

# Attention Is All You Need

Authors: Ashish Vaswani, Noam Shazeer, Niki Parmar, Jakob Uszkoreit, Llion Jones, Aidan N. Gomez, Łukasz Kaiser, Illia Polosukhin

From: Google brain Google research

NIPS 2017

Slides prepared by: Chandra Thapa (July 2020)

# Applications of Attention mechanism

- Cyber-security:
  - Attention-Based Automated Feature Extraction for Malware Analysis [1]
  - Neural Malware Analysis with Attention Mechanism [2]
  - How to Make Attention Mechanisms More Practical in Malware Classification [3]
  - I-MAD: A Novel Interpretable Malware Detector Using Hierarchical Transformer [4]
- Phishing email detection
  - Phishing Email Detection Using Improved RCNN Model With Multilevel Vectors and Attention Mechanism [5]
- Fake news detection
  - Self Multi-Head Attention-based Convolutional Neural Networks for fake news detection [6]
  - HGAT: Hierarchical Graph Attention Network for Fake News Detection [7]

[1] <https://www.mdpi.com/1424-8220/20/10/2893>

[2] <https://www.sciencedirect.com/science/article/pii/S0167404819300264>

[3] <https://ieeexplore.ieee.org/stamp/stamp.jsp?arnumber=8876839>

[4] <https://arxiv.org/abs/1909.06865>

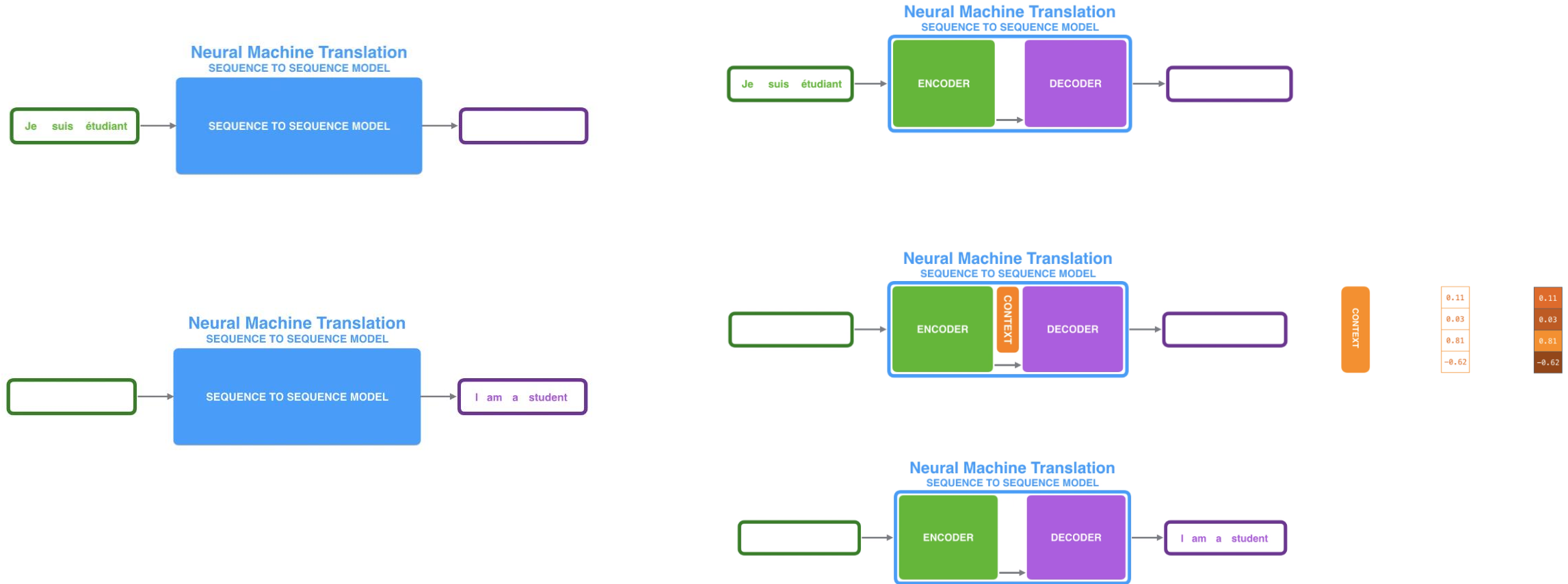
[5] <https://ieeexplore.ieee.org/abstract/document/8701426>

[6] <https://journals.plos.org/plosone/article?id=10.1371/journal.pone.0222713>

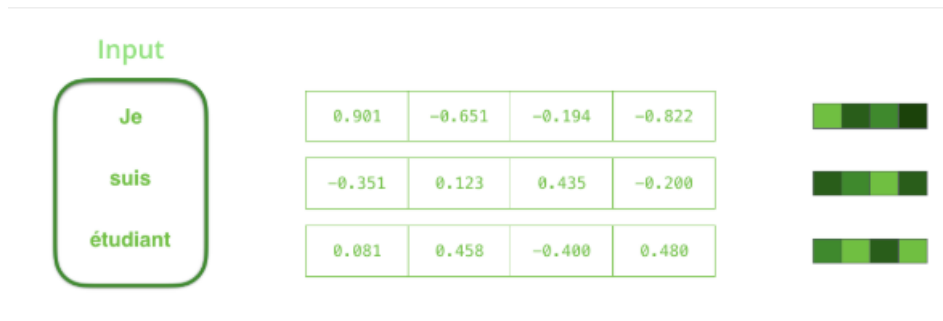
[7] <https://arxiv.org/abs/2002.04397>

# Background

- In machine translation: Sequence-to-sequence model



- Word embeddings



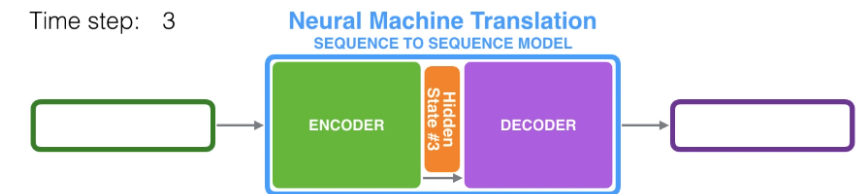
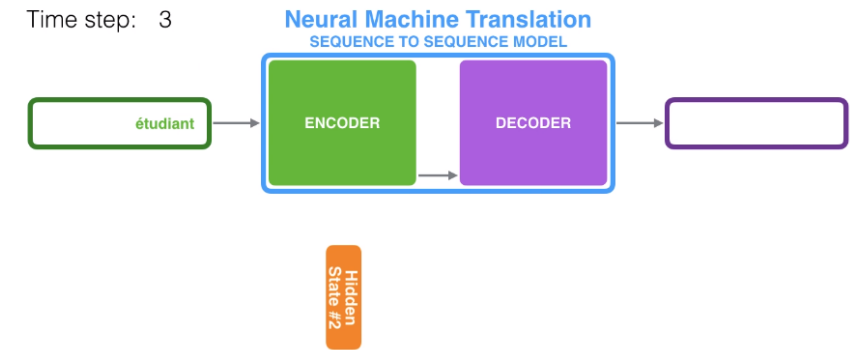
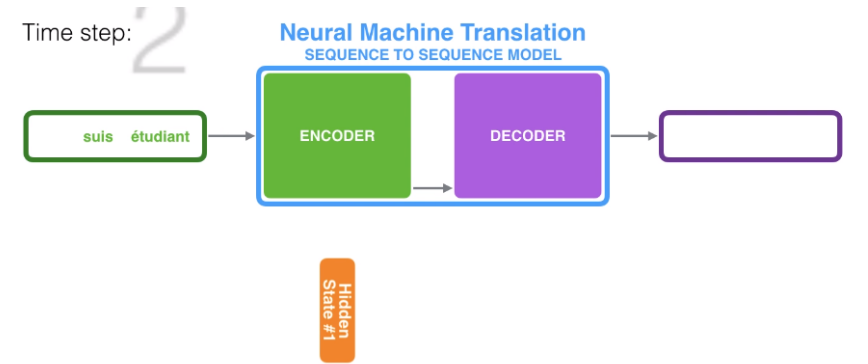
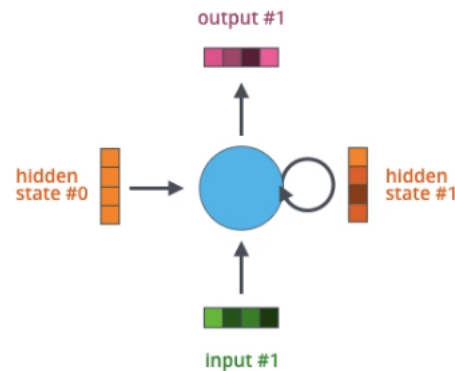
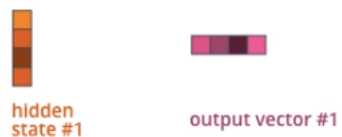
- A RNN takes two inputs at each time step:
  - Input (e.g., one word of the input sentence, if encoder)
  - Hidden state

An RNN takes two input vectors:



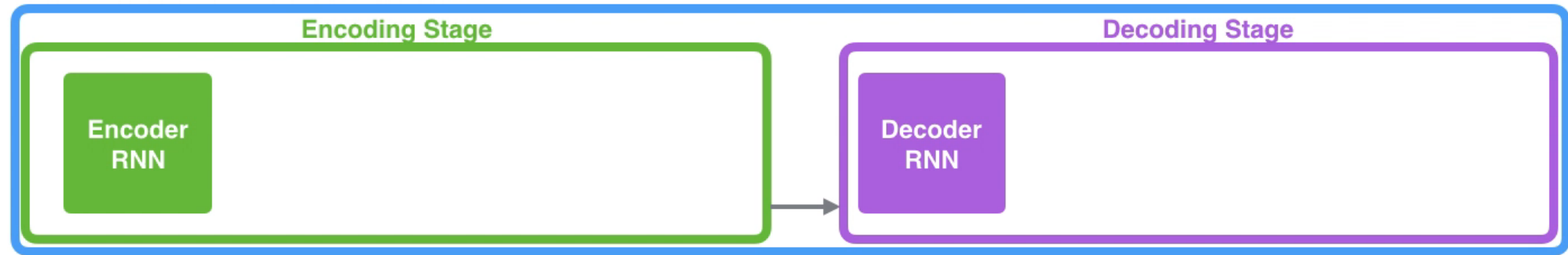
Processes them

Then produces two output vectors:



# Neural Machine Translation

## SEQUENCE TO SEQUENCE MODEL

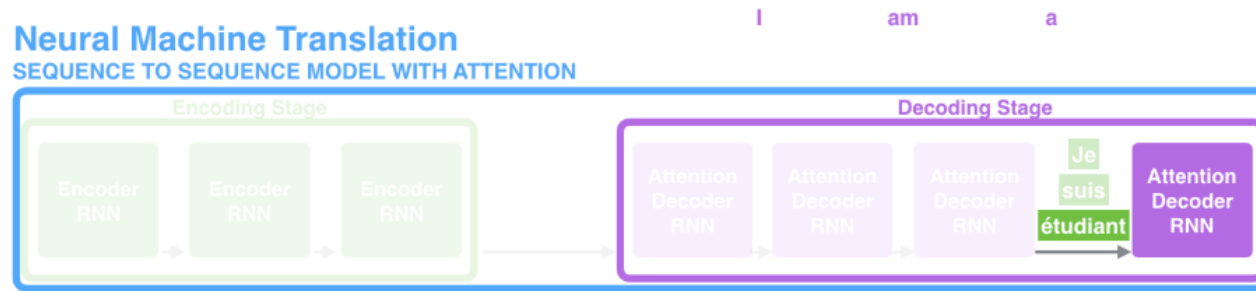


Je

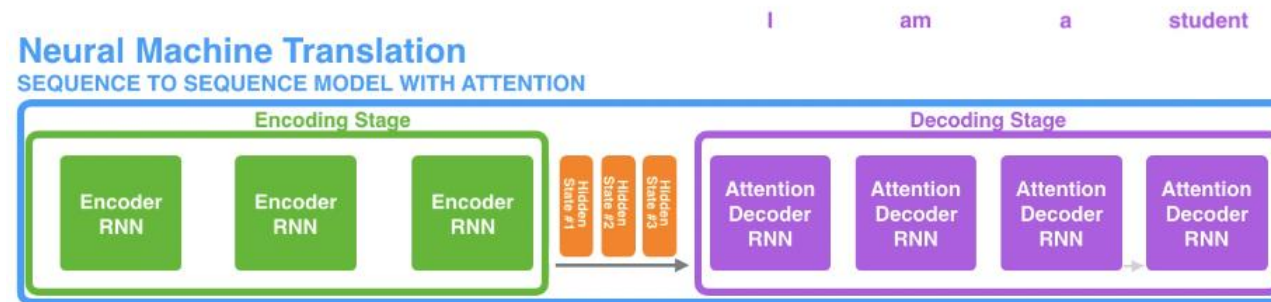
suis

étudiant

- Challenging if long sequence – context vector is a bottleneck
- (Classic) Attention is introduced [1,2] – it allows the model to focus on the relevant parts of the input sequences as needed



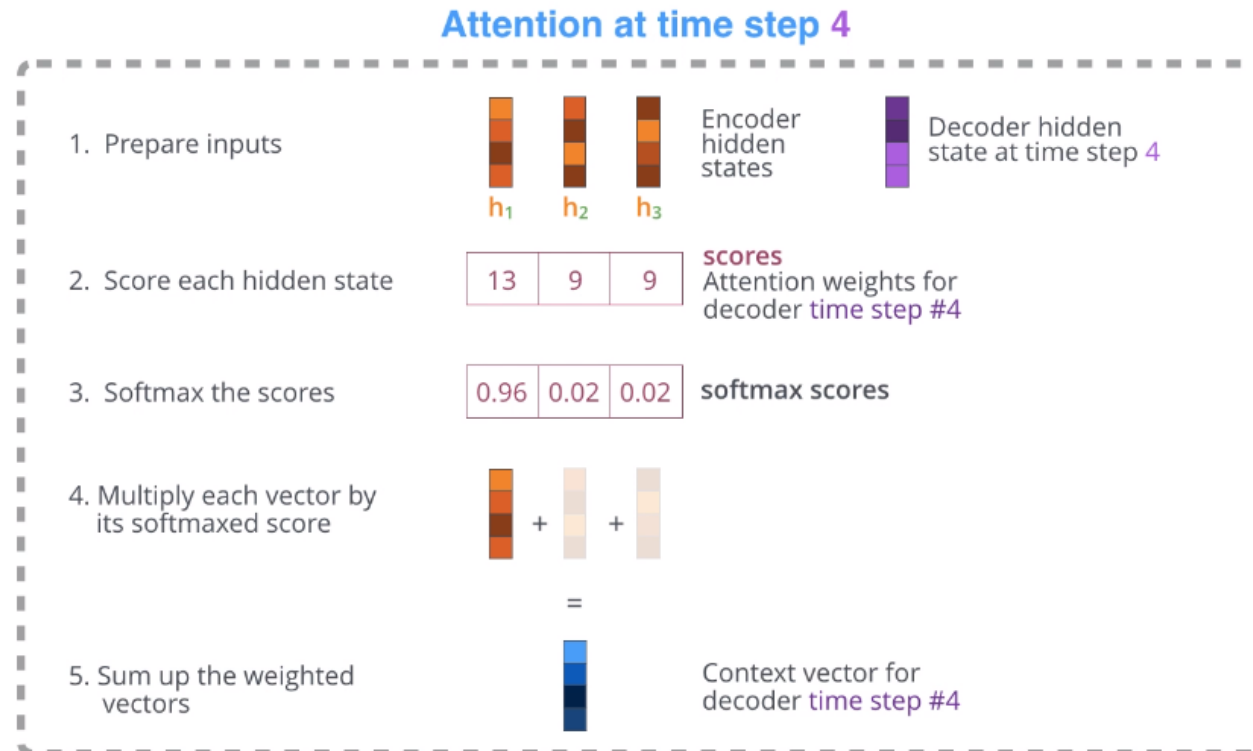
- An attention model differs from a classic sequence-to-sequence model in the following ways:
  - the encoder passes *all* the hidden states to the decoder



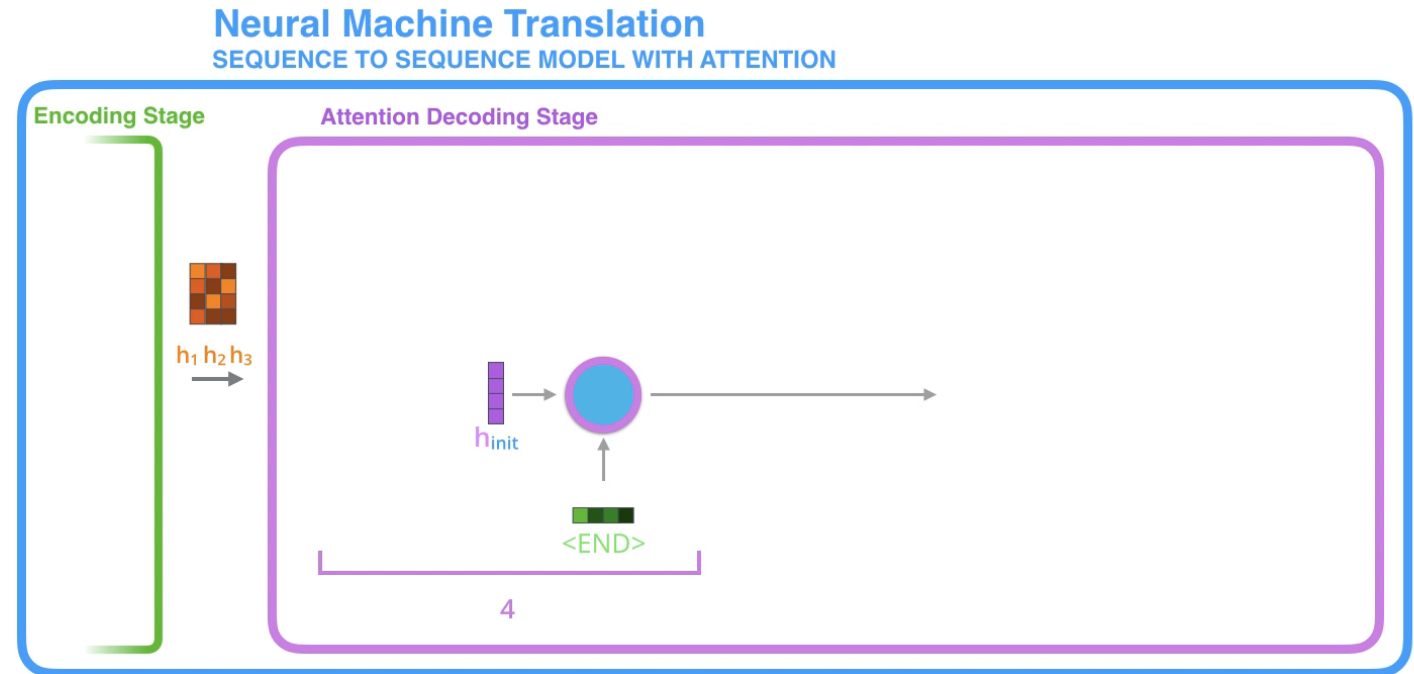
[1] [Dzmitry Bahdanau, Kyunghyun Cho, Yoshua Bengio](https://arxiv.org/abs/1409.0473), Neural Machine Translation by Jointly Learning to Align and Translate, <https://arxiv.org/abs/1409.0473>

[2] [Minh-Thang Luong, Hieu Pham, Christopher D. Manning](https://arxiv.org/abs/1508.04025), Effective Approaches to Attention-based Neural Machine Translation, <https://arxiv.org/abs/1508.04025>

- Attention decoder does the following:
  - Look at the set of encoder hidden states it received – each encoder hidden states is most associated with a certain word in the input sentence
  - Give each hidden states a score
  - Multiply each hidden states by its softmaxed score, thus amplifying hidden states with high scores, and drowning out hidden states with low scores



- How the attention process works?
  - The attention decoder RNN takes in the embedding of the <END> token, and an initial decoder hidden state.
  - The RNN processes its inputs, producing an output and a new hidden state vector ( $h_4$ ). The output is discarded.
  - Attention Step: Use the encoder hidden states and the  $h_4$  vector to calculate a context vector ( $C_4$ ) for this time step.
  - Concatenate  $h_4$  and  $C_4$  into one vector.
  - Pass this vector through a feedforward neural network (one trained jointly with the model).
  - The output of the feedforward neural networks indicates the output word of this time step.
  - Repeat for the next time steps

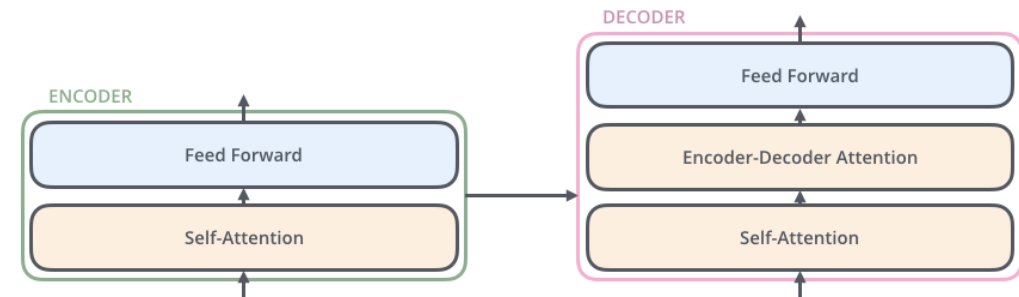
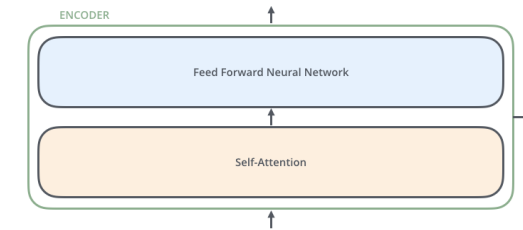
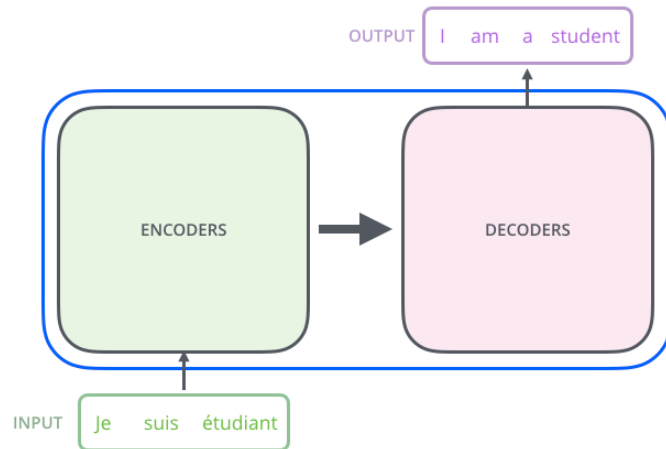
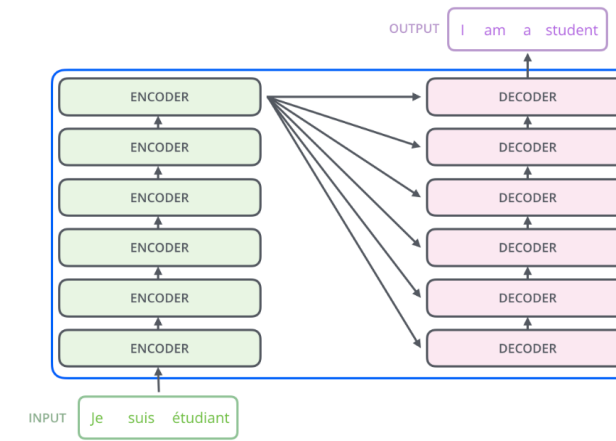




Transformer

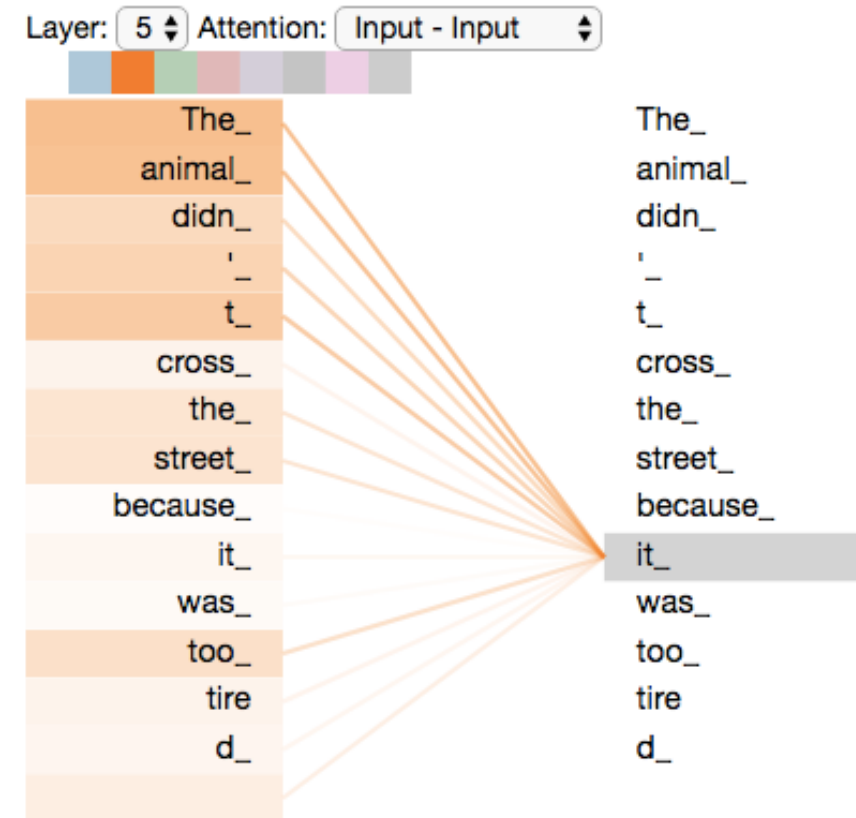


# High-level look



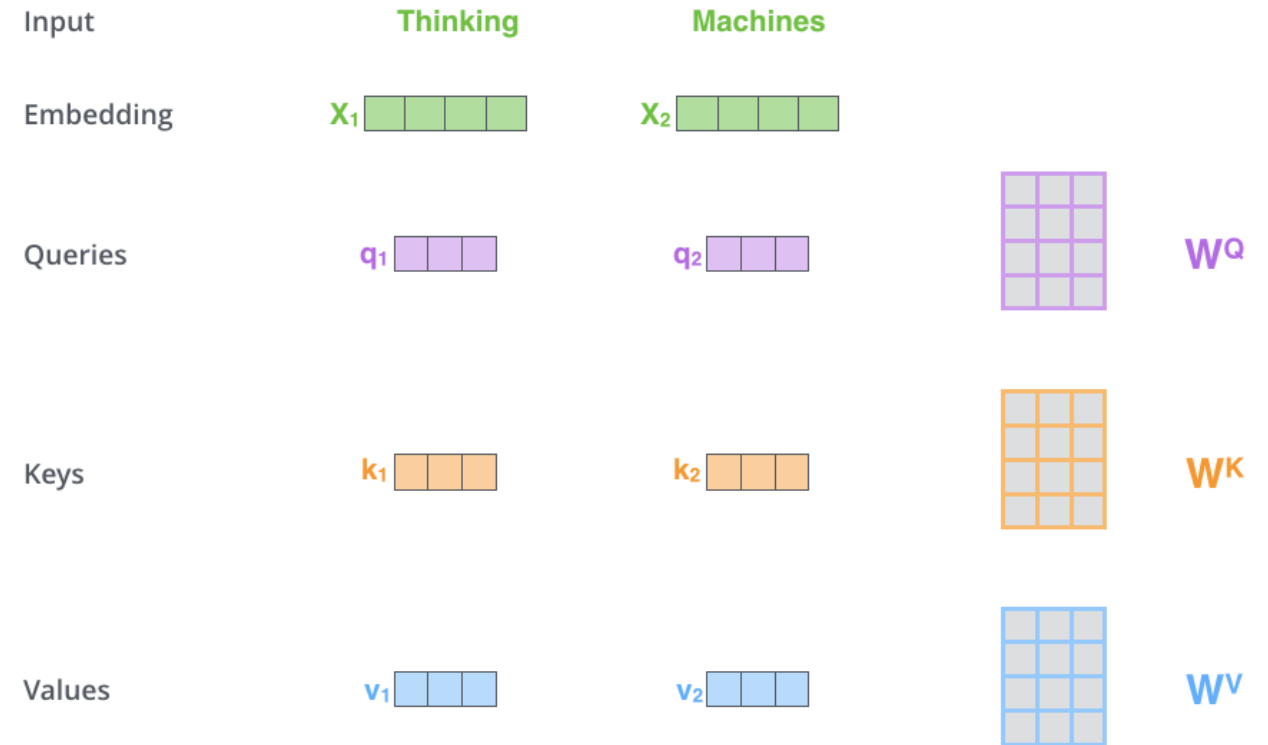
# Self-Attention at a high level

- Translating: “The animal didn't cross the street because **it** was too tired”
- What does “**it**” in this sentence refer to?
- When the model is processing the word “it”, self-attention allows it to associate “it” with “animal”.



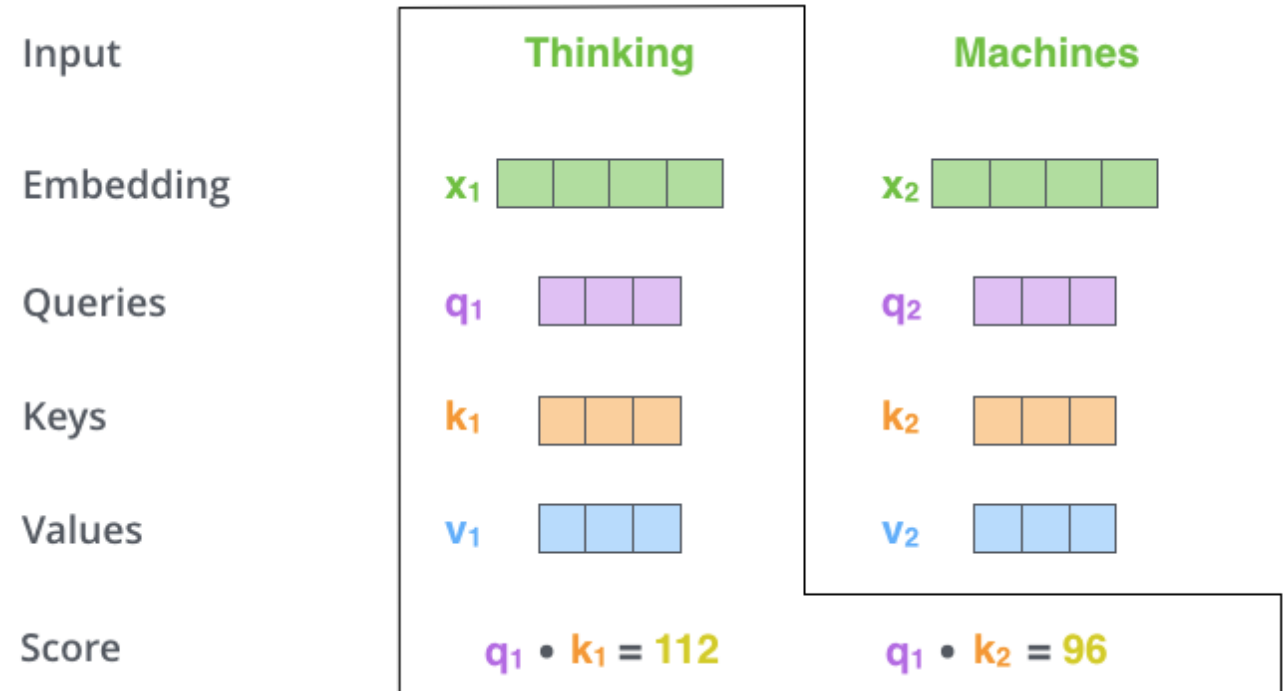
# Self-attention

- Q, K , V are created for each of the encoder's input vector (i.e., embedding of each word)
- These vectors are created by multiplying the embedding by three matrices  $W^Q$ ,  $W^K$  and  $W^V$  – trained during the training process



# Self-attention

- Calculate a score
  - calculated by taking the dot product of the query vector with the key vector of the respective word we're scoring.
  - So if processing the self-attention for the word in position #1, the first score would be the dot product of  $q_1$  and  $k_1$ . The second score would be the dot product of  $q_1$  and  $k_2$ .



# Self-attention

- Divide the score by the square root of the dimension of the key vectors
- Take softmax operation – determines how much each word will be expressed at this position

Input

Embedding

Queries

Keys

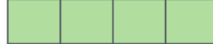
Values

Score

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

Softmax

Thinking

$x_1$  

$q_1$  

$k_1$  

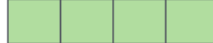
$v_1$  

$q_1 \cdot k_1 = 112$

14

0.88

Machines

$x_2$  

$q_2$  

$k_2$  

$v_2$  

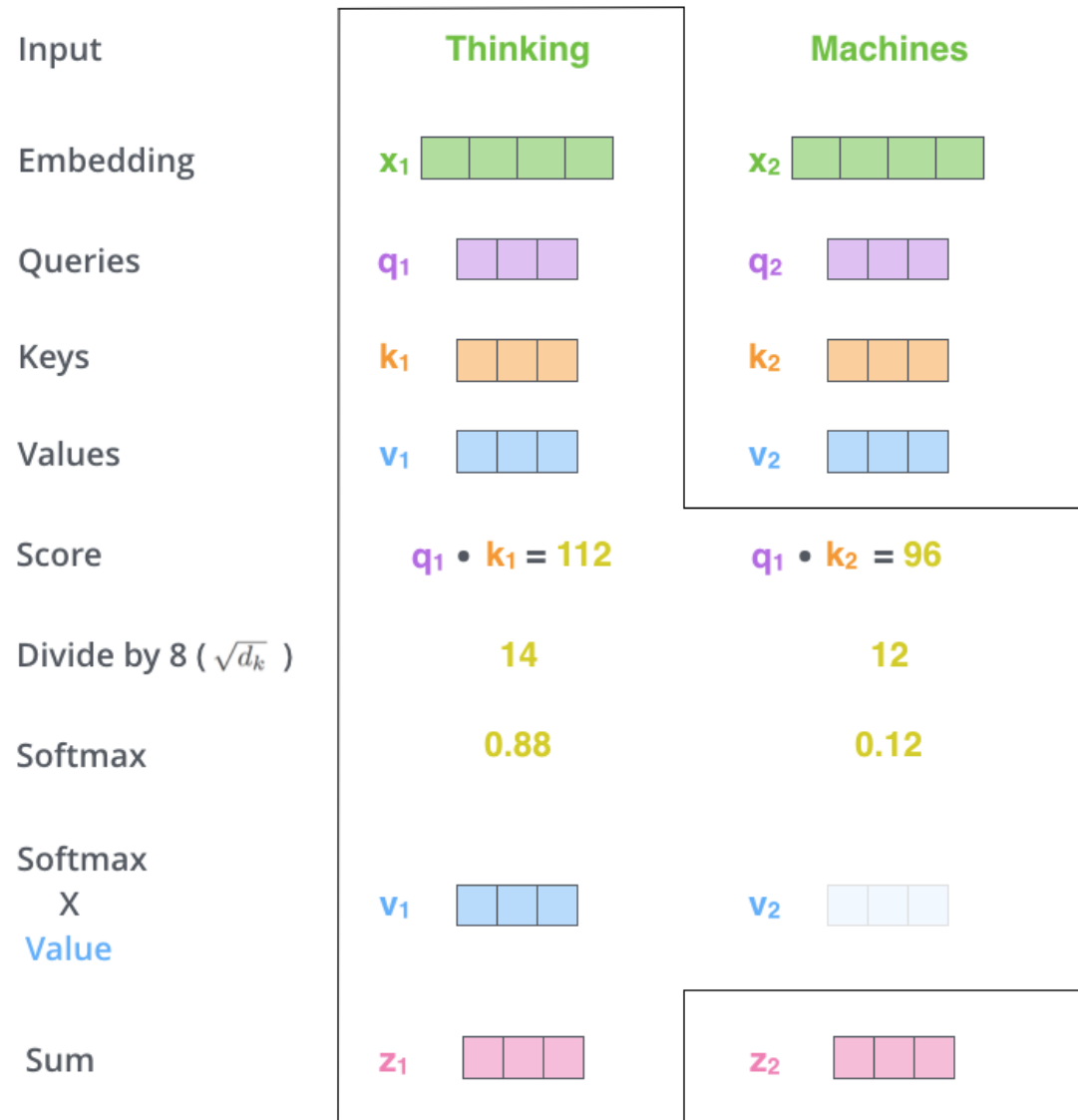
$q_1 \cdot k_2 = 96$

12

0.12

# Self-attention

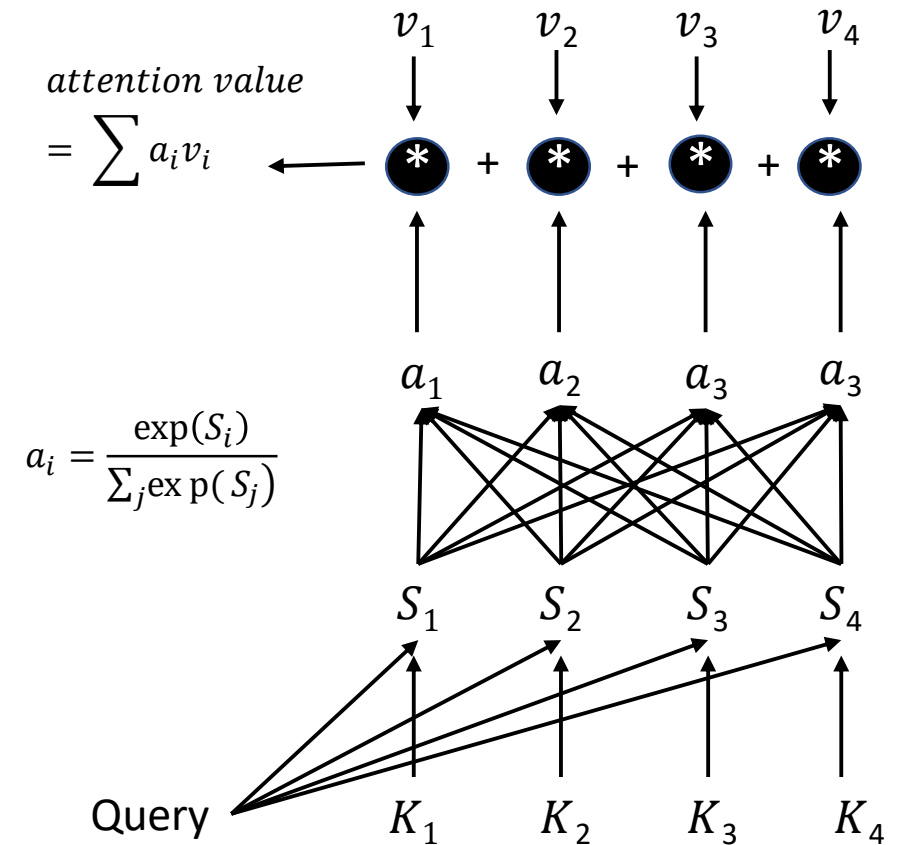
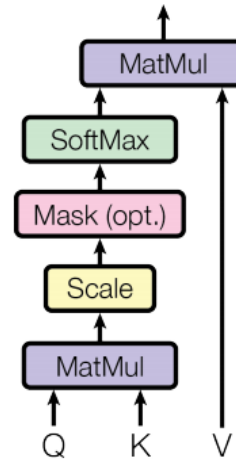
- Multiply each value by the softmax score.
- Sum up the weighted value of vectors – this is attention value.



# Attention mechanism

- Mimics the retrieval of a **value**  $v_i$  for a **query**  $q_i$  based on a **key**  $k_i$  in database.

$$\text{attention}(q, \mathbf{k}, \mathbf{v}) = \sum_i \text{similarity}(q, k_i) \times v_i$$



$$S_i = F(q, k_i) = \begin{cases} q^T k_i, & \text{dot product similarity} \\ \frac{q^T k_i}{\sqrt{d}}, & \text{scaled dot product} \\ q^T W k_i, & \text{general dot product} \end{cases}$$



# Multihead attention

- Compute multiple attentions per query with different weights

$$\text{multihead}(Q, K, V) = W^O \text{concat}(\text{head}_1, \text{head}_2, \dots, \text{head}_h)$$

$$\text{head}_i = \text{attention}(W_i^Q Q, W_i^K K, W_i^V V)$$

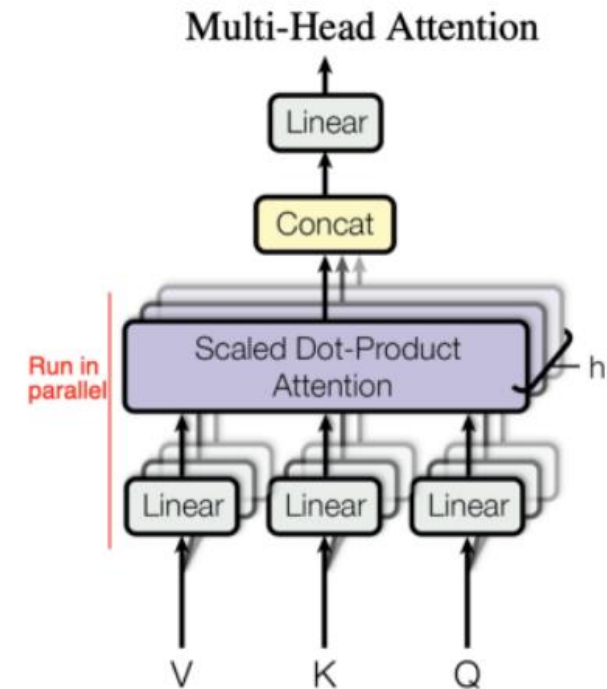
$$\text{attention}(Q, K, V) = \text{softmax}\left(\frac{Q^T K}{\sqrt{d_k}}\right) V$$

- Masked multi-head attention: multi-head where some values are masked (i.e., probabilities of masked values are nullified to prevent them from being selected).
- When decoding, an output value should only depend on previous outputs (not future outputs). Hence we mask future outputs.

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

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

where  $M$  is a mask matrix of 0's and  $-\infty$ 's

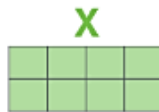


# Multihead attention

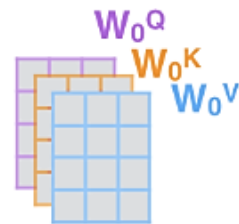
1) This is our input sentence\*

Thinking  
Machines

2) We embed each word\*



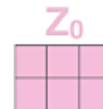
3) Split into 8 heads.  
We multiply  $X$  or  $R$  with weight matrices



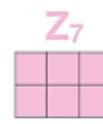
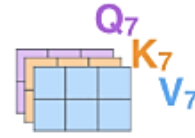
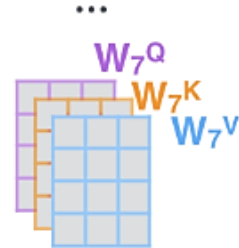
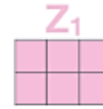
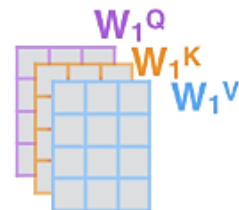
4) Calculate attention using the resulting  $Q/K/V$  matrices



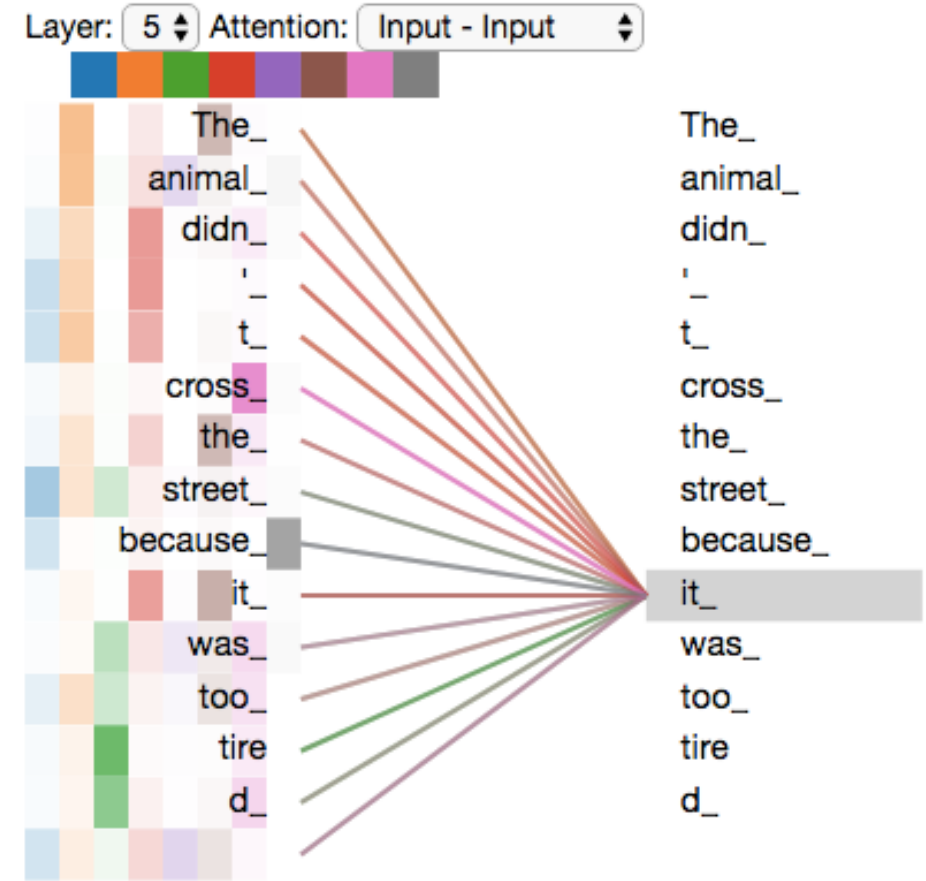
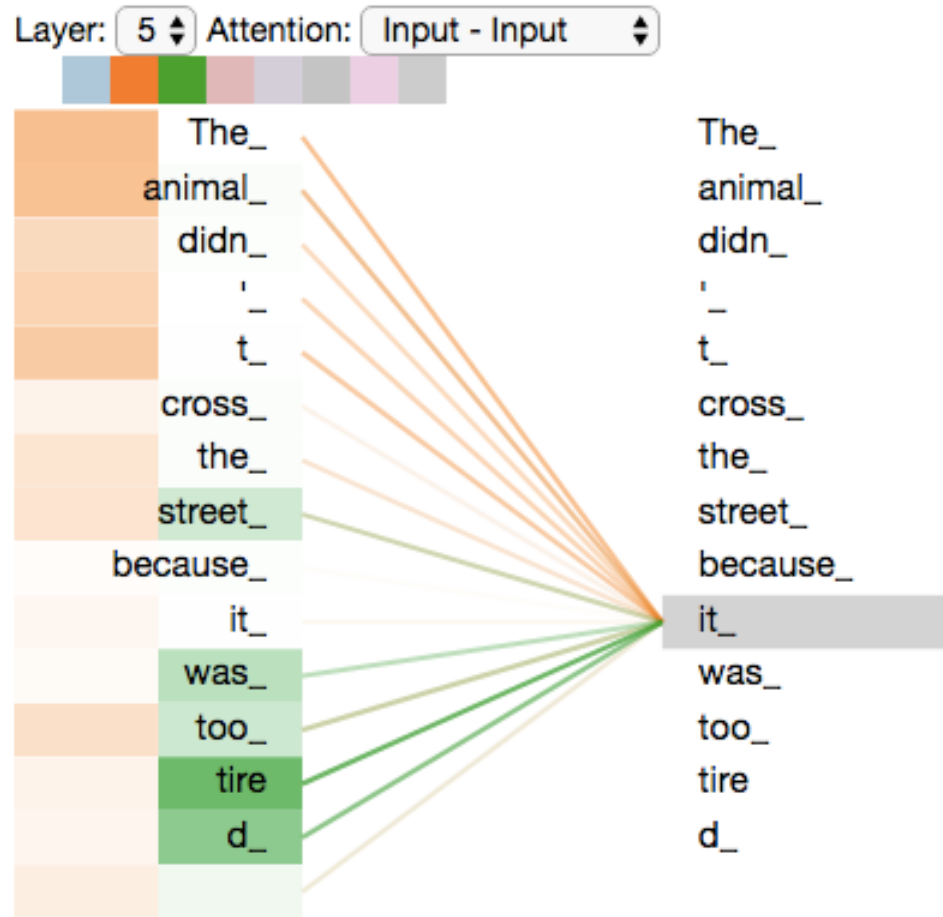
5) Concatenate the resulting  $Z$  matrices, then multiply with weight matrix  $W^O$  to produce the output of the layer



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



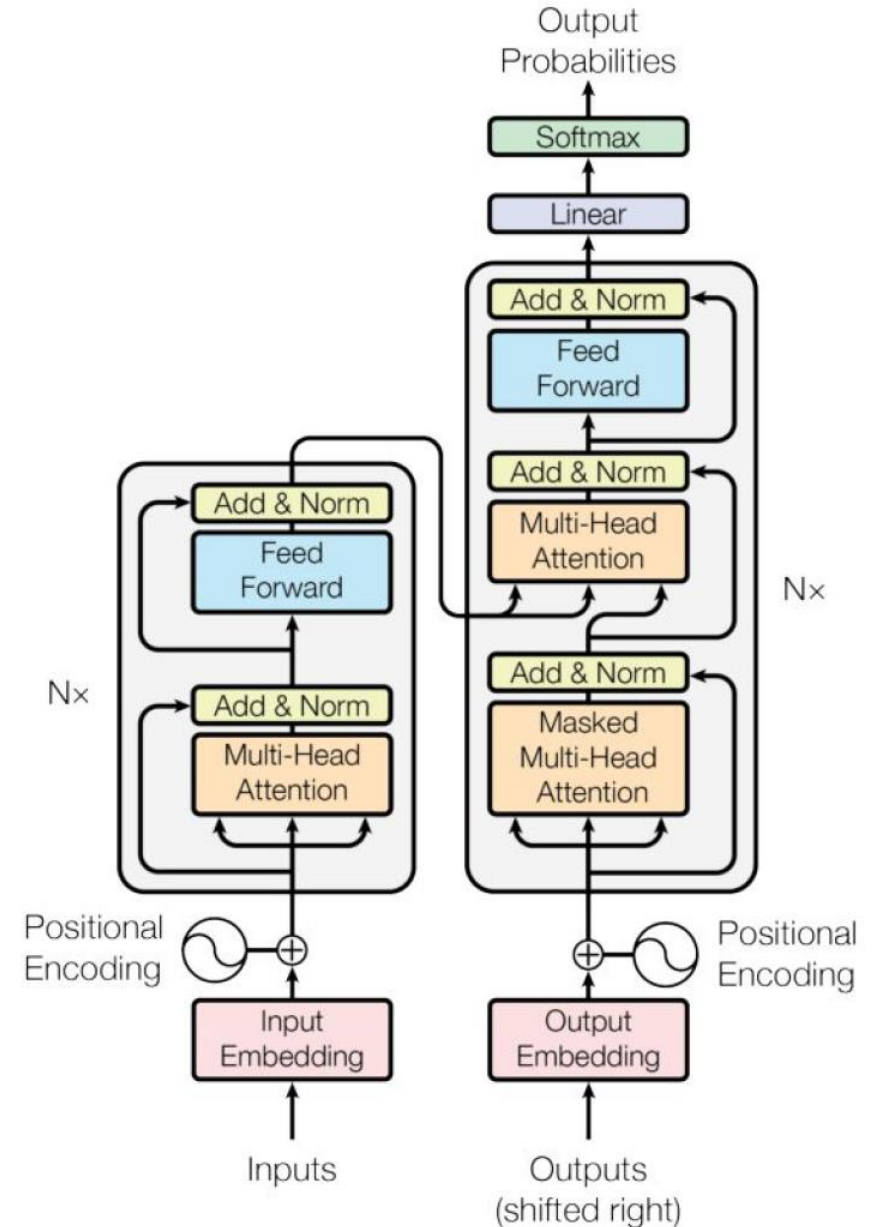
# Multihead attention



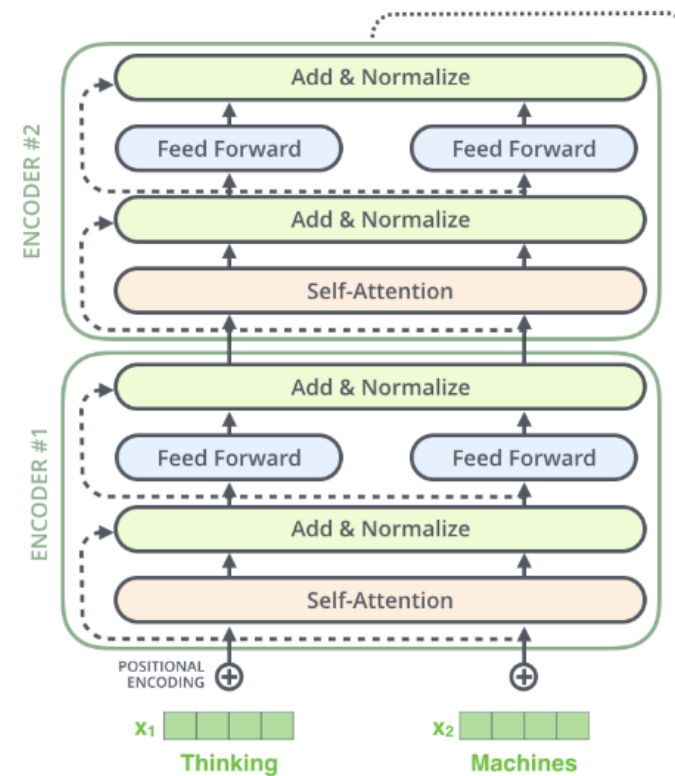
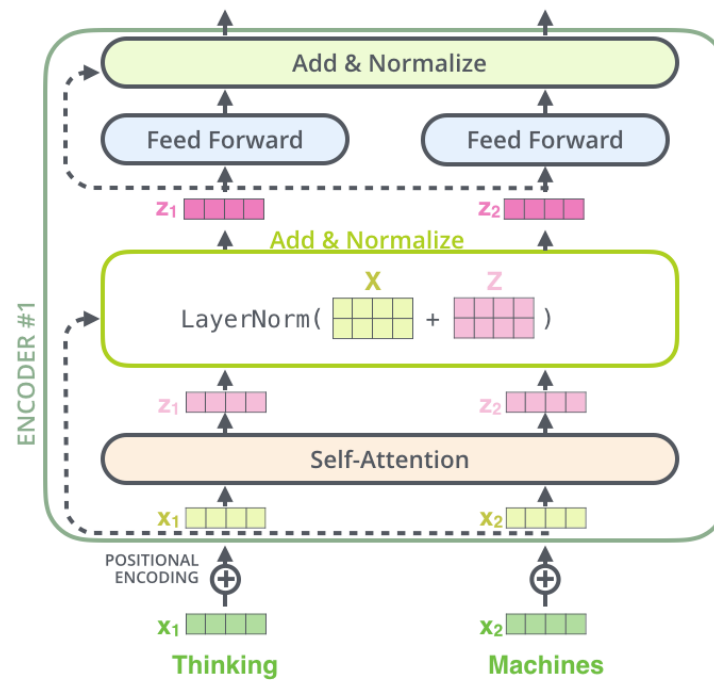
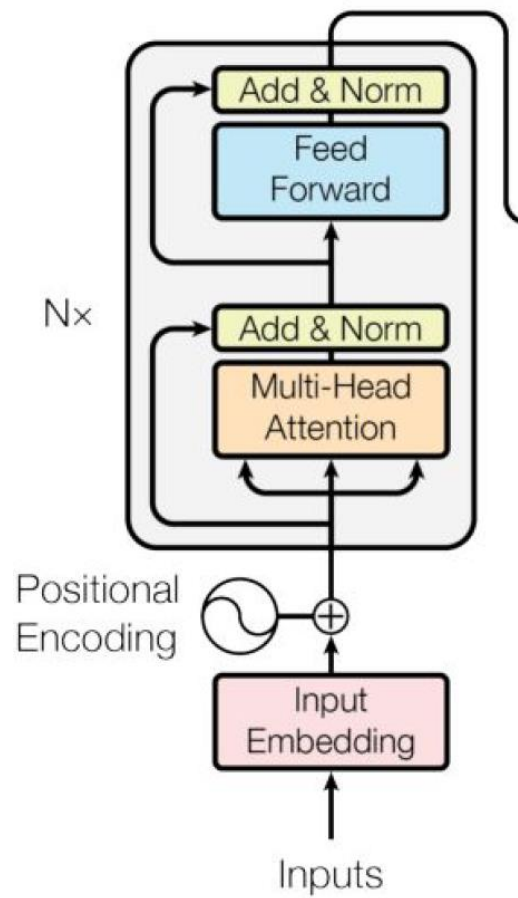
# Transformer architecture

- Encoder-decoder based on attention (no recurrence)
- Positional encoding:
  - use sine and cosine functions of different frequencies to encode the position information

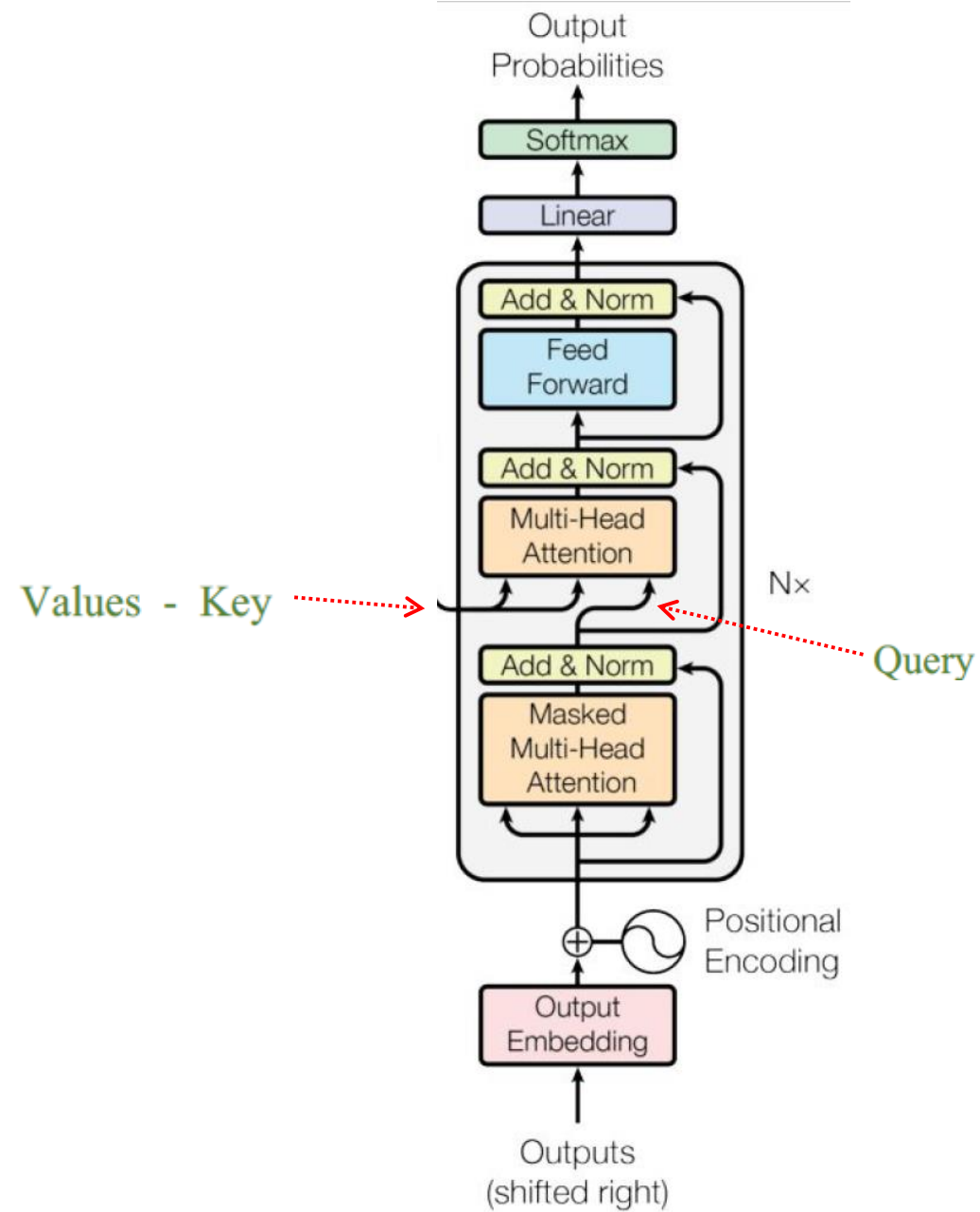
$$PE_{(pos, 2i)} = \sin(pos/10000^{2i/d_{model}})$$
$$PE_{(pos, 2i+1)} = \cos(pos/10000^{2i/d_{model}})$$



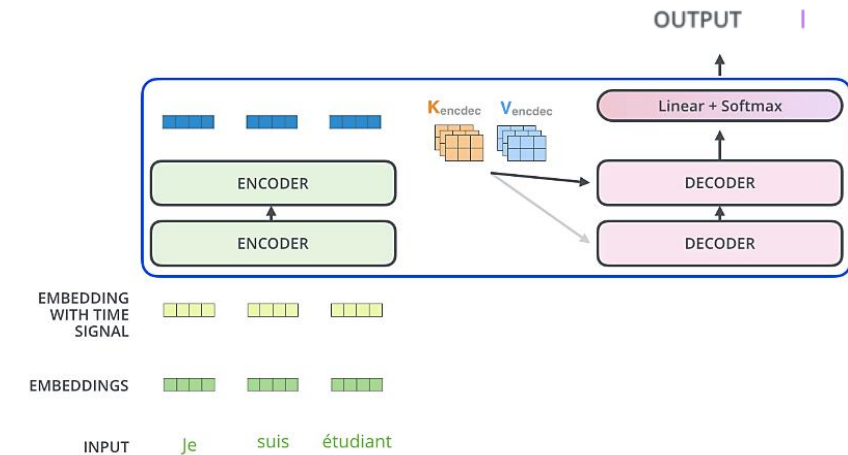
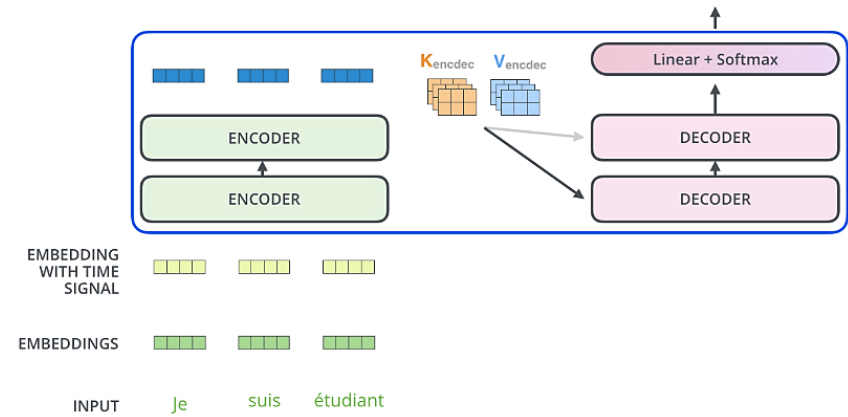
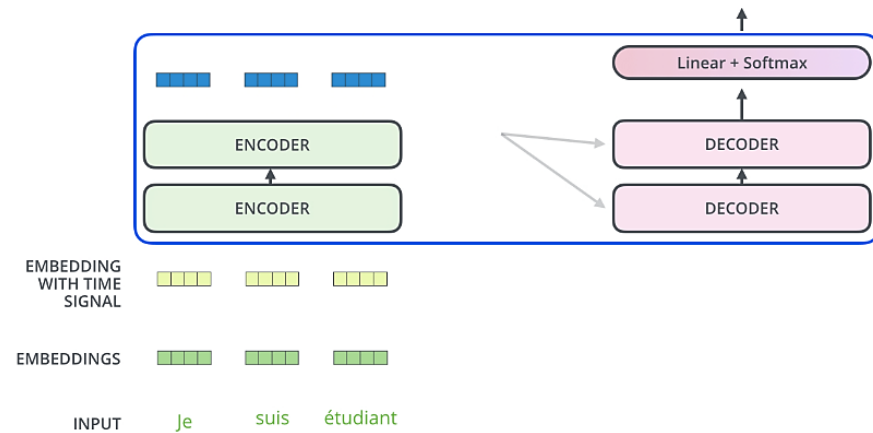
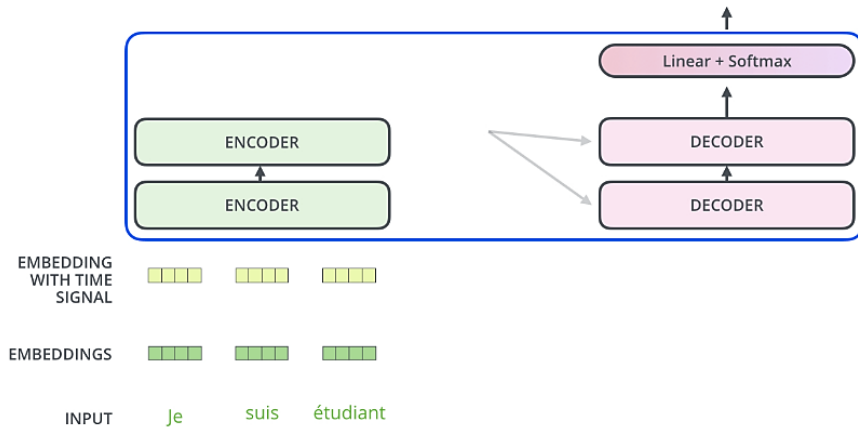
# Encoder



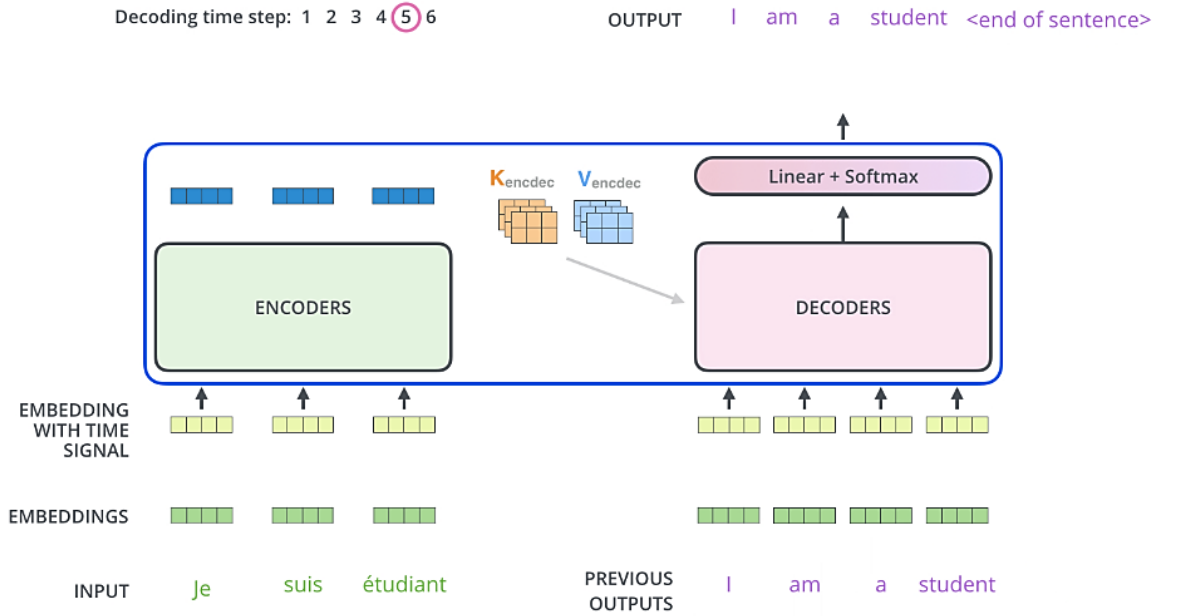
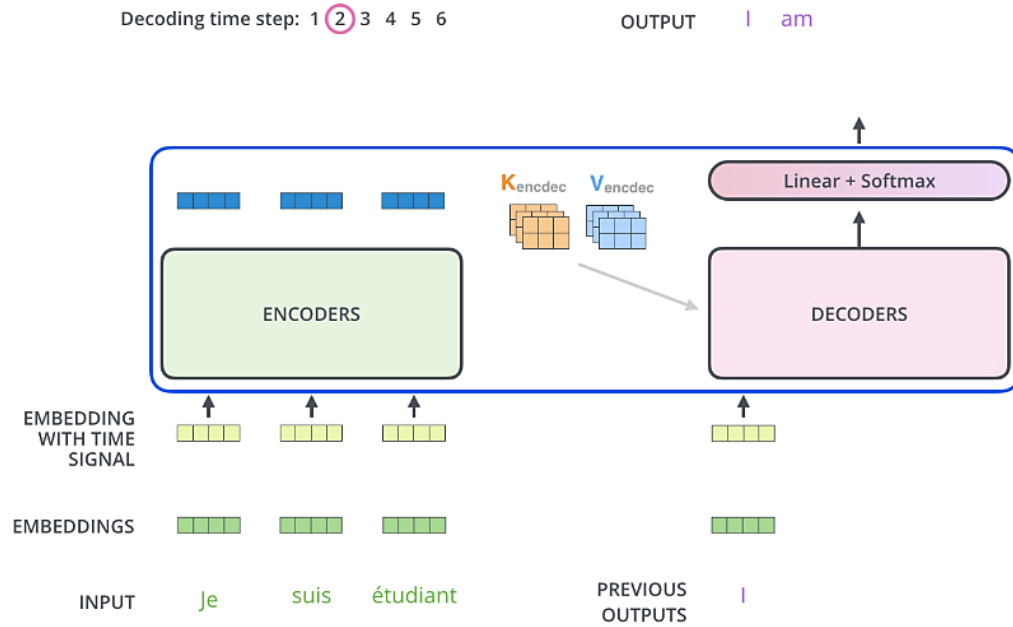
# Decoder



# Decoder

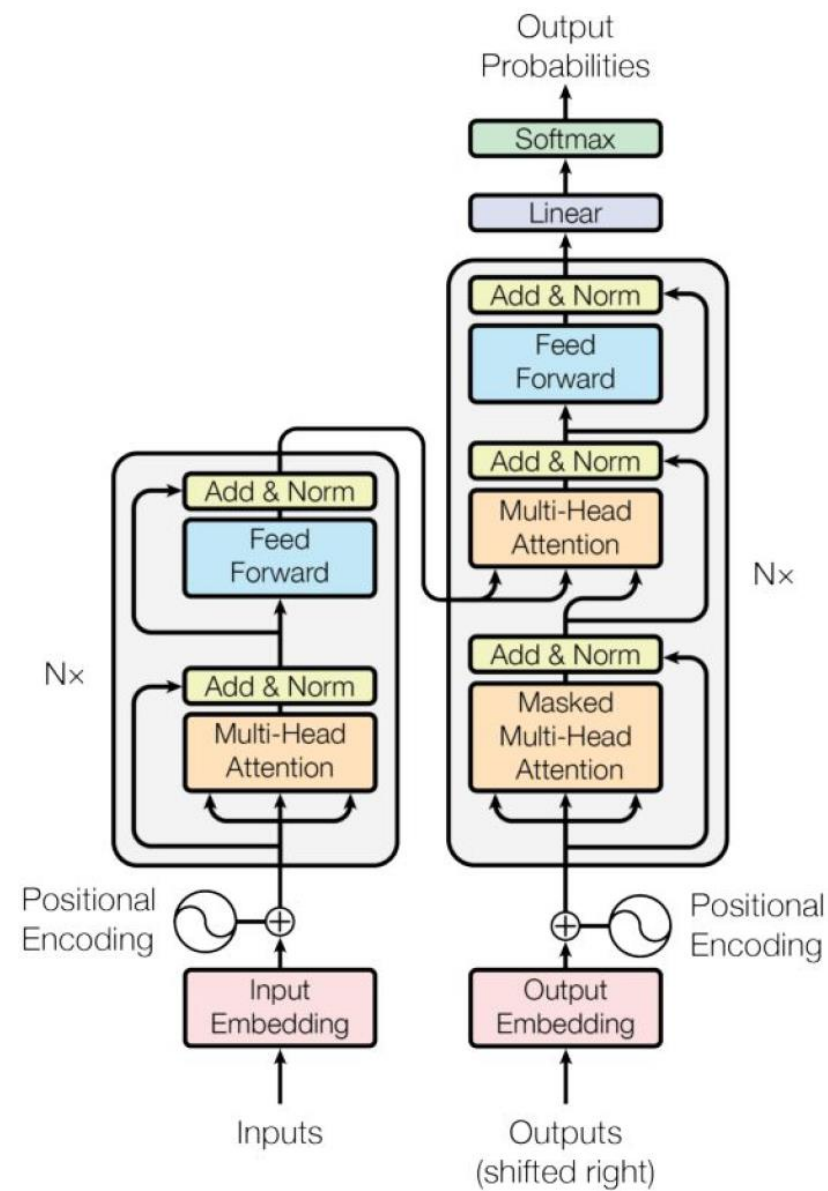


# Decoder





# Transformer



# Results

Table 1: Maximum path lengths, per-layer complexity and minimum number of sequential operations for different layer types.  $n$  is the sequence length,  $d$  is the representation dimension,  $k$  is the kernel size of convolutions and  $r$  the size of the neighborhood in restricted self-attention.

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

# Results

Table 2: The Transformer achieves better BLEU scores than previous state-of-the-art models on the English-to-German and English-to-French newstest2014 tests at a fraction of the training cost.

Model	BLEU		Training Cost (FLOPs)	
	EN-DE	EN-FR	EN-DE	EN-FR
ByteNet [15]	23.75			
Deep-Att + PosUnk [32]		39.2		$1.0 \cdot 10^{20}$
GNMT + RL [31]	24.6	39.92	$2.3 \cdot 10^{19}$	$1.4 \cdot 10^{20}$
ConvS2S [8]	25.16	40.46	$9.6 \cdot 10^{18}$	$1.5 \cdot 10^{20}$
MoE [26]	26.03	40.56	$2.0 \cdot 10^{19}$	$1.2 \cdot 10^{20}$
Deep-Att + PosUnk Ensemble [32]		40.4		$8.0 \cdot 10^{20}$
GNMT + RL Ensemble [31]	26.30	41.16	$1.8 \cdot 10^{20}$	$1.1 \cdot 10^{21}$
ConvS2S Ensemble [8]	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>41.0</b>	$2.3 \cdot 10^{19}$	

# Variants of Transformers

- BERT
- XLNet
- RoBERTa
- GPT-2
- GPT-3
- ALBERT

Thank you for your attention!!