

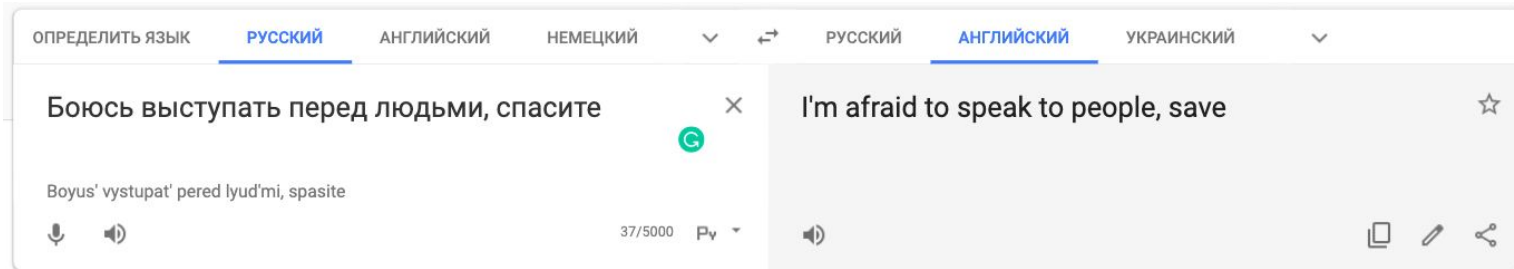
Attention Is All You Need



Rak Arina, AMI171
HSE Research Seminar
2019

● **Sequence transduction** — transformation of input sequences into output sequences:

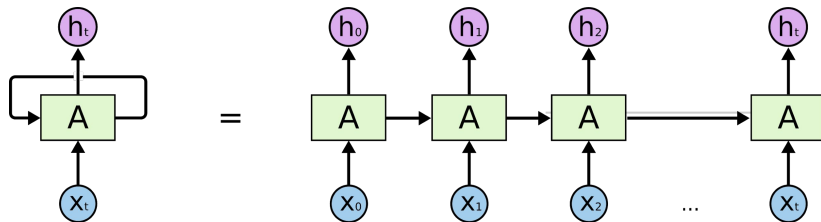
- Machine translation
- Speech recognition
- Spelling correction
- Part of speech tagging



Recap | RNNs

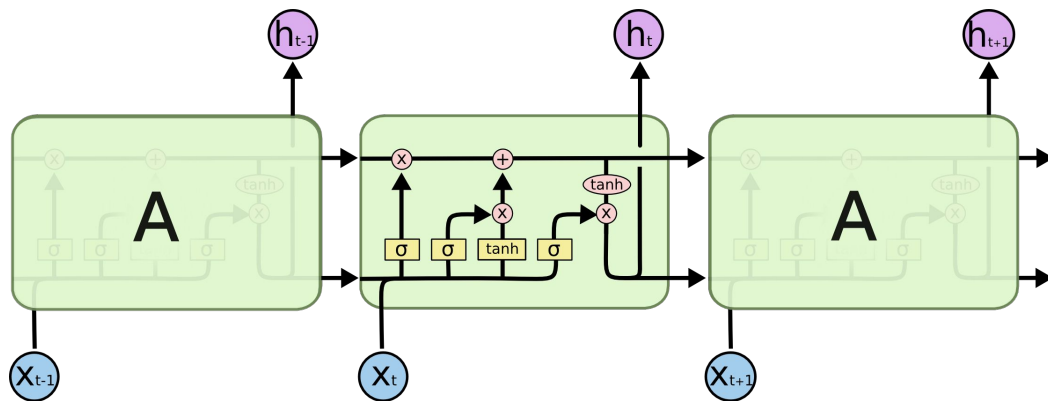


Vanilla
RNN

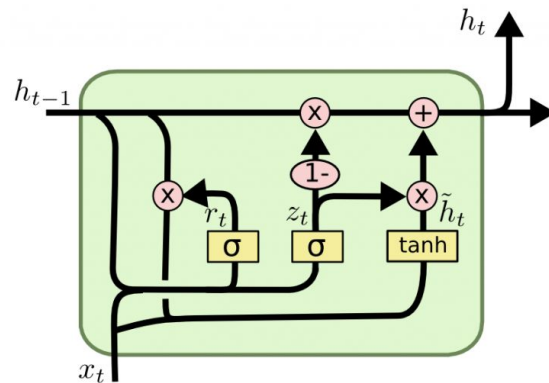


Problems:

- Learns slow
- Vanishing/exploding gradients
- Difficult to learn dependencies between distant positions



LSTM

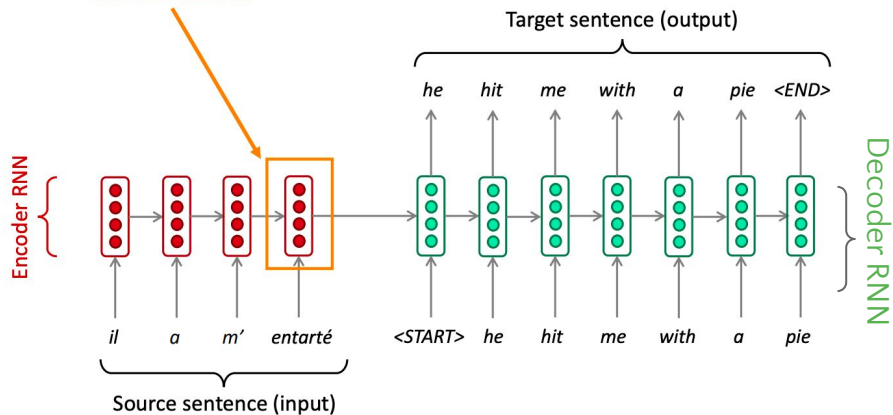


GRU

Recap | Attention



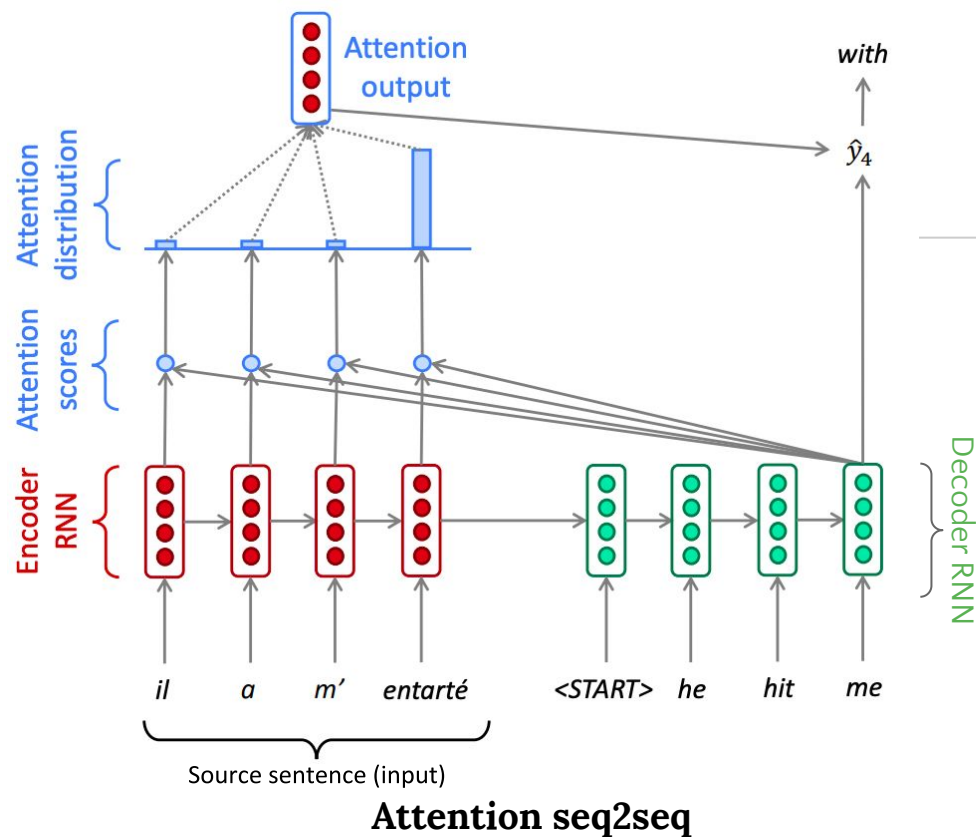
Encoding of the source sentence.



Classical seq2seq

Attention:

- Improves performance
- Helps with vanishing gradients
- Solves the bottleneck problem
- Helps with interpretability



BUT: Models get more and more complex and the computations still can not be done in parallel => **SLOW**



Transformer

Was proposed in 2017 by Google

In WMT (MT conference + competition):
The summary report in 2016 contains the word 'RNN' 44 times

The summary report in 2018 contains the word 'RNN' 9 times and the word 'Transformer' 63 times

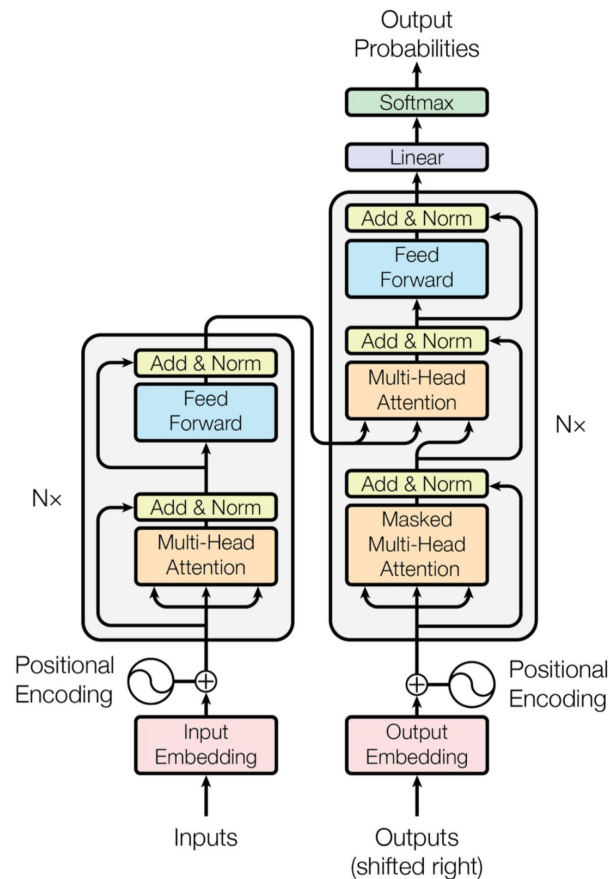
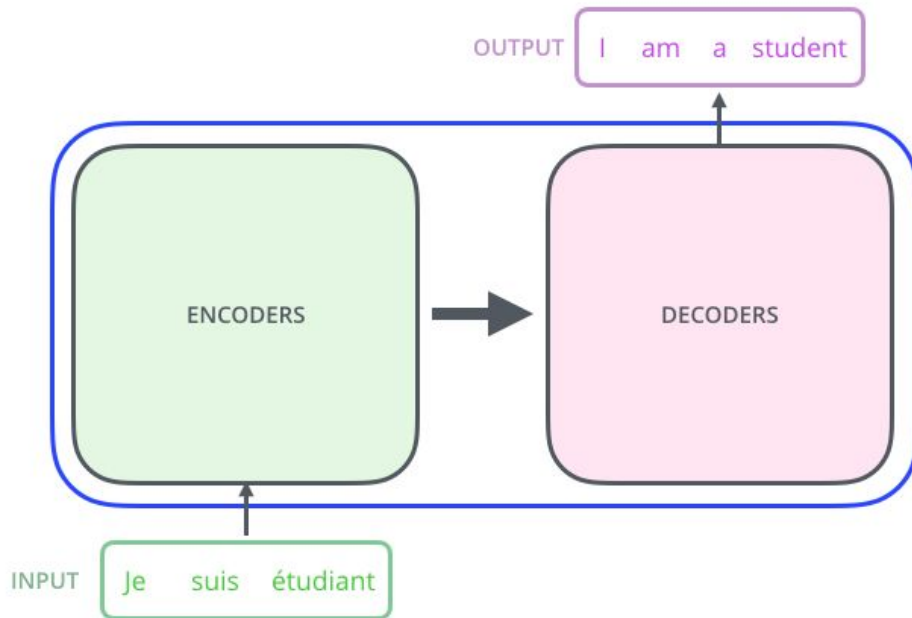


Figure 1: The Transformer - model architecture.



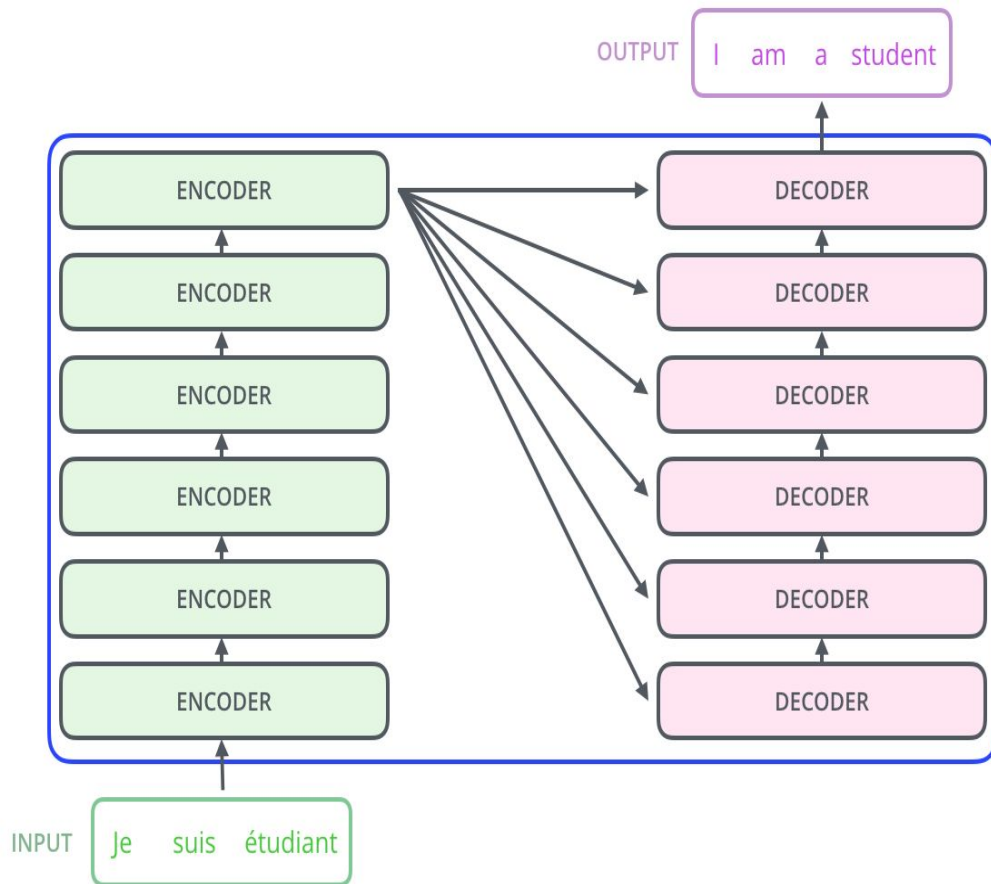
Architecture

Encoder receives a list of fixed size of vectors each of the size of the embeddings dimensionality

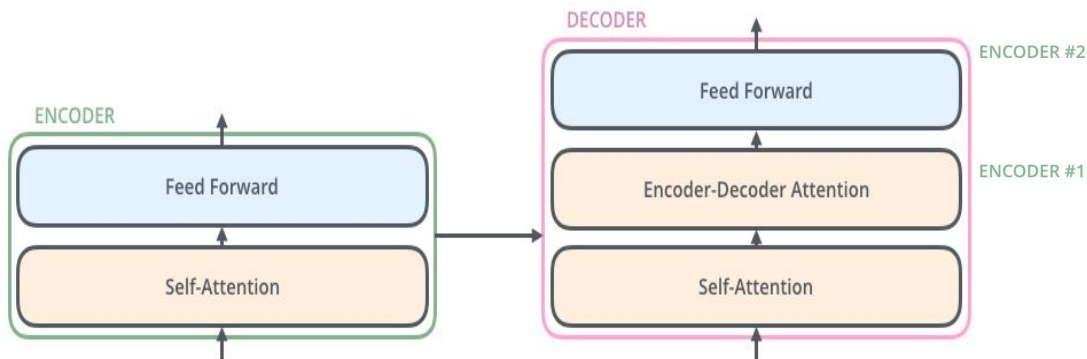


Architecture

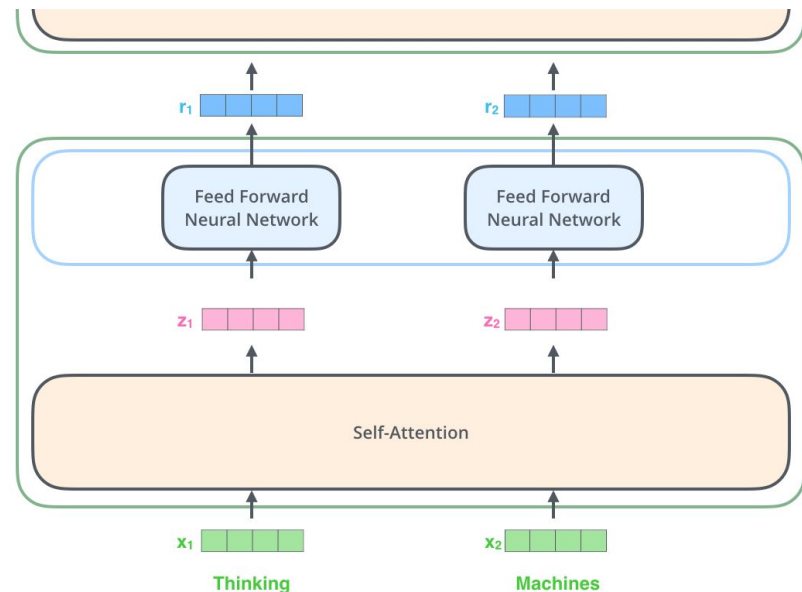
All Encoder and Decoder blocks have the same architecture, but they do not share weights.



Self-attention is an attention mechanism relating different positions of a single sequence in order to compute a representation of the sequence.



Self-attention in Encoder and Decoder are the same, except for the fact that Decoder can look only on the words previous to the current one

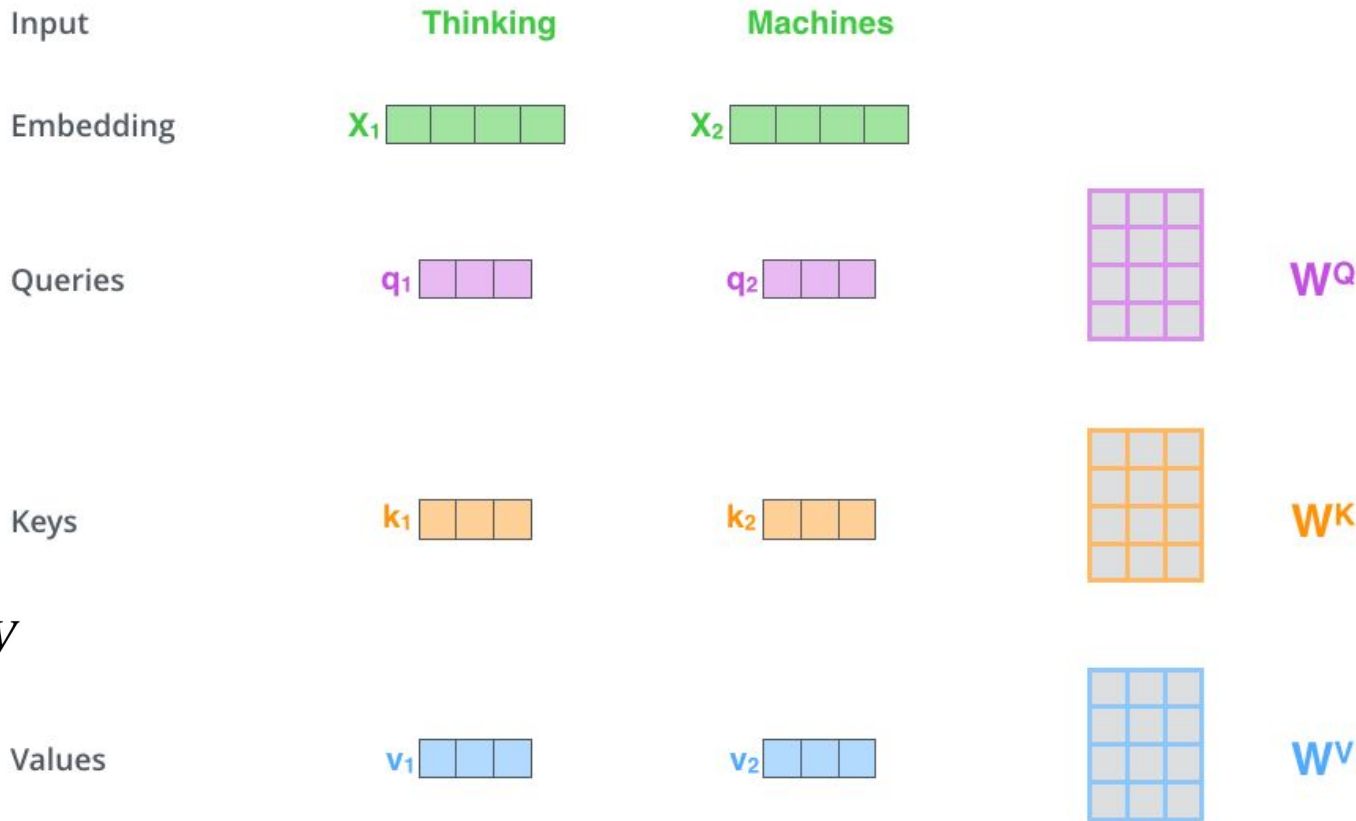


Self-Attention



For each word x
a query q , a key
 k and a value v
is calculated by
multiplying x by
matrices

W^Q , W^K and W^V
respectively



Self-Attention



For every x_i
 $\langle q_i, k_j \rangle, j \in \{1, \dots, input_len\}$
 is calculated, which is a score
 reflecting the weight of v_j
 in the representation of x_i

$$z_i = \sum_{j=1}^{input_len} w_j v_j$$

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

Input

Embedding

Queries

Keys

Values

Score

Divide by 8 ($\sqrt{d_k}$)

Softmax

Softmax
X
Value

Sum

Thinking

x_1

q_1

k_1

v_1

$q_1 \cdot k_1 = 112$

14

w_1 0.88

v_1

z_1

Machines

x_2

q_2

k_2

v_2

$q_2 \cdot k_2 = 96$

12

w_2 0.12

v_2

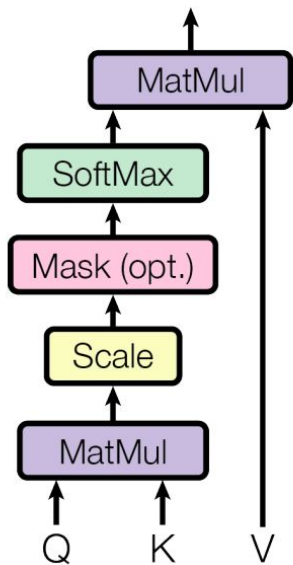
z_2

Self-Attention

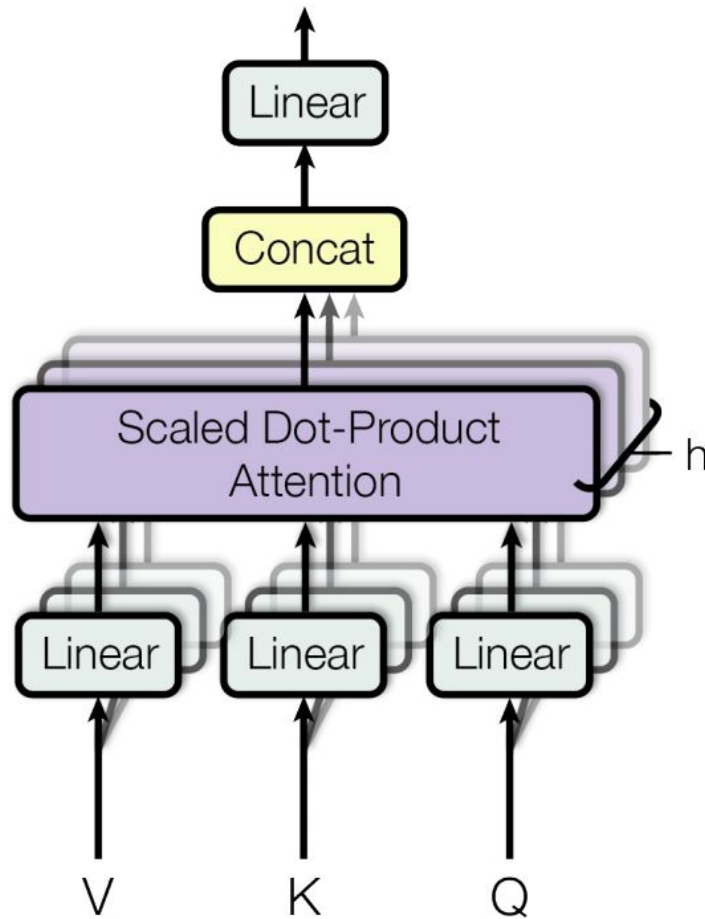


Multi-Head Attention just does the same thing h times and then concatenates the results and projects it back to the dimension of x with a linear layer

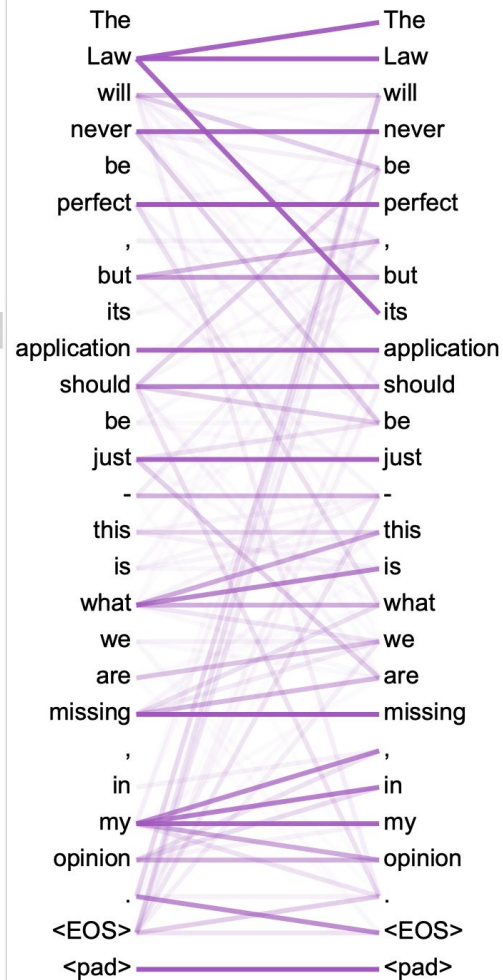
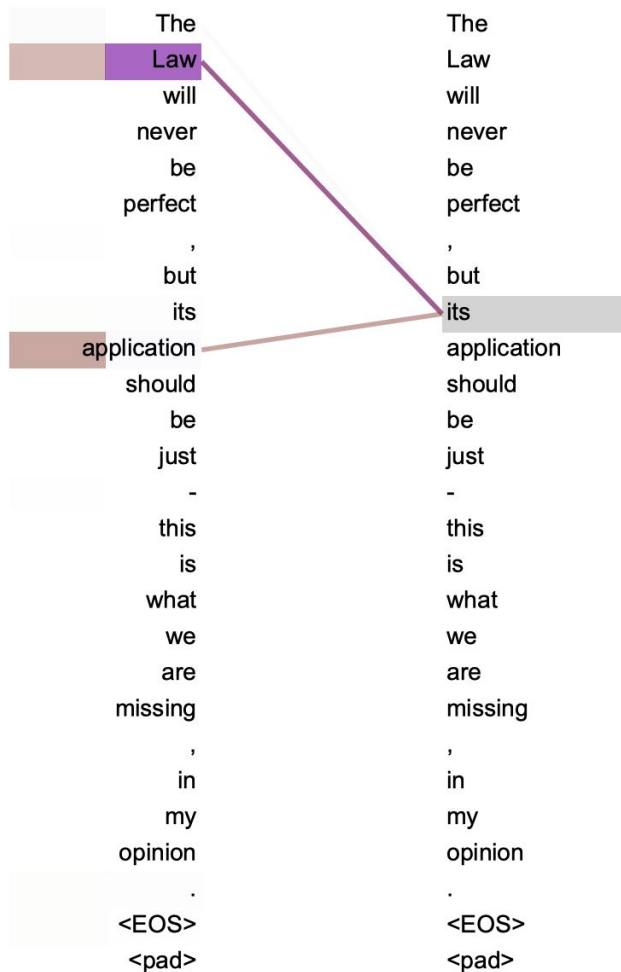
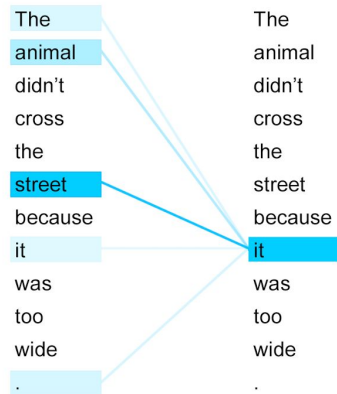
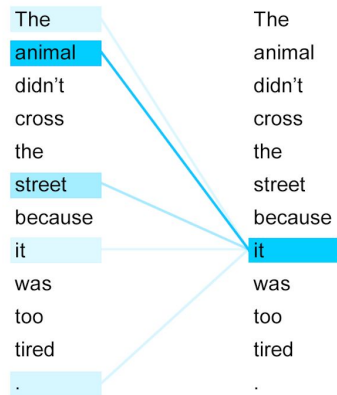
Scaled Dot-Product Attention



Multi-Head Attention



Self-Attention examples



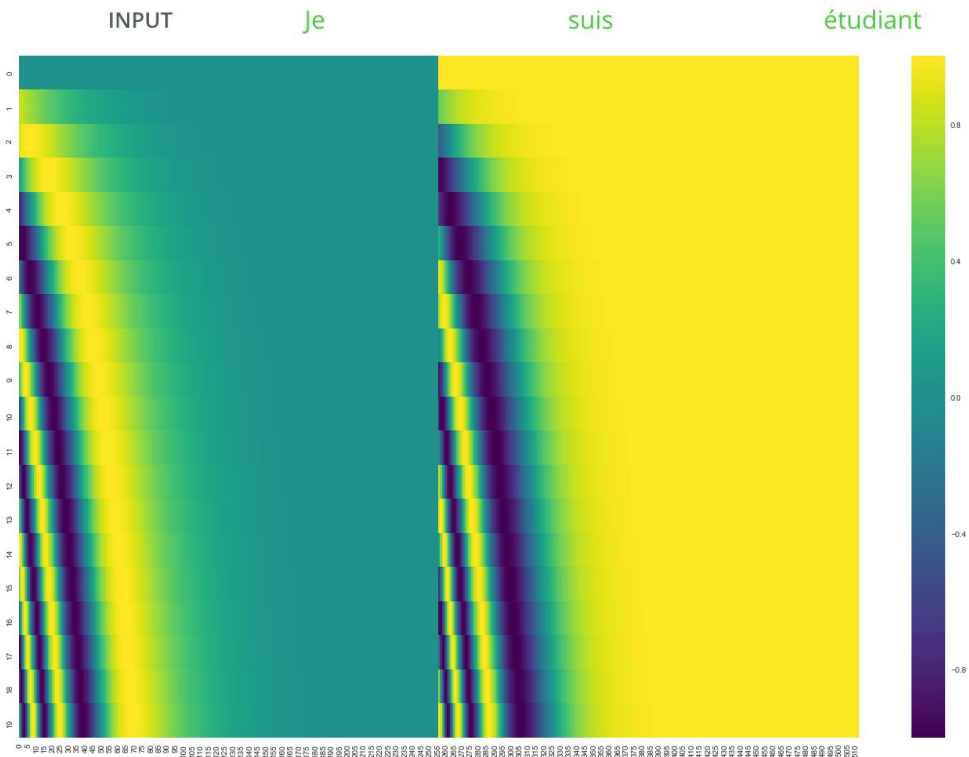
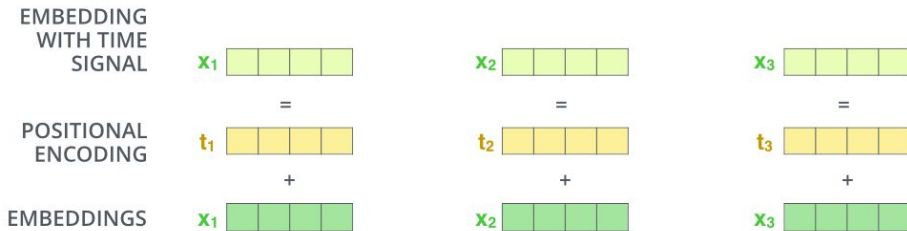
Positional encoding



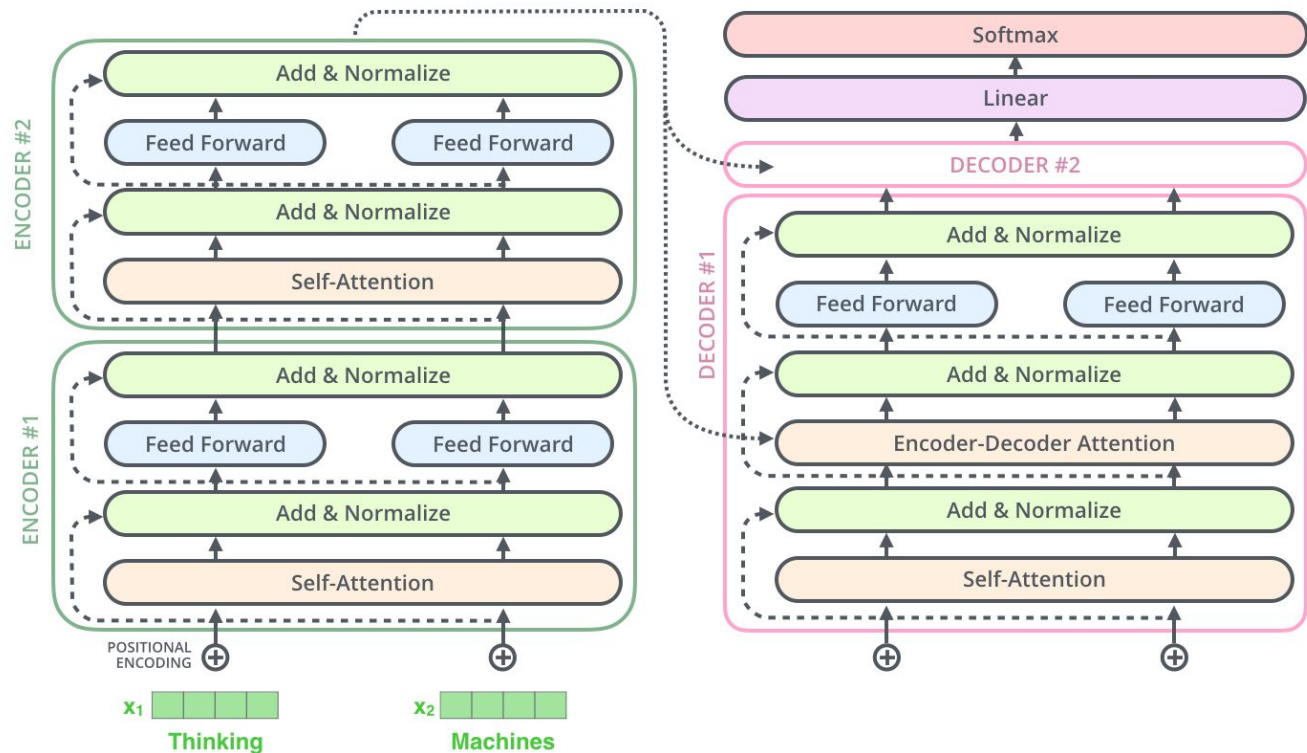
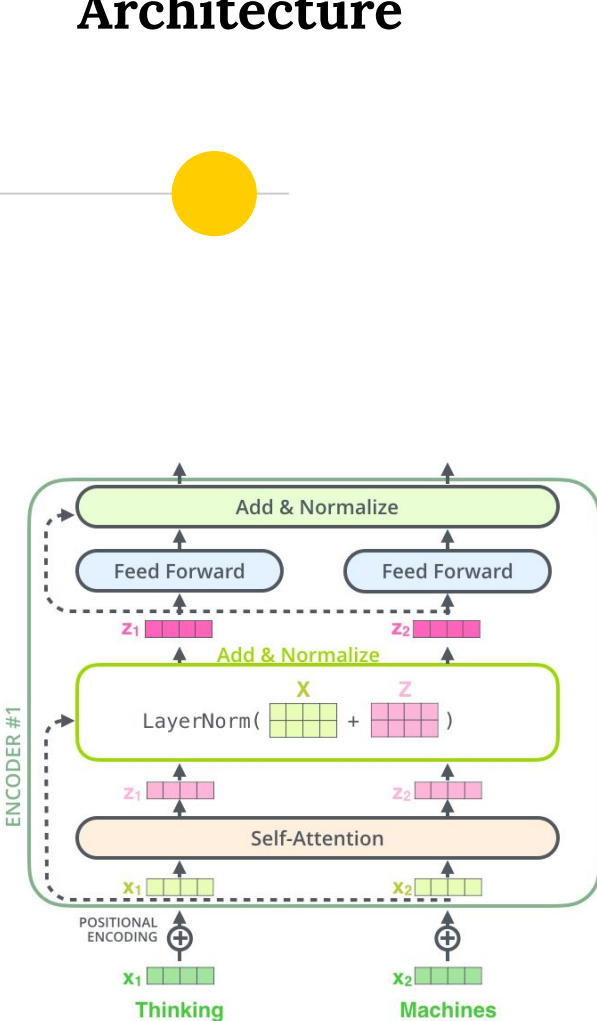
In order for the model to make use of the order of the sequence, some information about the relative or absolute position of the tokens of the sequence must be injected

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

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



Architecture



There are skip connections around each of the sublayers of the blocks and layer normalization after each sublayer.

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)	$O(r \cdot n \cdot d)$	$O(1)$	$O(n/r)$

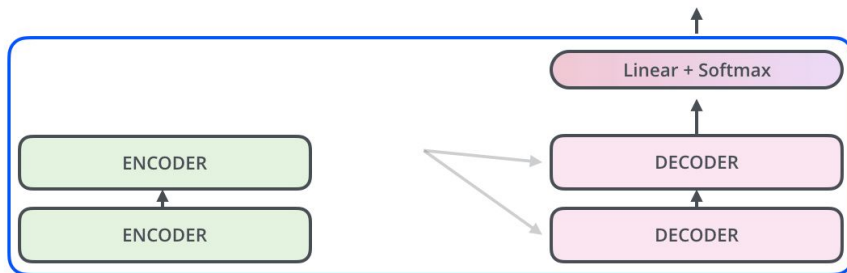
Motivating our use of self-attention we consider **three desiderata**.

- One is the total computational complexity per layer.
- Another is the amount of computation that can be parallelized, as measured by the minimum number of sequential operations required.
- The third is the path length between long-range dependencies in the network. Learning long-range dependencies is a key challenge in many sequence transduction tasks.

“

Decoding time step: 1 2 3 4 5 6

OUTPUT



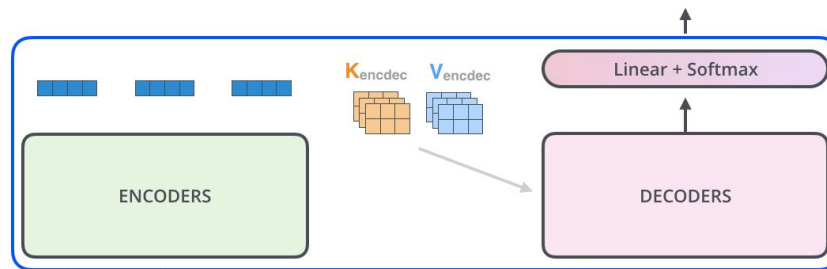
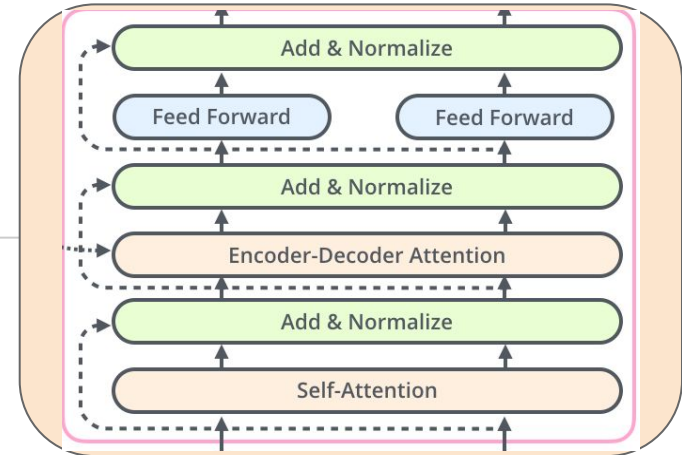
EMBEDDING WITH TIME SIGNAL
[][][][] [][][][] [][][][]

EMBEDDINGS
[][][][] [][][][] [][][][]

INPUT Je suis étudiant

Decoding time step: 1 2 3 4 5 6

OUTPUT |



EMBEDDING WITH TIME SIGNAL
[][][][] [][][][] [][][][]

EMBEDDINGS
[][][][] [][][][] [][][][]

INPUT Je suis étudiant

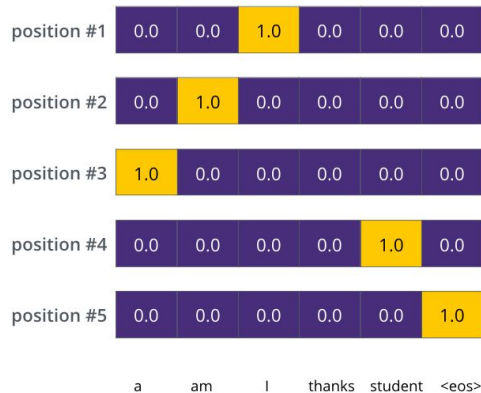
PREVIOUS OUTPUTS |



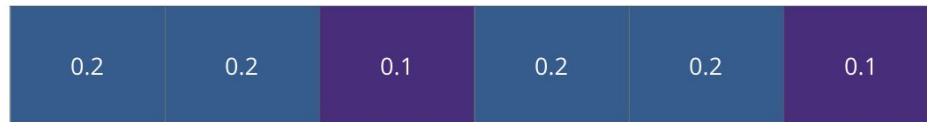
Training

Target Model Outputs

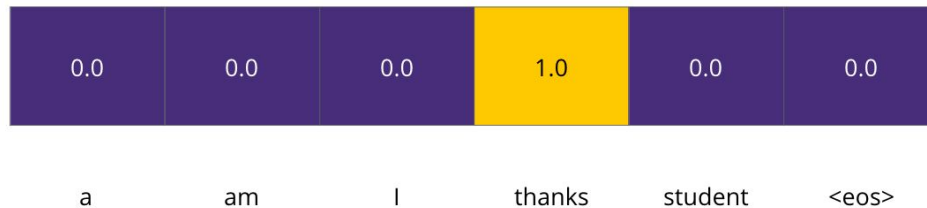
Output Vocabulary: a am I thanks student <eos>



Untrained Model Output



Correct and desired output



The output is a probability distribution.

Loss:

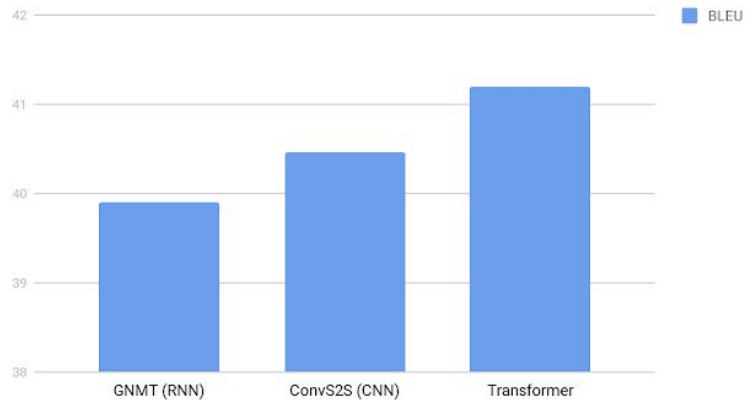
- Cross-entropy:

$$-\sum_i y_{true_i} \log(p_{pred_i})$$

- KL - divergence:

$$-\sum_i p_{true_i} \log\left(\frac{p_{true_i}}{p_{pred_i}}\right)$$

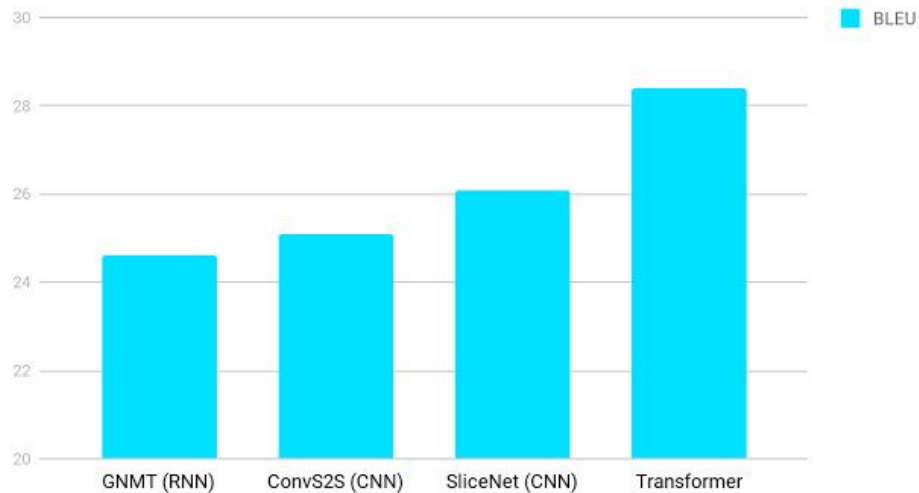
English French Translation Quality

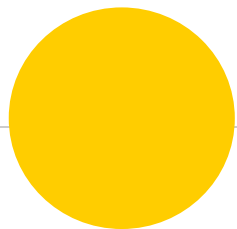


Transformers are used in both Google and Yandex Translate

Results in 2017

English German Translation quality





Вопросы

- Что подается на вход Encoder-у трансформера?
- Что такое q , k , v в слое self-attention?
- Чем отличаются слои self-attention у Encoder-а и Decoder-а?



Resources

- ◉ <https://arxiv.org/pdf/1706.03762.pdf> (original paper)
- ◉ <https://habr.com/ru/post/341240/> (less papers but in Russian)
- ◉ <http://jalammar.github.io/illustrated-transformer/> (best pics and explained quite nicely)
- ◉ https://www.youtube.com/watch?v=S0KakHcj_rs&t=1132s (video on the paper)
- ◉ <https://www.youtube.com/watch?v=QEw0qEa0E50&feature=youtu.be> (CS224n, Stanford's course on NLP)
- ◉ <https://ai.googleblog.com/2017/08/transformer-novel-neural-network.html>