

# Attention Is All You Need

Прокопьева Дарья

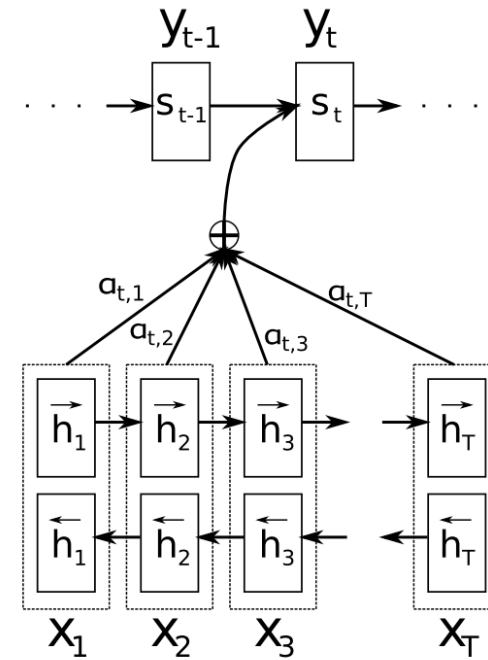
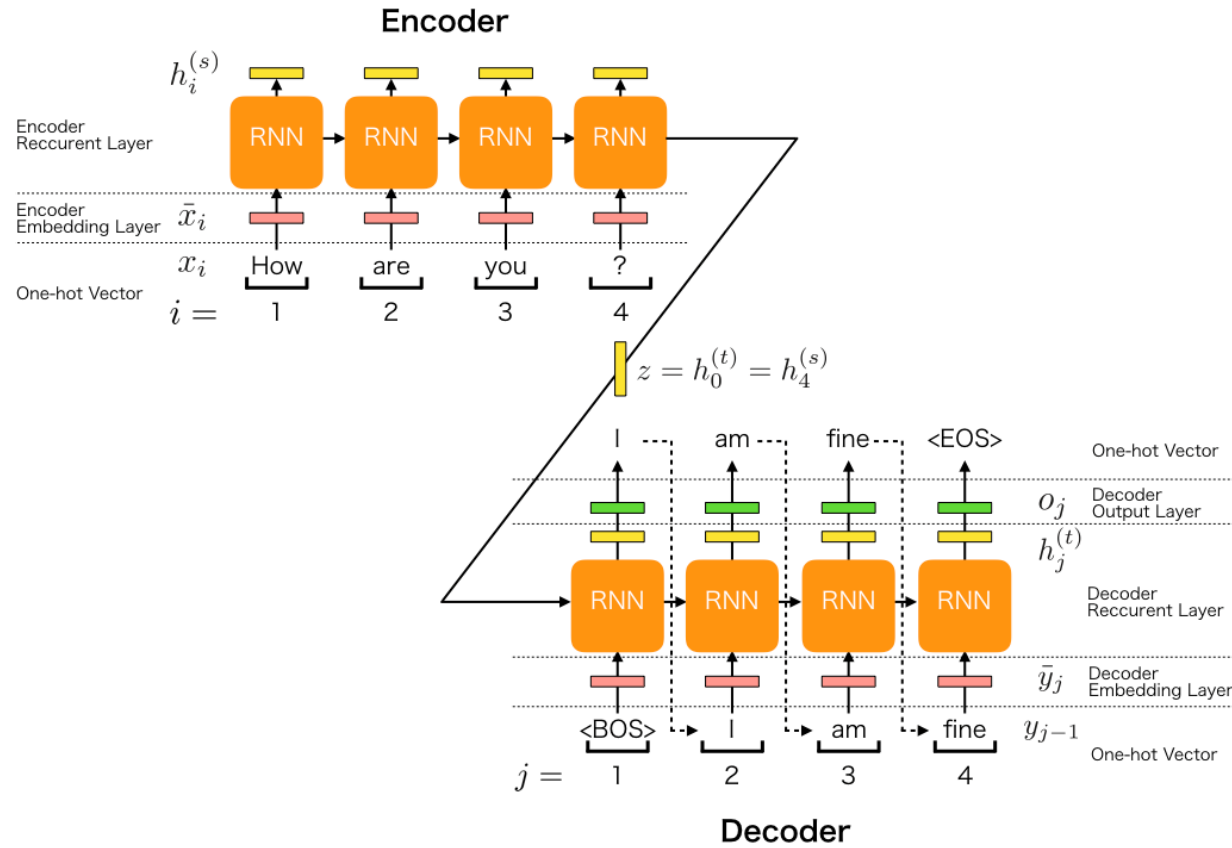
НИУ ВШЭ

2019

# Содержание

- Seq2seq, attention
- Transformer
- Результаты

# Seq2seq, attention



# Transformer – model architecture

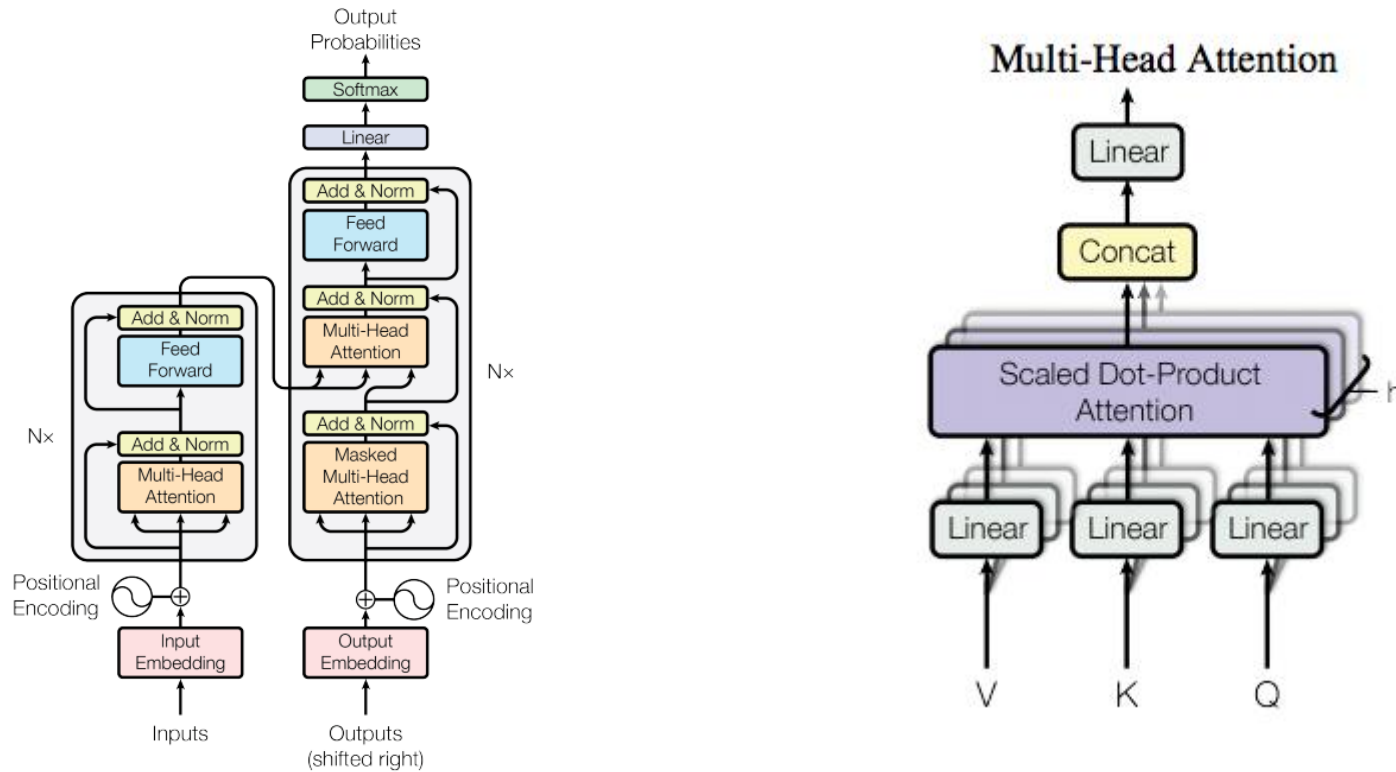
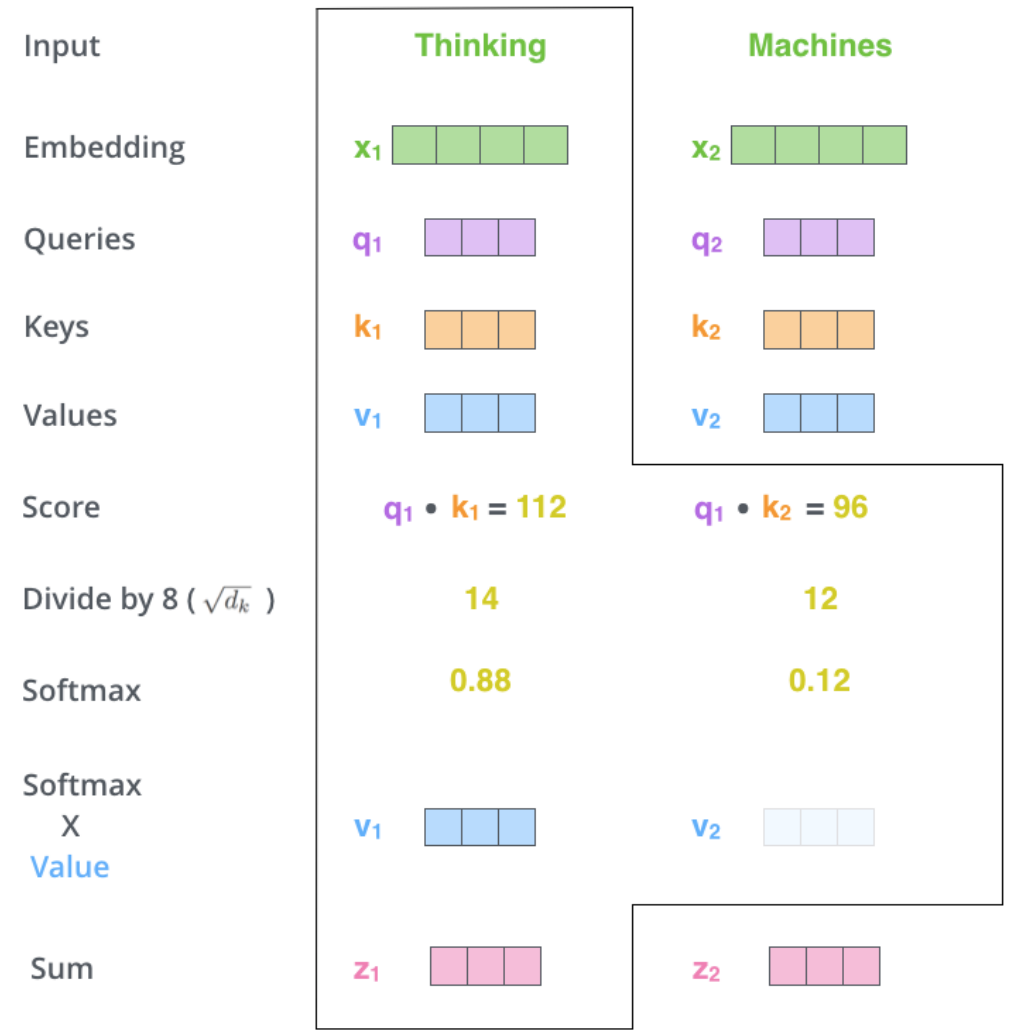


Figure 1: The Transformer - model architecture.

# Self-attention



# Self-attention

$$\mathbf{X} \times \mathbf{W}^Q = \mathbf{Q}$$


The diagram illustrates the calculation of the Query matrix  $\mathbf{Q}$ . It shows a green input matrix  $\mathbf{X}$  (2 rows by 4 columns) multiplied by a purple weight matrix  $\mathbf{W}^Q$  (4 rows by 3 columns) to produce a purple output matrix  $\mathbf{Q}$  (2 rows by 3 columns).

$$\mathbf{X} \times \mathbf{W}^K = \mathbf{K}$$


The diagram illustrates the calculation of the Key matrix  $\mathbf{K}$ . It shows a green input matrix  $\mathbf{X}$  (2 rows by 4 columns) multiplied by an orange weight matrix  $\mathbf{W}^K$  (4 rows by 3 columns) to produce an orange output matrix  $\mathbf{K}$  (2 rows by 3 columns).

$$\mathbf{X} \times \mathbf{W}^V = \mathbf{V}$$


The diagram illustrates the calculation of the Value matrix  $\mathbf{V}$ . It shows a green input matrix  $\mathbf{X}$  (2 rows by 4 columns) multiplied by a blue weight matrix  $\mathbf{W}^V$  (4 rows by 3 columns) to produce a blue output matrix  $\mathbf{V}$  (2 rows by 3 columns).

# Self-attention

$$\text{softmax}\left(\frac{\begin{matrix} \text{Q} \\ \begin{array}{|c|c|c|} \hline & & \\ \hline \end{array} \end{matrix} \times \begin{matrix} \text{K}^T \\ \begin{array}{|c|c|} \hline & \\ \hline \end{array} \end{matrix}}{\sqrt{d_k}}\right) \begin{matrix} \text{V} \\ \begin{array}{|c|c|c|} \hline & & \\ \hline \end{array} \end{matrix}$$

=

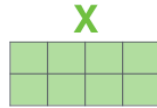
$$\begin{matrix} \text{Z} \\ \begin{array}{|c|c|c|} \hline & & \\ \hline \end{array} \end{matrix}$$

# Multi-headed self-attention

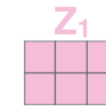
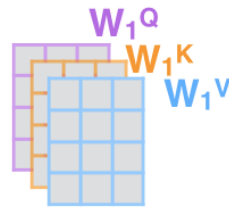
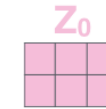
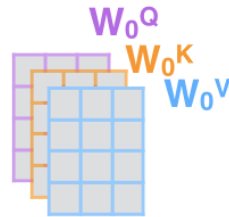
1) This is our input sentence\*

Thinking  
Machines

2) We embed each word\*



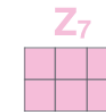
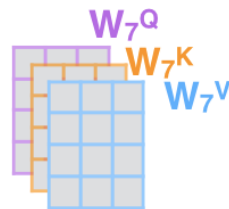
3) Split into 8 heads.  
We multiply  $X$  or  $R$  with weight matrices



...

...

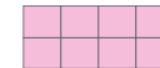
...



$W^O$



$Z$

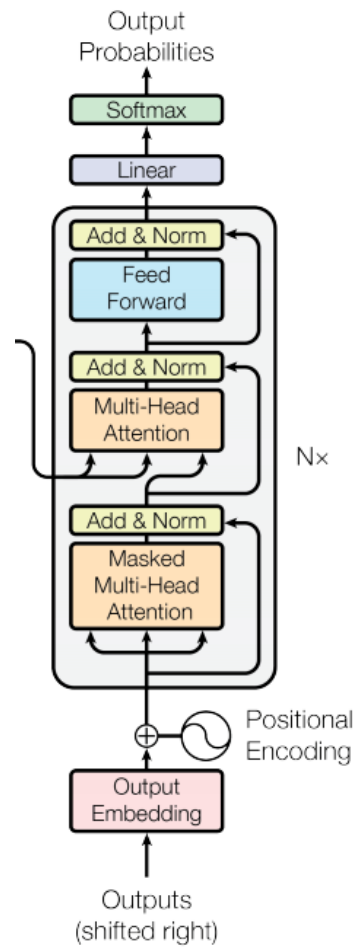


\* In all encoders other than #0,  
we don't need embedding.  
We start directly with the output  
of the encoder right below this one



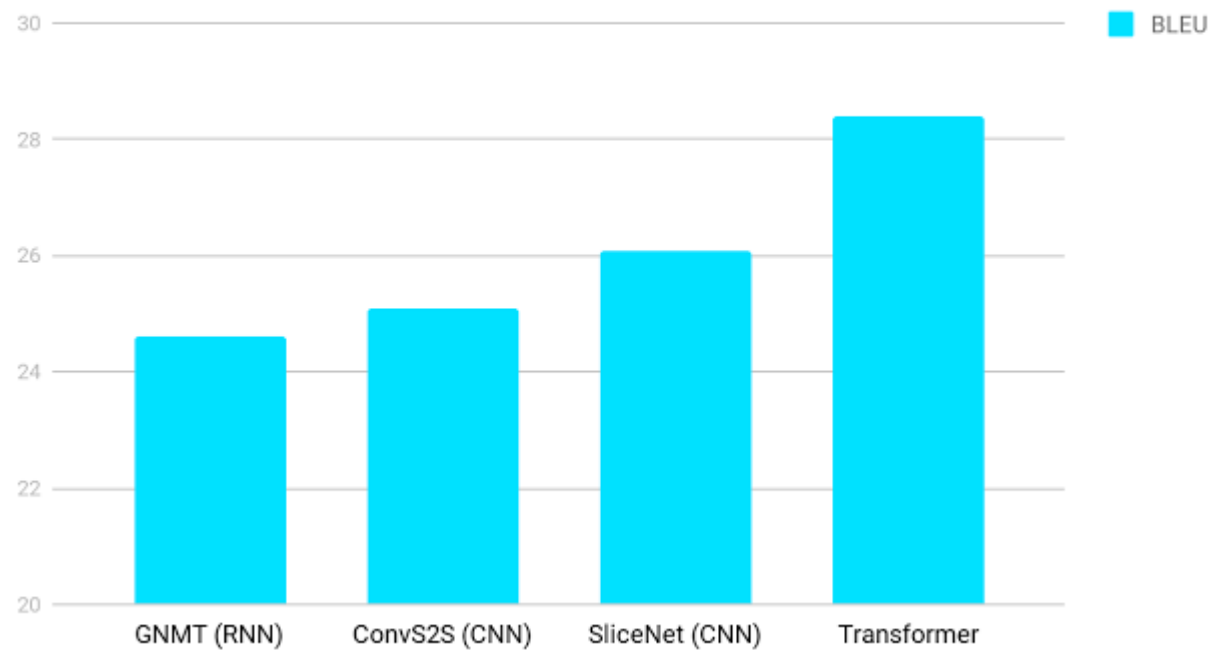


# Decoder



# Результаты

English German Translation quality



# Результаты

Layer Type	Complexity per Layer	Sequential Operations	Maximum Path Length
Self-Attention	$O(n^2 \cdot d)$	$O(1)$	$O(1)$
Recurrent	$O(n \cdot d^2)$	$O(n)$	$O(n)$
Convolutional	$O(k \cdot n \cdot d^2)$	$O(1)$	$O(\log_k(n))$
Self-Attention (restricted)	$O(r \cdot n \cdot d)$	$O(1)$	$O(n/r)$

Model	BLEU		Training Cost (FLOPs)	
	EN-DE	EN-FR	EN-DE	EN-FR
ByteNet [18]	23.75			
Deep-Att + PosUnk [39]		39.2		$1.0 \cdot 10^{20}$
GNMT + RL [38]	24.6	39.92	$2.3 \cdot 10^{19}$	$1.4 \cdot 10^{20}$
ConvS2S [9]	25.16	40.46	$9.6 \cdot 10^{18}$	$1.5 \cdot 10^{20}$
MoE [32]	26.03	40.56	$2.0 \cdot 10^{19}$	$1.2 \cdot 10^{20}$
Deep-Att + PosUnk Ensemble [39]		40.4		$8.0 \cdot 10^{20}$
GNMT + RL Ensemble [38]	26.30	41.16	$1.8 \cdot 10^{20}$	$1.1 \cdot 10^{21}$
ConvS2S Ensemble [9]	26.36	<b>41.29</b>	$7.7 \cdot 10^{19}$	$1.2 \cdot 10^{21}$
Transformer (base model)	27.3	38.1	<b><math>3.3 \cdot 10^{18}</math></b>	
Transformer (big)	<b>28.4</b>	<b>41.8</b>	$2.3 \cdot 10^{19}$	

# ИСТОЧНИКИ

- <https://arxiv.org/pdf/1706.03762.pdf>
- <http://jalammar.github.io/illustrated-transformer/>
- <https://habr.com/ru/post/341240/>