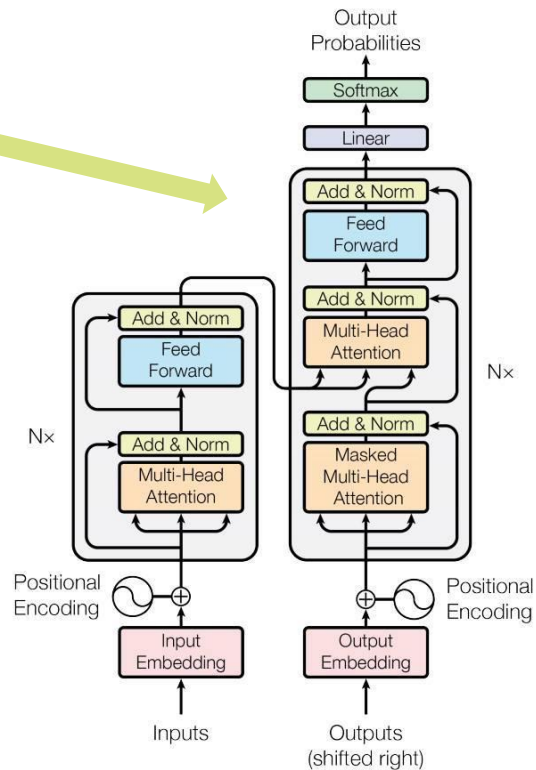| POSITIONAL ENCODING | 0 | 0 | 1 | 1 | | 0.84 | 0.0001 | 0.54 | 1 | | 0.91 | 0.0002 | -0.42 | 1 |

EMBEDDINGS $x_1$    $x_2$    $x_3$

INPUT    Je    suis    étudiant

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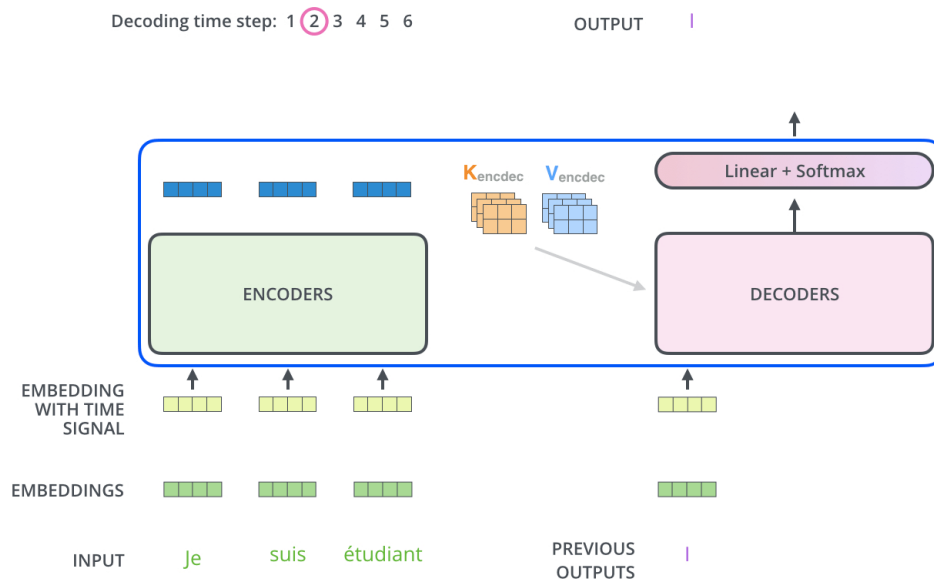
Skoltech

# Decoder is similar with another N=6 blocks

Additional sublayer to take into account attentions from encoder

We generate one token and proceed to the next token generation

"Attention is all you need" paper

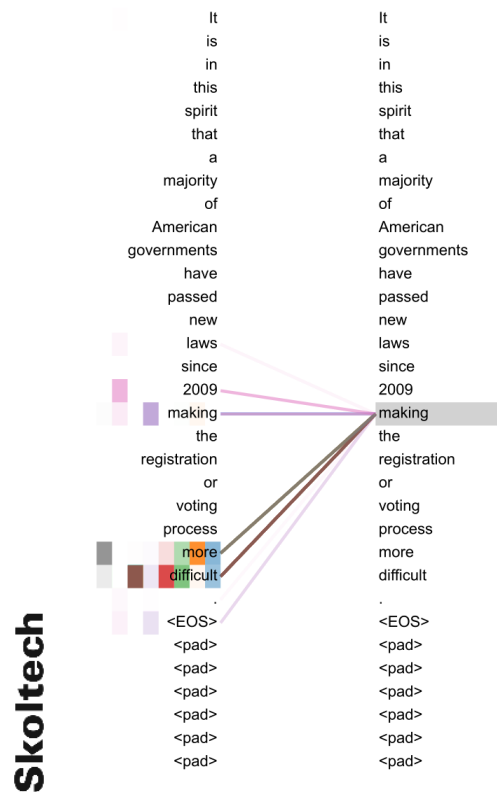# Decoder inference

We generate one token and proceed to the
next token generation

# Attention visualizations

Skoltech

# Attention visualizations
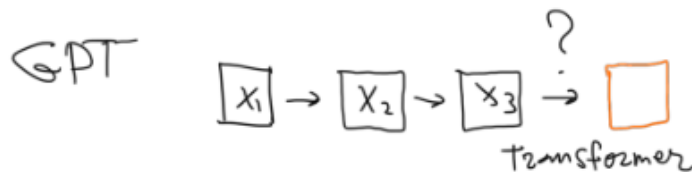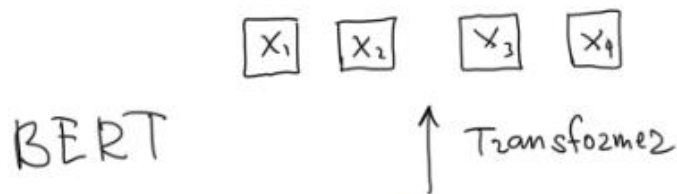
Skoltech

# Transformer training

- Masked language model:

Replace random tokens with masks, try to reconstruct them using a Neural network

- Next token prediction

Predict next token

We don't need labeled examples, we just create them
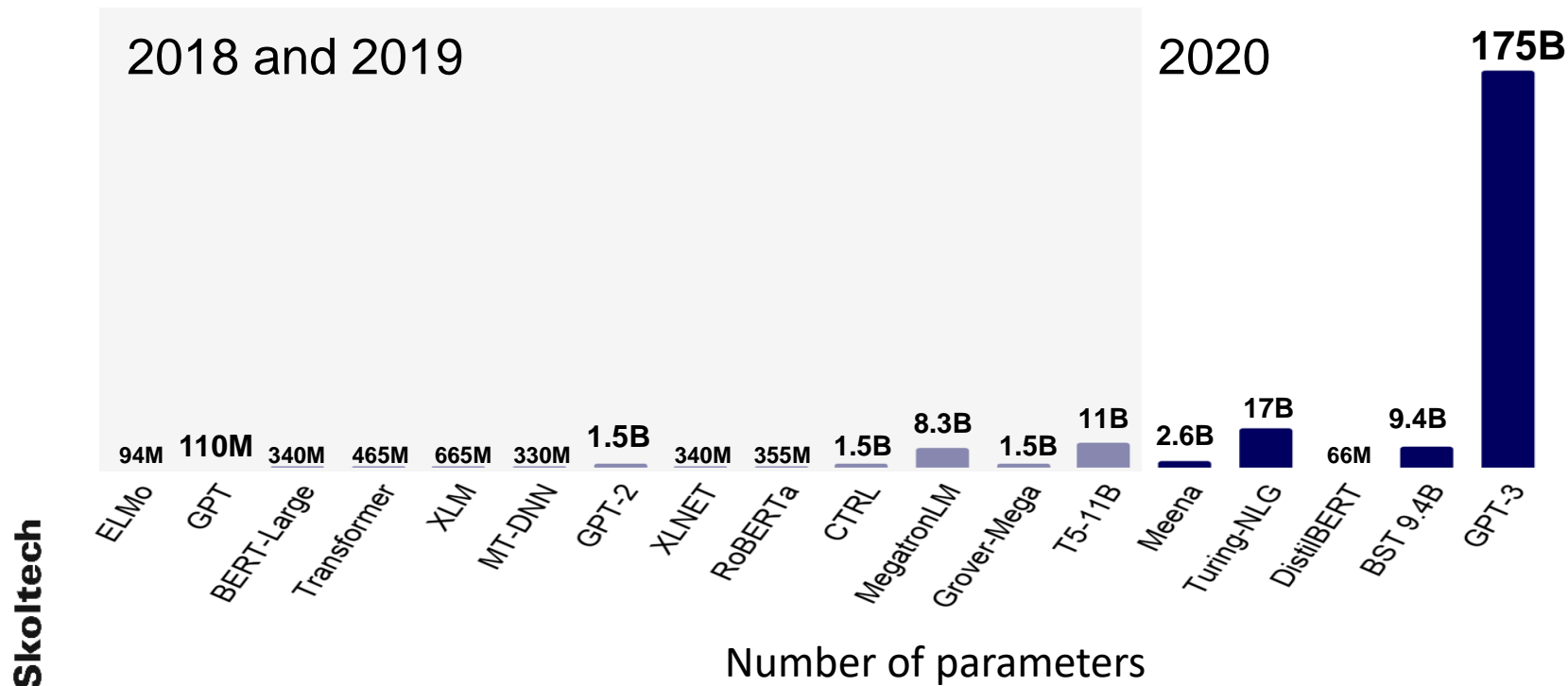
# Efficiency of transformers

| Model Name | $n_{\text{params}}$ | $n_{\text{layers}}$ | $d_{\text{model}}$ | $n_{\text{heads}}$ | $d_{\text{head}}$ | 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 (~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

| Model | Parameters |
|-------|-----------|
| ELMo | 94M |
| GPT | 110M |
| BERT-Large | 340M |
| Transformer | 465M |
| XLM | 665M |
| MT-DNN | 330M |
| GPT-2 | 1.5B |
| XLNET | 340M |
| RoBERTa | 355M |
| CTRL | 1.5B |
| MegatronLM | 8.3B |
| Grover-Mega | 1.5B |
| T5-11B | 11B |
| Meena | 2.6B |
| Turing-NLG | 17B |
| DistilBERT | 66M |
| BST 9.4B | 9.4B |
| GPT-3 | 175B |

Number of parameters

**Skoltech**

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

# Conclusions

- Attention + positional encoding are all you need

- The attention idea is pretty general

**Skoltech**

# Sources

- Stanford CS224N: NLP with Deep Learning | Winter 2019 | Lecture 8 – Translation, Seq2Seq, Attention

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

**Skoltech**