

# Attention becomes transformer

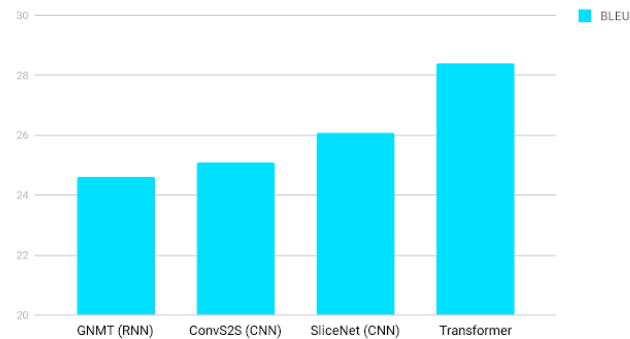
8 May, 2020

Alexey Zaytsev, head of Laboratory LARSS, PhD

Foundations of Data Science

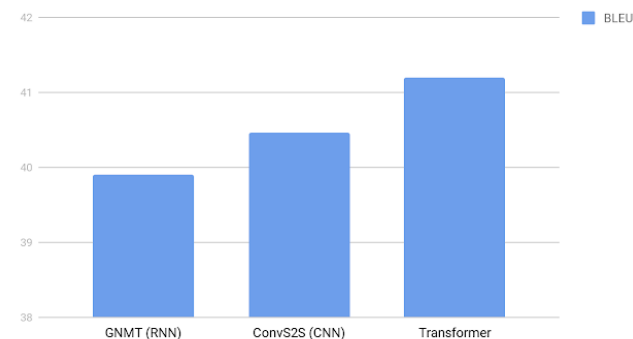
# New state of the art: **attention** is all we need

English German Translation quality



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

English French Translation Quality



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

# Machine translation

Translate a sentence from one language to another

source language

target language

$x$



$y$

*A la guerre comme a la guerre*

*На войне как на войне*

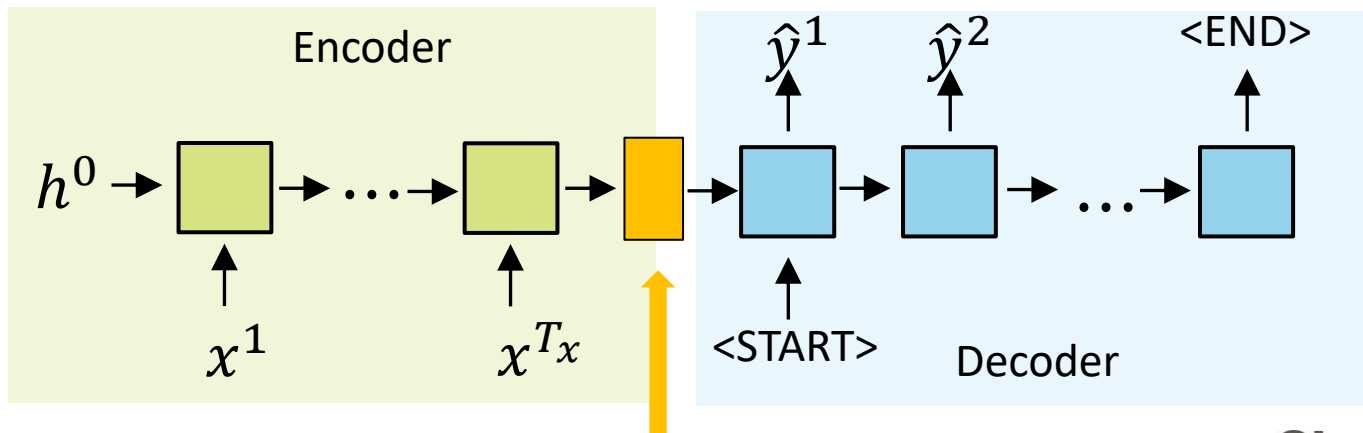
# New world of seq2seq models

Summarization: long text – text summary

Dialogue: one phrase – another phrase

Parsing: input – output parse as a sequence

Code generation: task description – python code

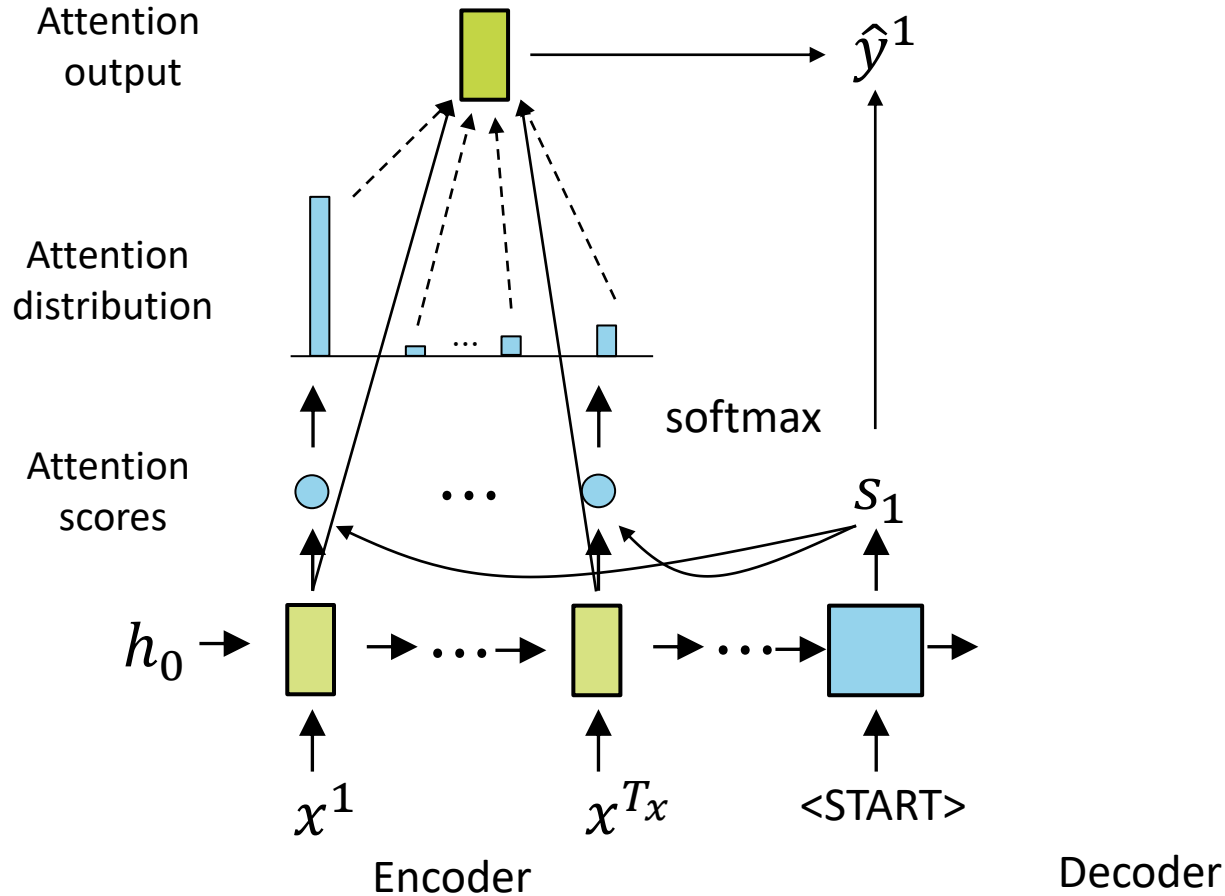


All information about the sequence is in this vector

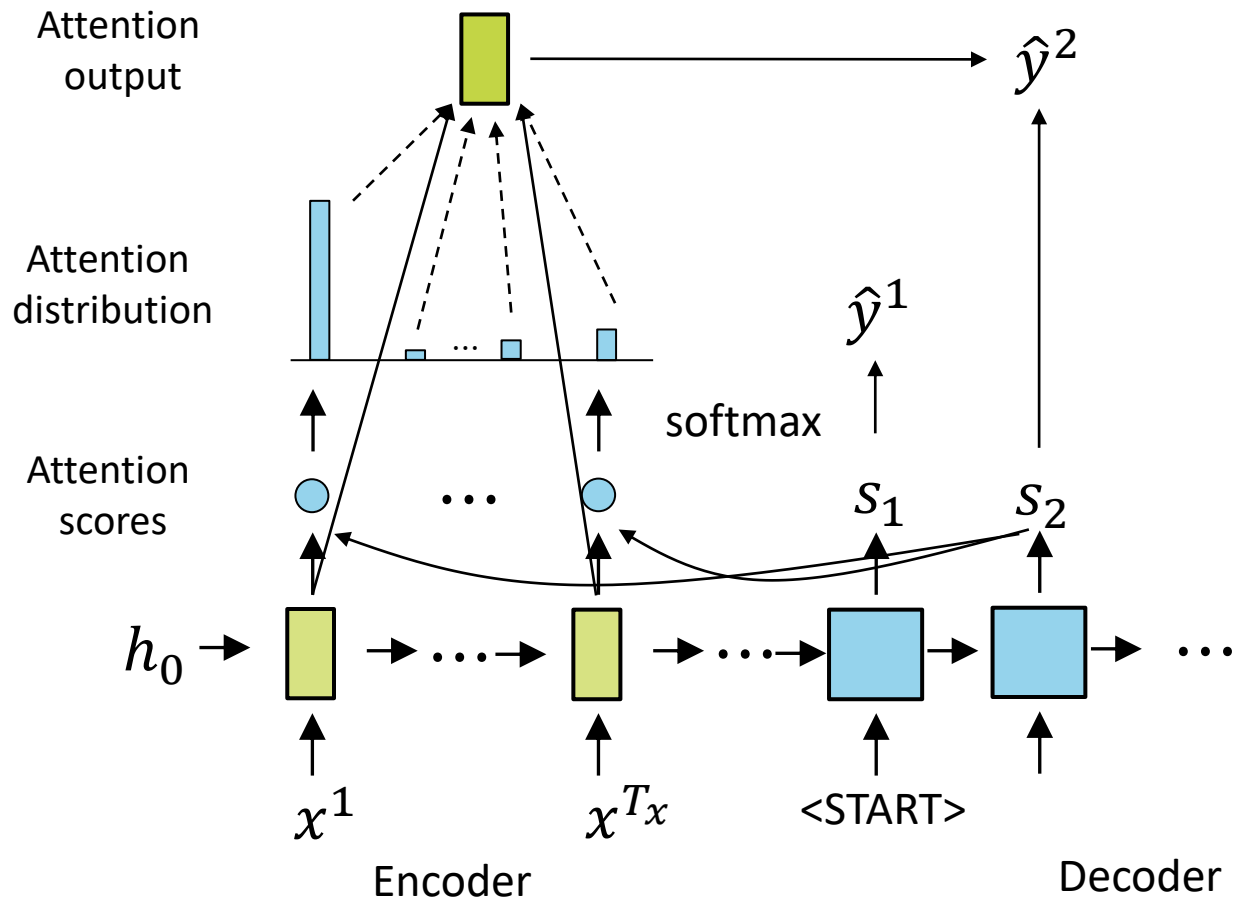
# Attention

- Solution to the bottleneck problem
- Direction connection between parts of input and output sequence

## Sequence 2 sequence with attention



## Sequence 2 sequence with attention



## Attention: formulas

- First RNN produces encoder hidden states  $\mathbf{h}_1, \dots, \mathbf{h}_{T_x} \in \mathbb{R}^h$
- Decoder hidden state  $\mathbf{s}_t \in \mathbb{R}^h$  at time step  $t$
- Attention scores for step  $t$ :  
 $\mathbf{e}^t = [\mathbf{s}_t^T \mathbf{h}_1, \dots, \mathbf{s}_t^T \mathbf{h}_{T_x}] \in \mathbb{R}^{T_x}$
- Softmax to get attention distribution: all values are positive, sum of all values is 1:

$$\boldsymbol{\alpha}^t = \text{softmax}(\mathbf{e}^t) \in \mathbb{R}^{T_x}$$

- Attention output  $\mathbf{a}_t$  is the weighted sum of hidden states:

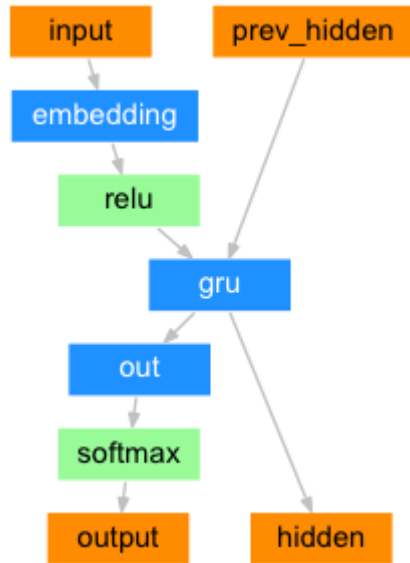
$$\mathbf{a}_t = \sum_{i=1}^{T_x} \alpha_i^t \mathbf{h}_i \in \mathbb{R}^h$$

- We concatenate the attention output  $\mathbf{a}_t$  with the decoder hidden state  $\mathbf{s}_t$  and proceed to the non-attention part of our seq2seq model

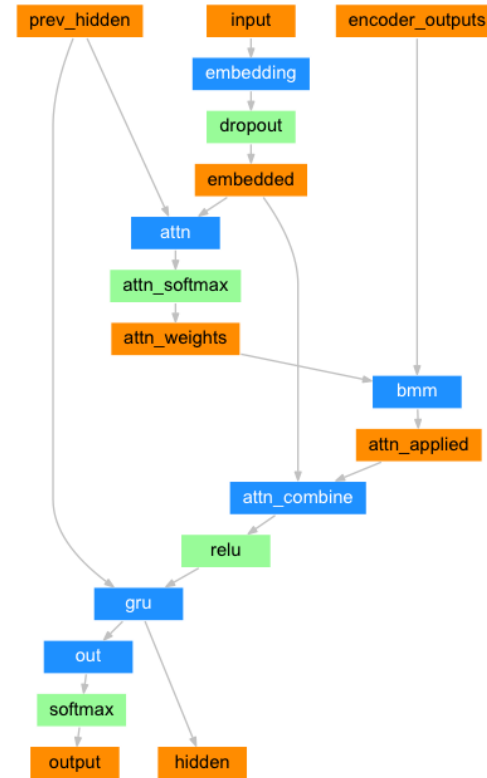
$$[\mathbf{a}_t, \mathbf{s}_t] \in \mathbb{R}^{2h}$$



# Attention: blocks



Simple decoder



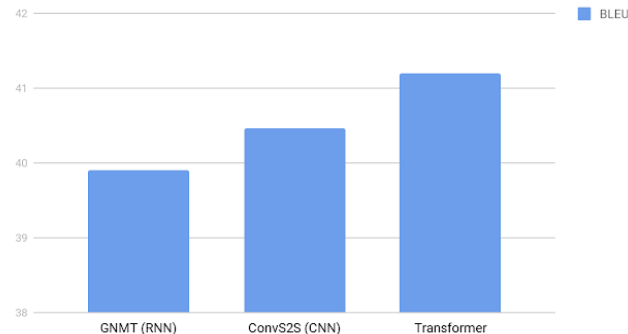
Decoder with attention

# Attention is just great

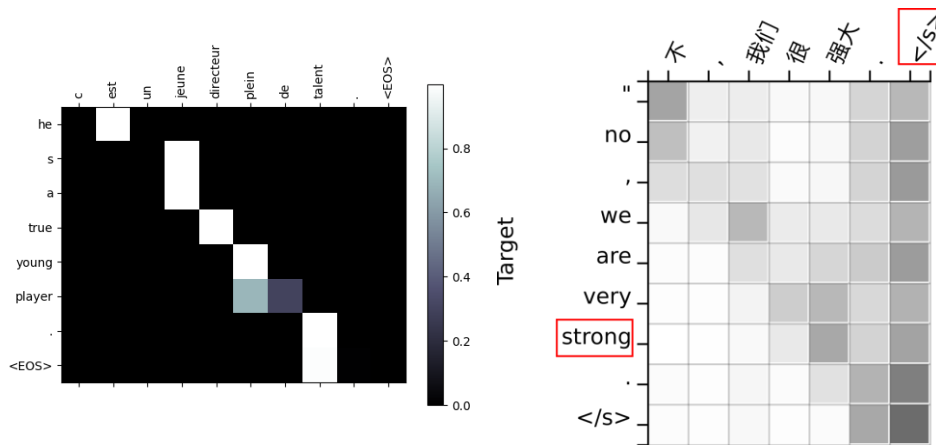
Similar to RNN seq2seq, but greater!

- 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

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
- Graph Neural networks

## Key value interpretation:

$s_i$  - query to a database

Hidden state of the decoder

$k_j$  - keys in the database

Hidden state of the encoder

$h_j$  - values in the database

Hidden state of the encoder

We calculate correspondences  $e(s_i, k_i)$

Then we extract the attention as weighted sum of values  $\mathbf{a}_i = \sum_{j=1}^{T_x} \alpha_j \mathbf{h}_j$

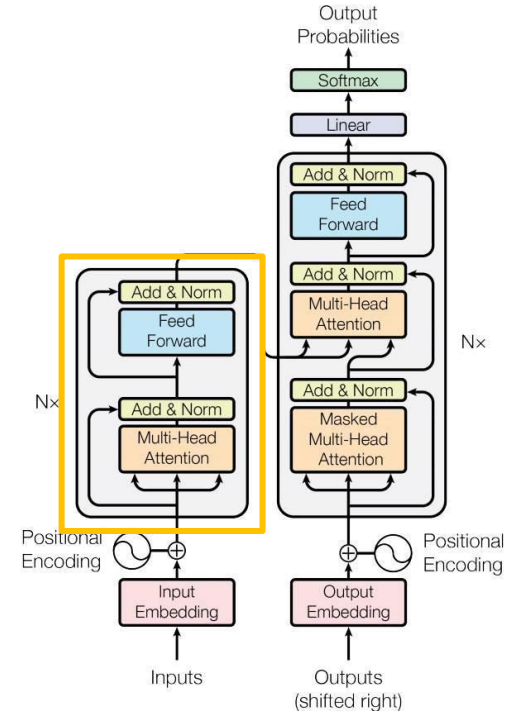
Transformer is based on the same idea

Now we completely drop RNN part

Also we repeat self-attention many times

Further we consider a separate block:

- Multihead attention
- Feed Forward



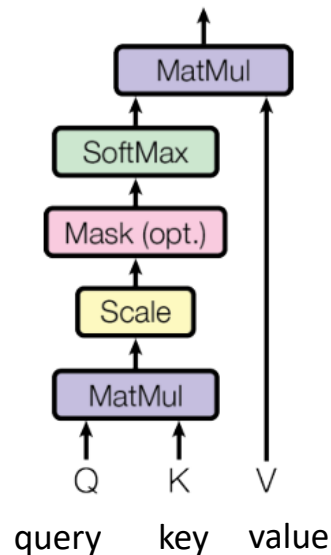
## Self-attention block

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

$d_k$  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 dot-product used here with a single-hidden layer neural network.

### Scaled Dot-Product Attention

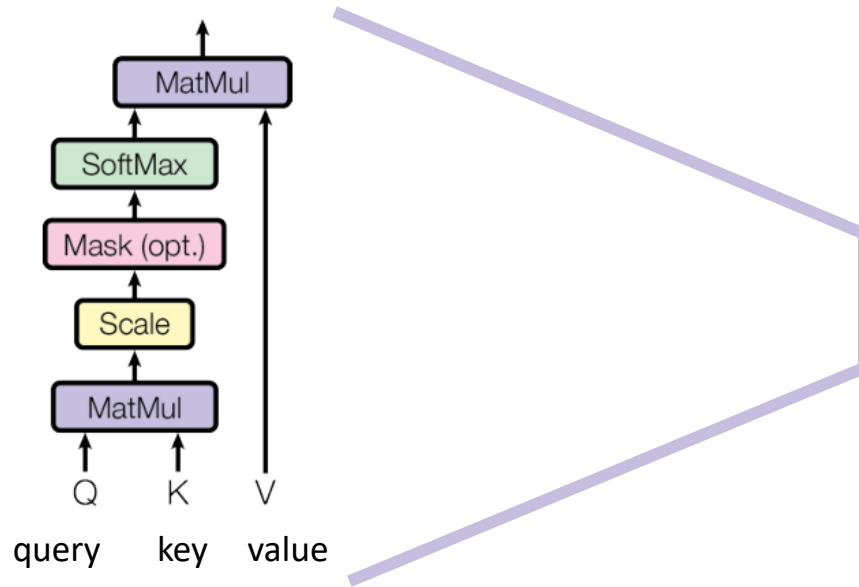


# Multi-Head attention

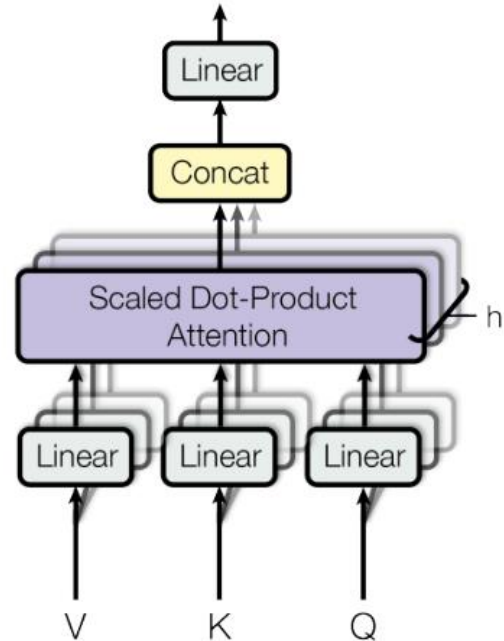
$$\text{MultiHead}(Q, K, V) = \text{Concat}(\text{head}_1, \dots, \text{head}_h) W^O$$

$$\text{where head}_i = \text{Attention}(QW_i^Q, KW_i^K, VW_i^V)$$

## Scaled Dot-Product Attention



## Multi-Head Attention



# Full block

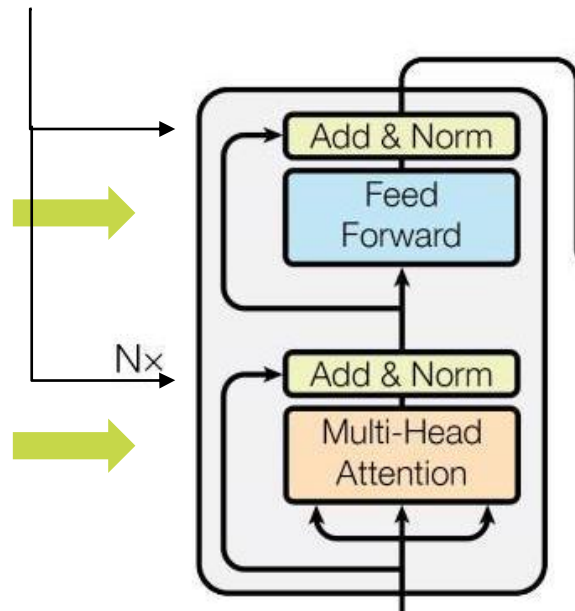
Two linear transformation with ReLU activation in between

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

$$\text{MultiHead}(Q, K, V) = \text{Concat}(\text{head}_1, \dots, \text{head}_h)W^O$$

where  $\text{head}_i = \text{Attention}(QW_i^Q, KW_i^K, VW_i^V)$

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



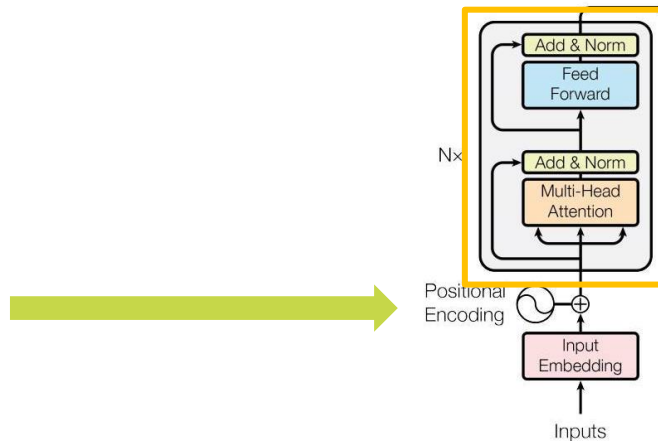


## Position encoding

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

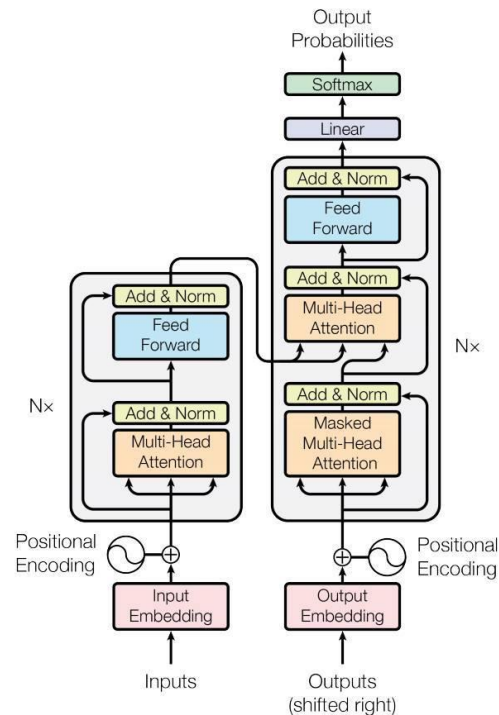
$$PE_{(pos, 2i)} = \sin(pos/10000^{2i/d_{\text{model}}})$$
$$PE_{(pos, 2i+1)} = \cos(pos/10000^{2i/d_{\text{model}}})$$



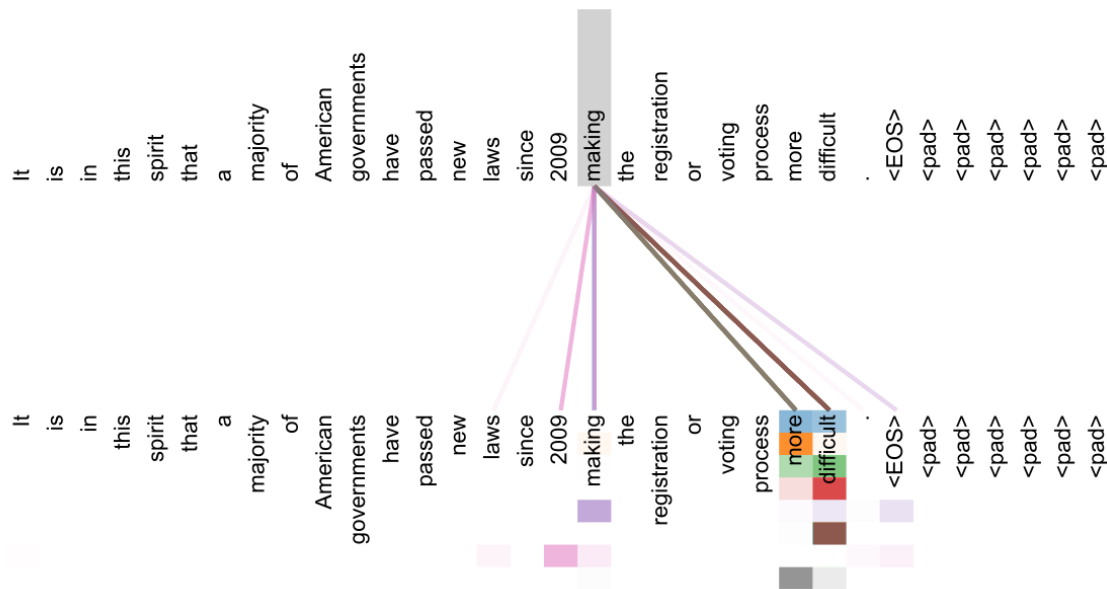
Decoder is similar

It has additional sublayer to take into account attentions from encoder

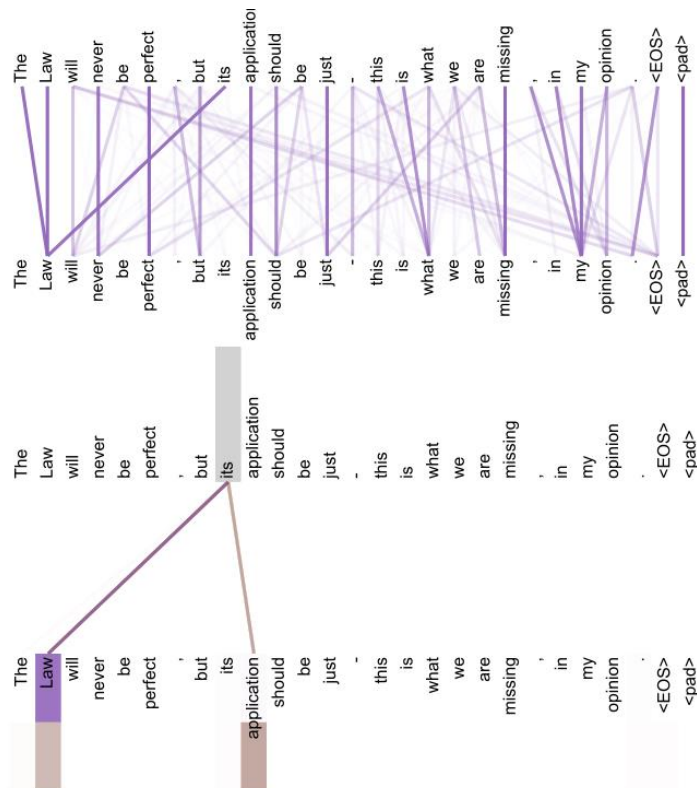
We generate one token and proceed to the next token generation



# Attention visualizations



# Attention visualizations



# Sources

“Attention is all you need” paper