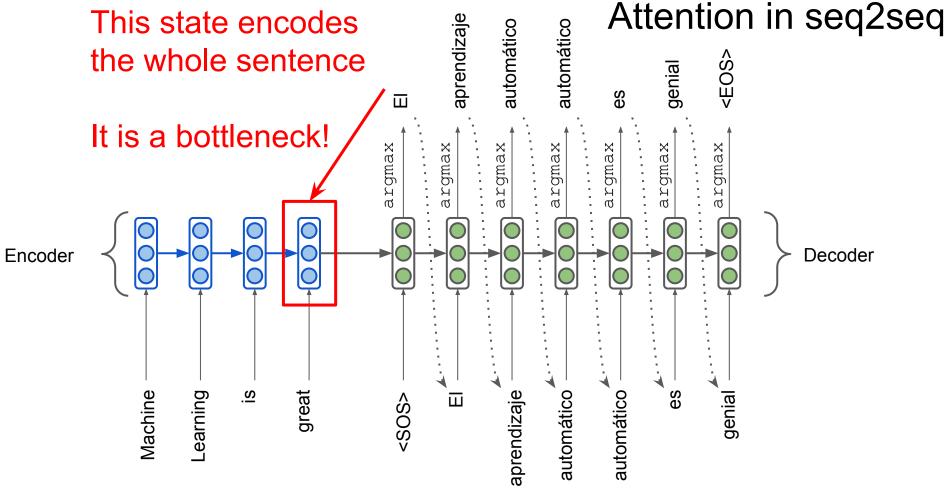
Lecture 05: Self-Attention & Transformer overview

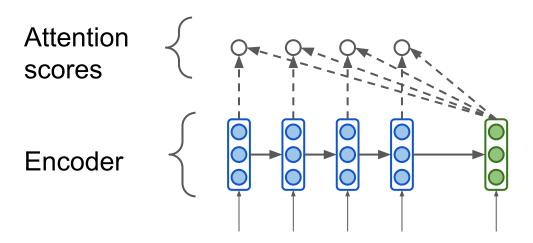
Radoslav Neychev

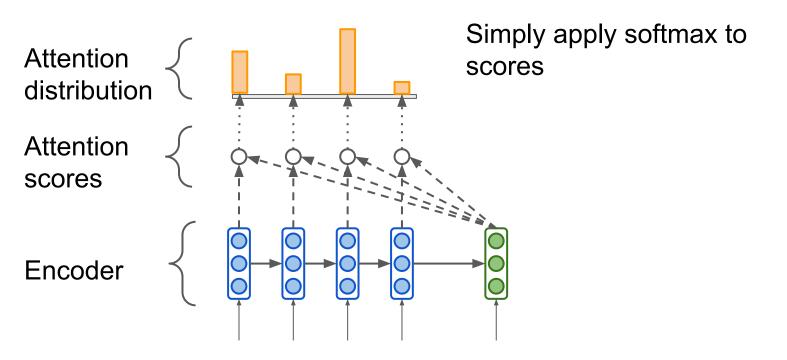
MADE, Moscow 07.04.2021

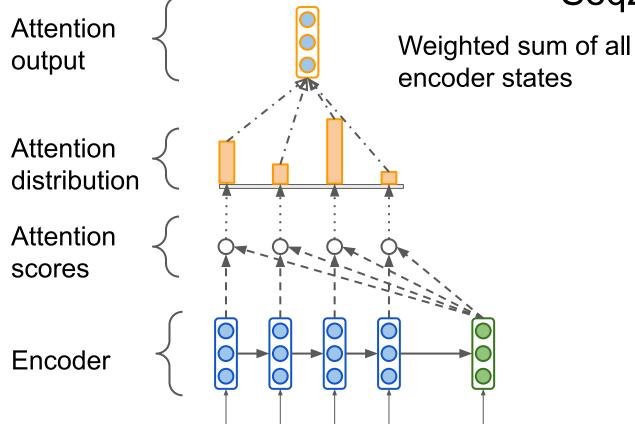
Outline

- 1. recap: Attention in seq2seq
- 2. Transformer architecture
- 3. Self-Attention
- 4. Positional encoding (optional)
- 5. Layer normalization (optional)
- 6. Decoder in Transformer (optional)



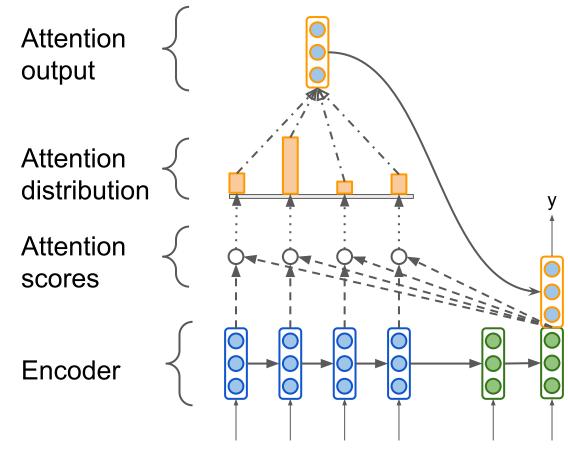


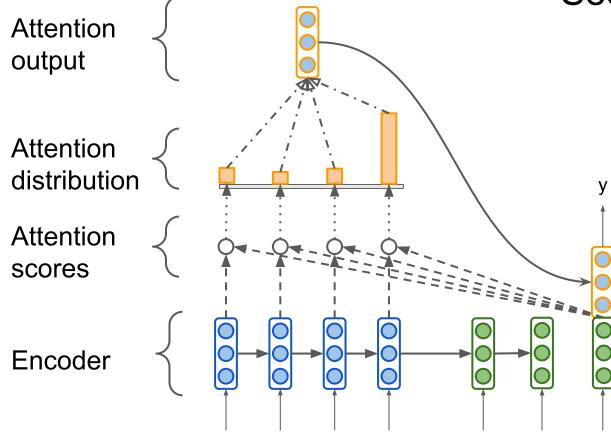




Attention output **Attention** Concatenate distribution **Attention** scores Encoder

Attention output **Attention** distribution Attention scores Encoder





Attention in equations

Denote encoder hidden states $\mathbf{h}_1,\dots,\mathbf{h}_N\in\mathbb{R}^k$ and decoder hidden state at time step t $\mathbf{s}_t\in\mathbb{R}^k$

The attention scores \mathbf{e}^t can be computed as dot product

$$\mathbf{e}^t = [\mathbf{s}^T \mathbf{h}_1, \dots, \mathbf{s}^T \mathbf{h}_N]$$

Then the attention vector is a linear combination of encoder states

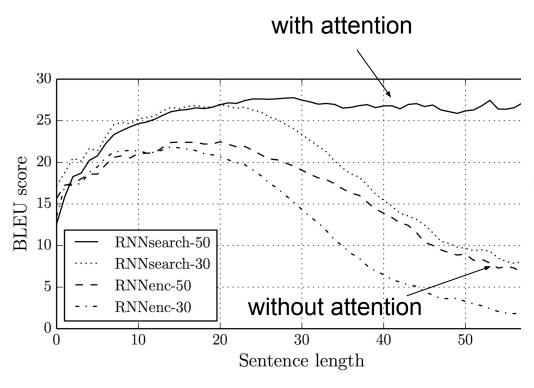
$$\mathbf{a}_t = \sum_{i=1}^N oldsymbol{lpha}_i^t \mathbf{h}_i \in \mathbb{R}^k$$
 , where $oldsymbol{lpha}_t = \operatorname{softmax}(\mathbf{e}_t)$

Attention variants

- Basic dot-product (the one discussed before): $e_i = s^T h_i \in \mathbb{R}$
- Multiplicative attention: $e_i = s^T W h_i \in \mathbb{R}$
 - \bigcirc $W \in \mathbb{R}^{d_2 \times d_1}$ weight matrix
- Additive attention: $e_i = v^T \tanh(W_1 h_i + W_2 s) \in \mathbb{R}$
 - \circ $extbf{W}_1 \in \mathbb{R}^{d_3 imes d_1}, extbf{W}_2 \in \mathbb{R}^{d_3 imes d_2}$ weight matrices
 - \circ $v \in \mathbb{R}^{d_3}$ weight vector

Attention advantages

- "Free" word alignment
- Better results on long sequences



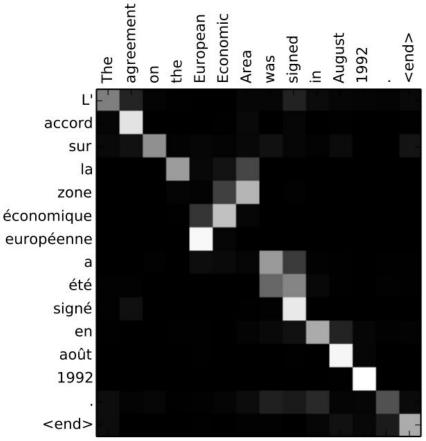
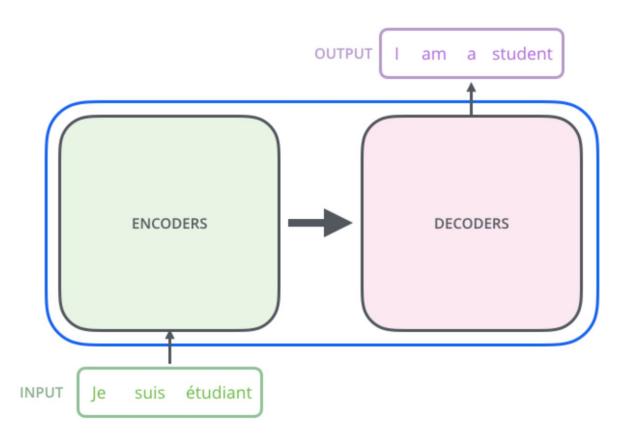
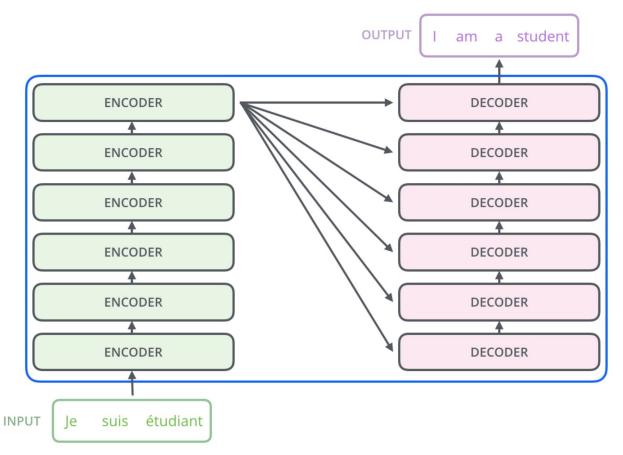


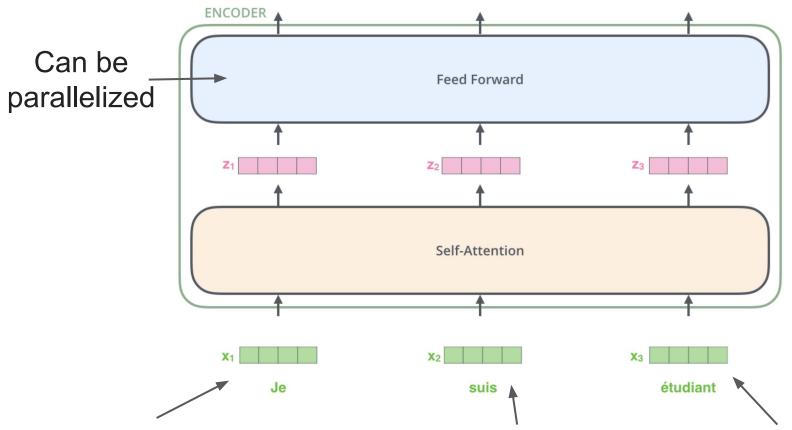
Image source: Neural Machine Translation by Jointly Learning to Align and Translate





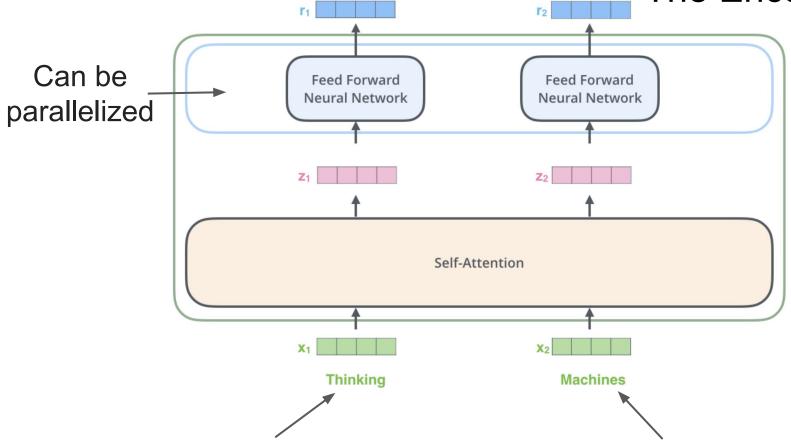


The Encoder Side



the word in each position flows through its own path in the encoder 18

The Encoder Side



the word in each position flows through its own path in the encoder 19

The Transformer: quick overview

- Proposed in 2017 in paper <u>Attention is All You Need</u> by Ashish Vaswani et al.
- No recurrent or convolutional layers, only attention
- Beats seq2seq in machine translation task
 - 28.4 BLEU on the WMT 2014 English-to-German translation task
- Much faster
- Uses **self-attention** concept

Self-Attention

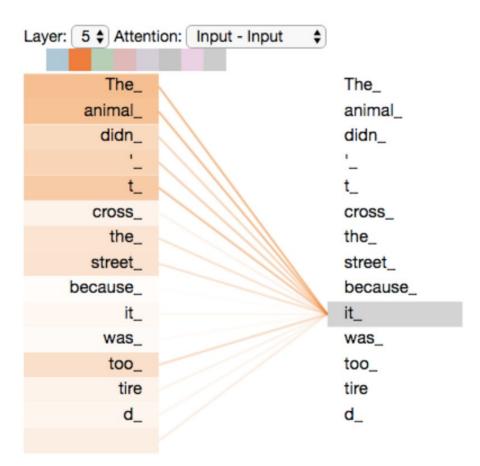
Self-Attention at a High Level

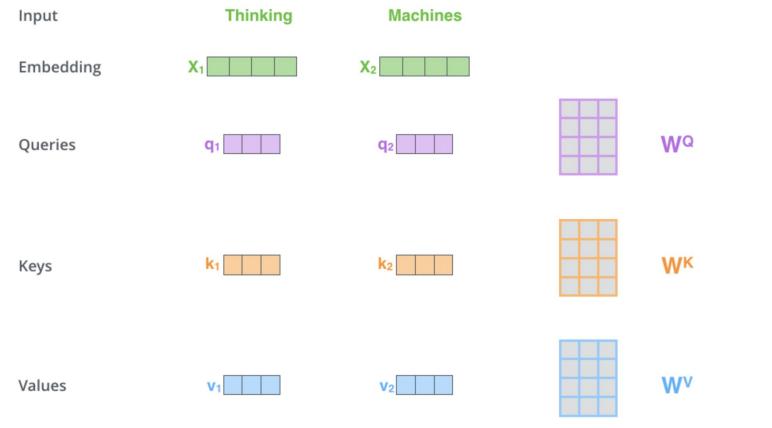
"The animal didn't cross the street because it was too tired"

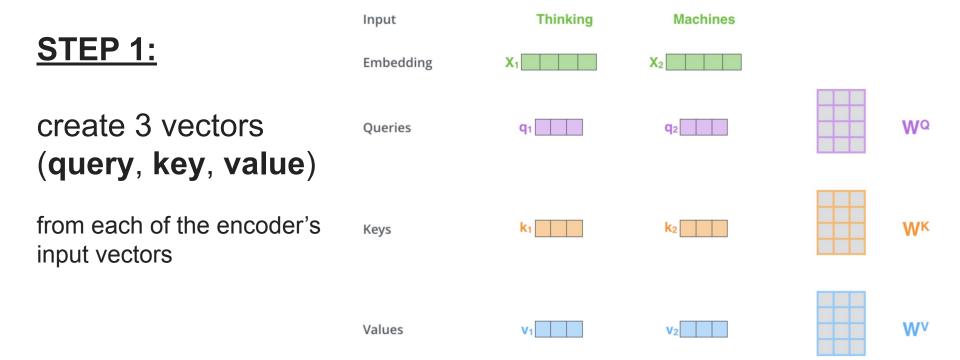
- What does "it" in this sentence refer to?
- We want self-attention to associate "it" with "animal"

 Self-attention is the method the Transformer uses to bake the "understanding" of other relevant words into the one we're currently processing

Self-Attention at a High Level







What are the query, key, value vectors?

They're abstractions that are useful for calculating and thinking about attention.

STEP 2:

calculate a score

(score each word of the input sentence against the current word) Input

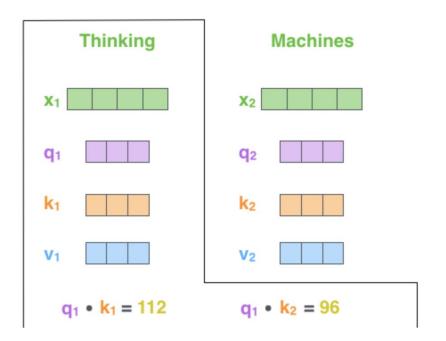
Embedding

Queries

Keys

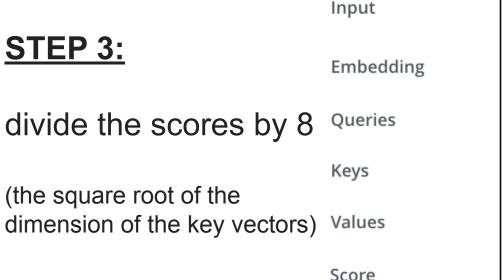
Values

Score



Thinking

14



Q1 k_1 V₁ $q_1 \cdot k_1 = 112$

Machines V2 $q_1 \cdot k_2 = 96$ 12

STEP 4:

softmax

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

Softmax

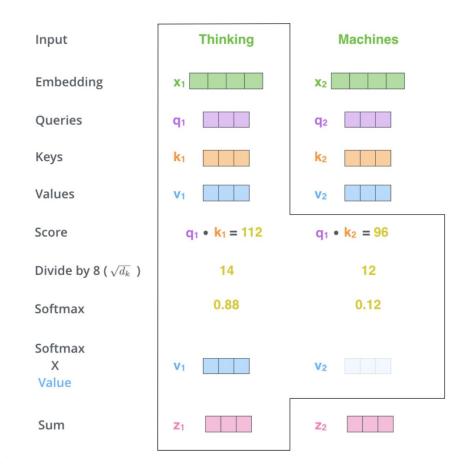
0.88 0.12

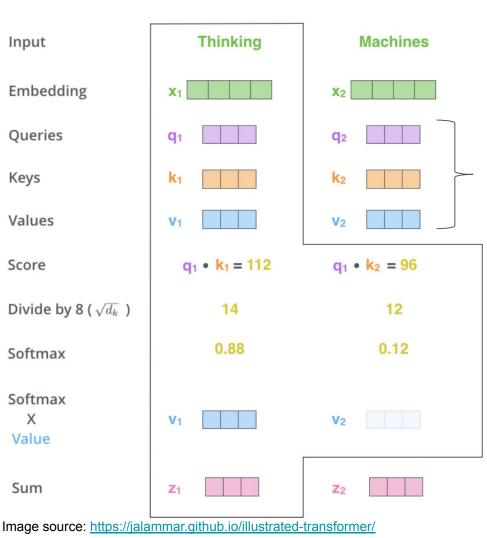
STEP 5:

multiply each value vector by the softmax score

STEP 6:

sum up the weighted value vectors





Self-Attention

STEP 1: create Query, Key, Value

STEP 3: divide by $\sqrt{d_k}$

STEP 2: calculate scores

STEP 4: softmax

STEP 5: multiply each value vector by the softmax score

STEP 6: sum up the weighted value vectors

Self-Attention: Matrix Calculation

Pack embeddings into matrix **X**

Multiply X by weight matrices we've trained (Wk, Wq, Wv)

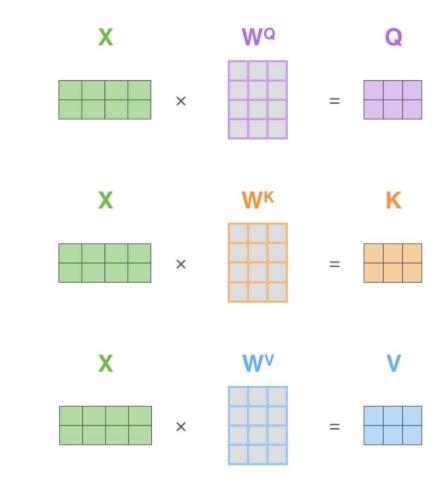
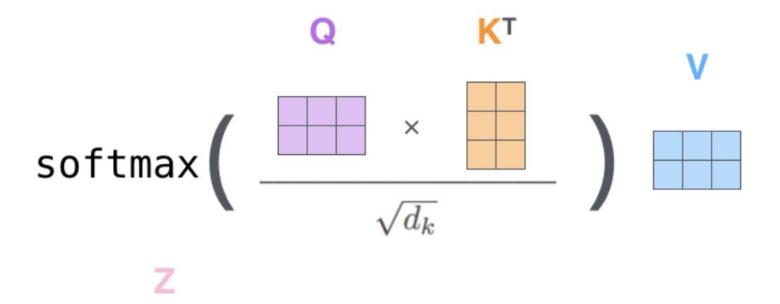


Image source: https://jalammar.github.io/illustrated-transformer/

Self-Attention: Matrix Calculation



Multi-Head Attention

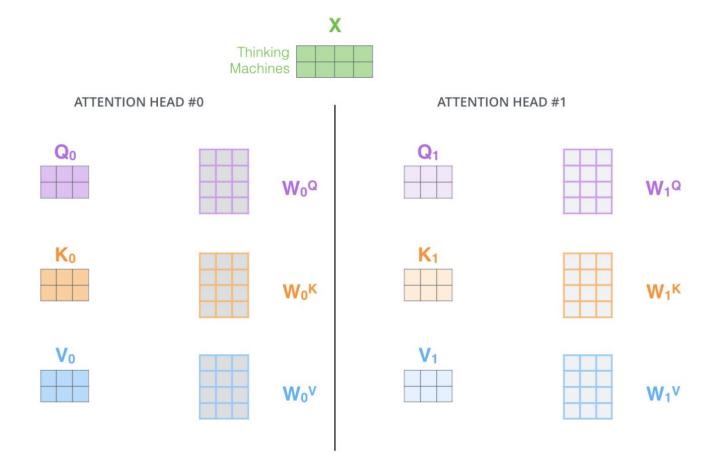


Image source: https://jalammar.github.io/illustrated-transformer/

Multi-Head Attention

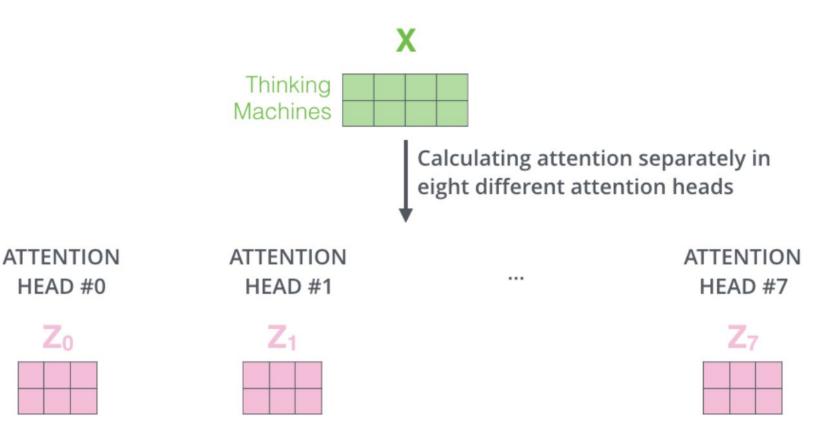


Image source: https://jalammar.github.io/illustrated-transformer/

Multi-Head Attention

1) Concatenate all the attention heads

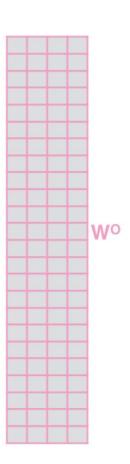


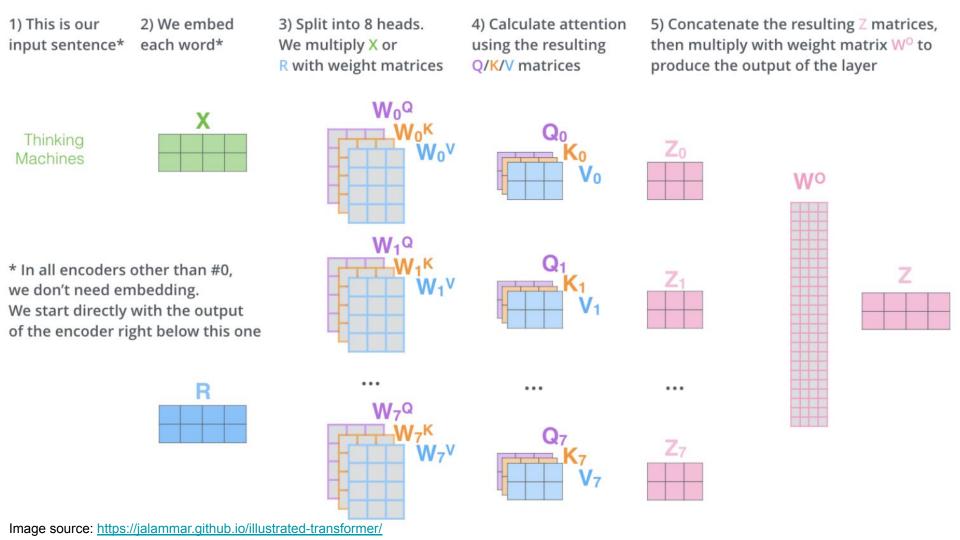
2) Multiply with a weight matrix W° that was trained jointly with the model

Χ

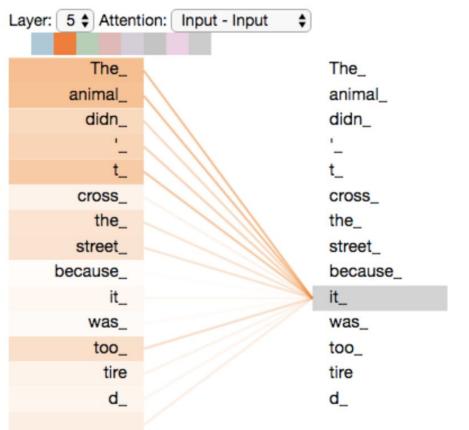
3) The result would be the Z matrix that captures information from all the attention heads. We can send this forward to the FFNN

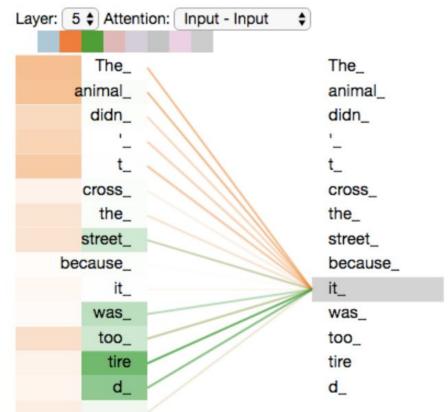




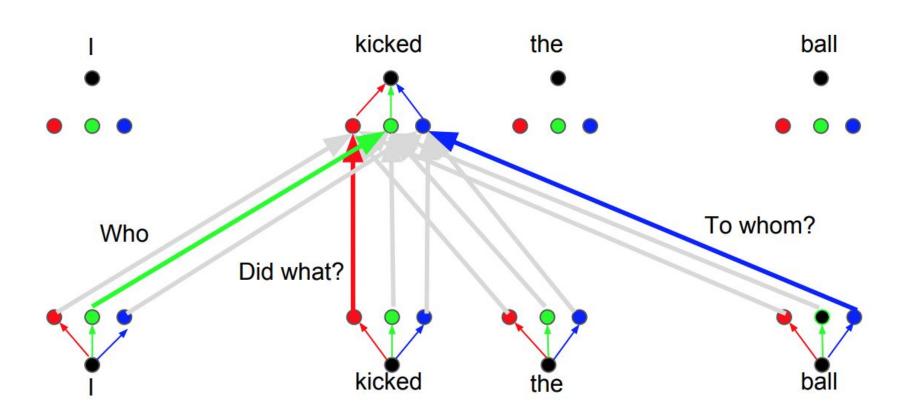


Multi-Head Attention

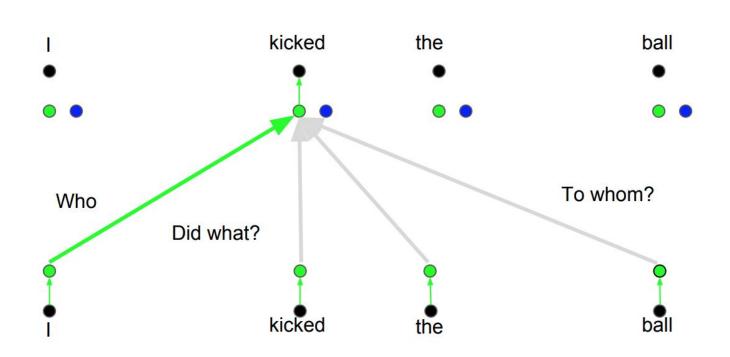




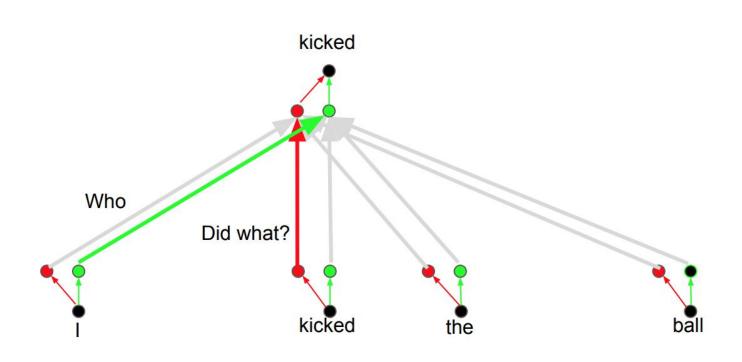
Why Multi-Head Attention?



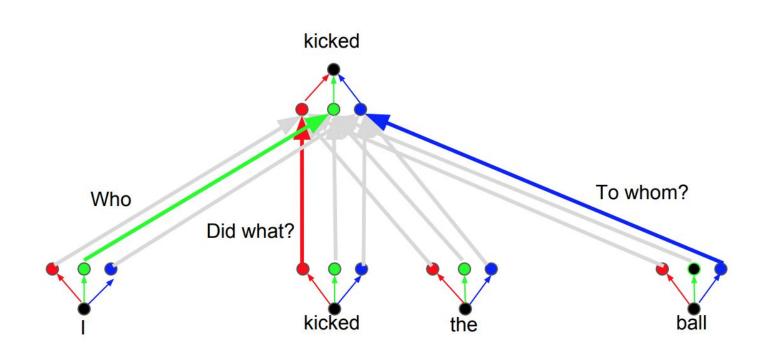
Attention head: Who



Attention head: Did What?

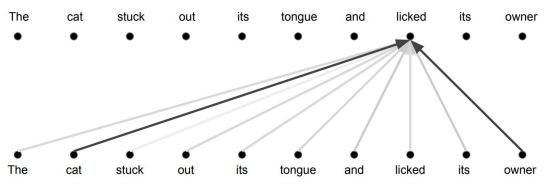


Attention head: To Whom?



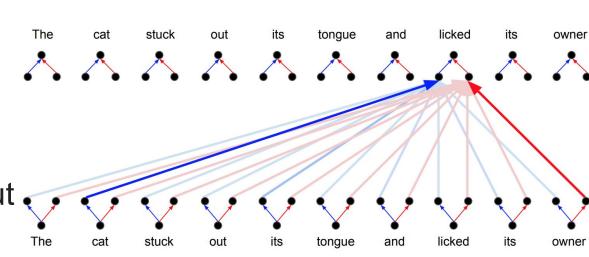
Attention vs. Multi-Head Attention

Attention: a weighted average



Multi-Head Attention:

parallel attention layers with different linear transformations on input and output.



Performance: WMT 2014 BLEU

	EN-DE	EN-FR
GNMT (orig)	24.6	39.9
ConvSeq2Seq	25.2	40.5
Transformer*	28.4	41.8

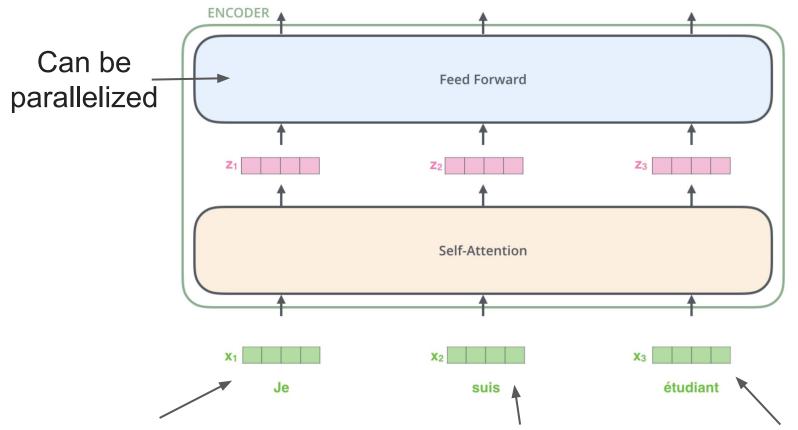
^{*}Transformer models trained >3x faster than the others.

Research Challenges

- Constant 'path length' between any two positions.
- Unbounded memory.
- Trivial to parallelize (per layer).
- Models Self-Similarity.
- Relative attention provides expressive timing, equivariance, and extends naturally to graphs.

Positional Encoding

The Encoder Side

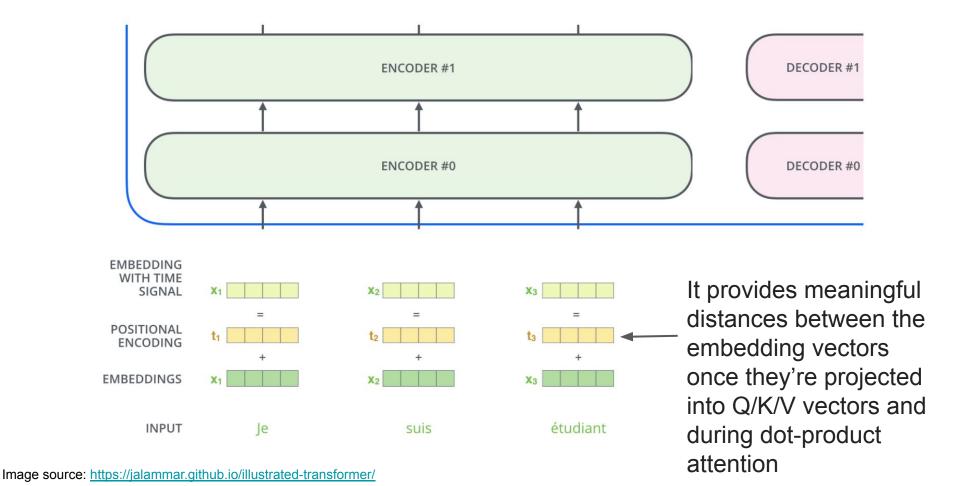


the word in each position flows through its own path in the encoder 46

Positional encoding requirements

- Positional encoding should be unique for every position in the sequence
- Distance between two same positions should be preserved with sequences of different length
- The positional encoding should be deterministic
- It would be great if it would work with long sequences (longer than any sequence in the training set)

Positional Encoding

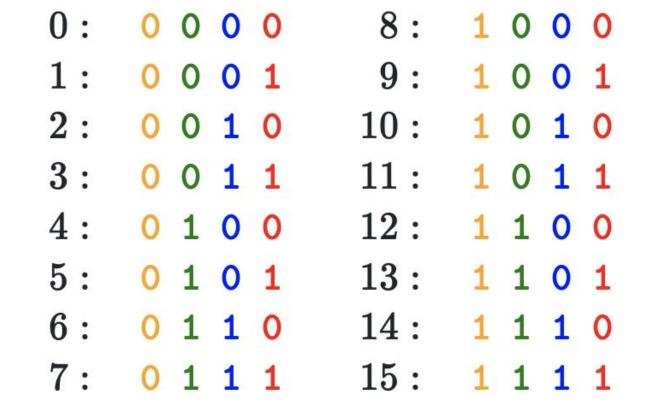


Positional Encoding: why sin and cos?

$$\vec{p_t}^{(i)} = f(t)^{(i)} = \begin{cases} \sin(\omega_k t), & \text{if } i = 2k \\ \cos(\omega_k t), & \text{if } i = 2k + 1 \end{cases}$$

$$\omega_k = \frac{1}{10000^{2k/d}} \qquad \vec{p_t} = \begin{cases} \sin(\omega_1 . t) \\ \cos(\omega_1 . t) \\ \sin(\omega_2 . t) \\ \cos(\omega_2 . t) \\ \vdots \\ \sin(\omega_{d/2} . t) \\ \cos(\omega_{d/2} . t) \\ \cos(\omega_{d/2} . t) \end{cases}$$
 t stays for position in the original sequence k is the index of the element in the positional vector

Positional Encoding



Positional Encoding

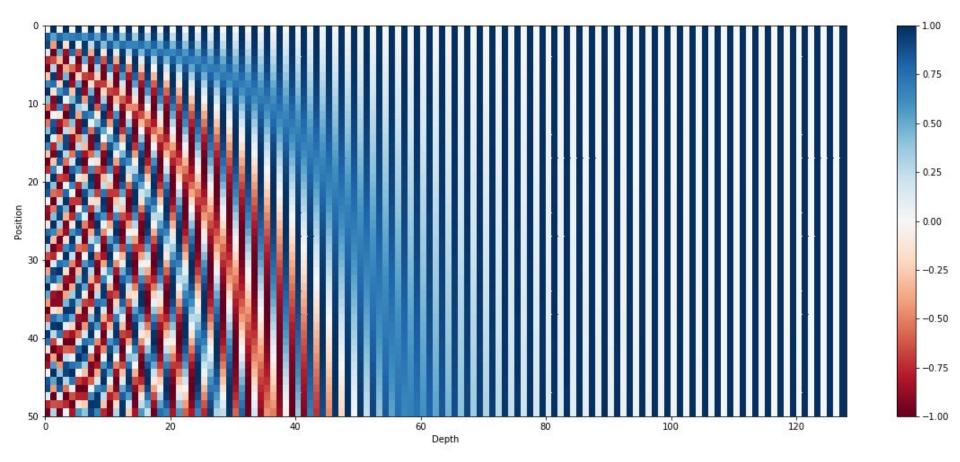


Image source: https://kazemnejad.com/blog/transformer_architecture_positional_encoding/

Positional Encoding: why sin and cos?

We chose this function because we hypothesized it would allow the model to easily learn to attend by relative positions, since for any fixed offset k, PEpos+k can be represented as a linear function of PEpos.

$$M \begin{bmatrix} \sin(\omega_k t) \\ \cos(\omega_k t) \end{bmatrix} = \begin{bmatrix} \sin(\omega_k (t + \phi)) \\ \cos(\omega_k (t + \phi)) \end{bmatrix}$$

Positional Encoding: why sin and cos?

$$\begin{bmatrix} u_1 & v_1 \\ u_2 & v_2 \end{bmatrix} \begin{bmatrix} \sin(\omega_k t) \\ \cos(\omega_k t) \end{bmatrix} = \begin{bmatrix} \sin(\omega_k (t + \phi)) \\ \cos(\omega_k (t + \phi)) \end{bmatrix}$$
$$\begin{bmatrix} u_1 & v_1 \\ u_2 & v_2 \end{bmatrix} \begin{bmatrix} \sin(\omega_k t) \\ \cos(\omega_k t) \end{bmatrix} = \begin{bmatrix} \sin(\omega_k t) \cos(\omega_k \phi) + \cos(\omega_k t) \sin(\omega_k \phi) \\ \cos(\omega_k t) \cos(\omega_k \phi) - \sin(\omega_k t) \sin(\omega_k \phi) \end{bmatrix}$$

$$M_{\phi,k} = \begin{bmatrix} \cos(\omega_k \phi) & \sin(\omega_k \phi) \\ -\sin(\omega_k \phi) & \cos(\omega_k \phi) \end{bmatrix}$$

Positional Encoding

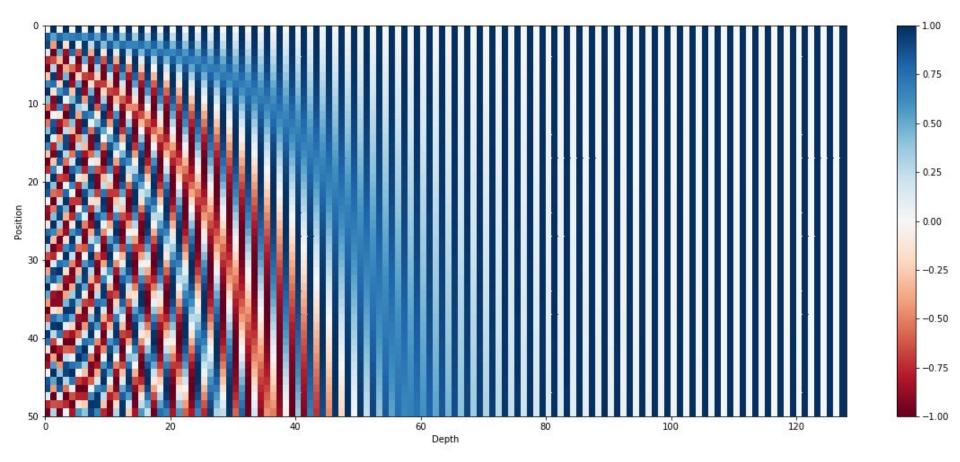
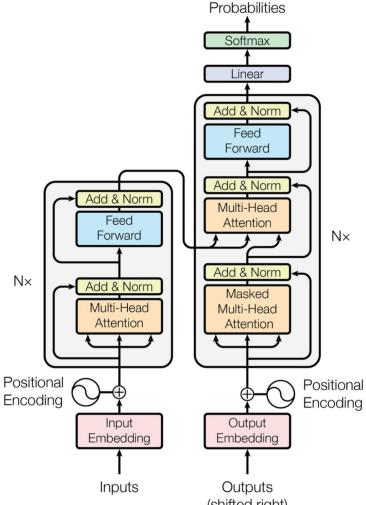


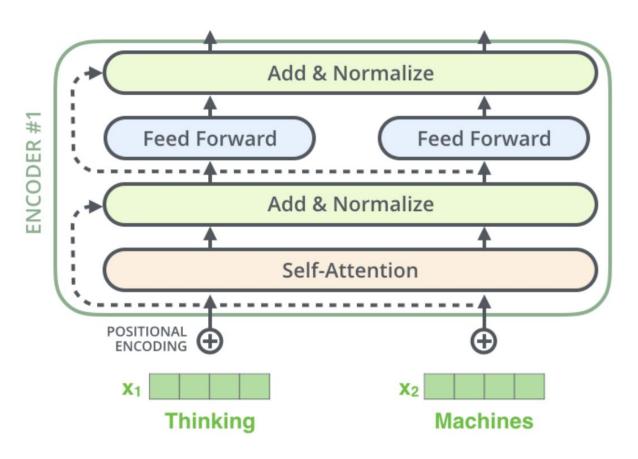
Image source: https://kazemnejad.com/blog/transformer_architecture_positional_encoding/



Output

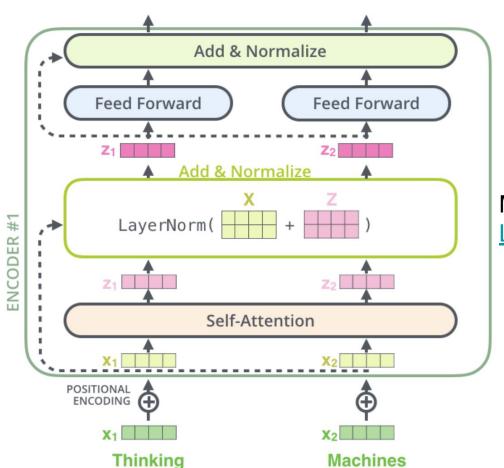
The Transformer: recap

56



Like BatchNorm

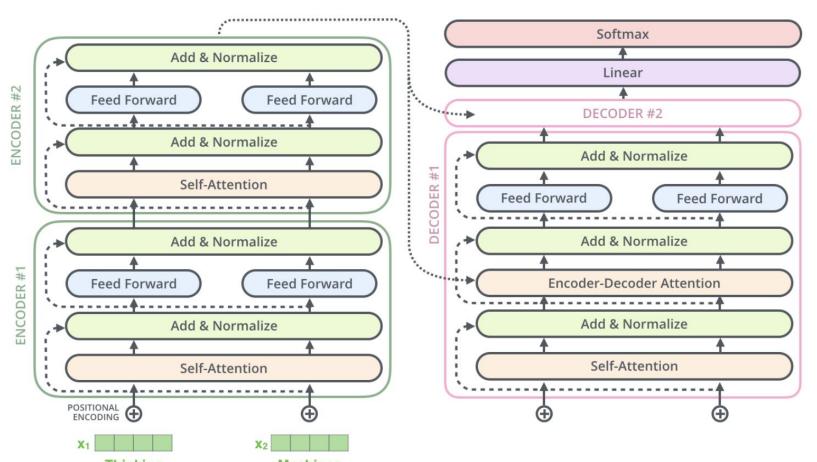
but normalize along all features representing latent vector



More info:

<u>Layer Normalization</u>

Image source: https://jalammar.github.io/illustrated-transformer/



Thinking Machines Image source: https://jalammar.github.fo/illustrated-transformer/

The Decoder

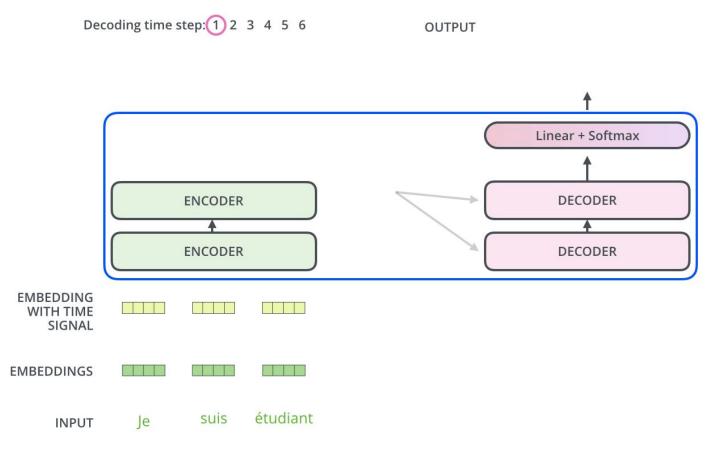


Image source: https://jalammar.github.io/illustrated-transformer/

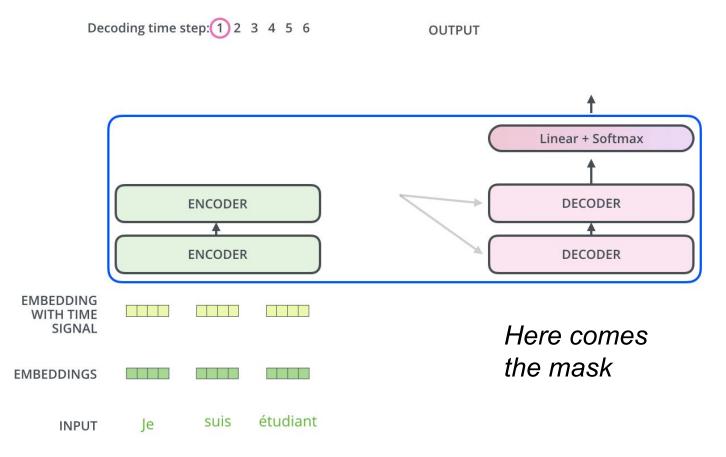
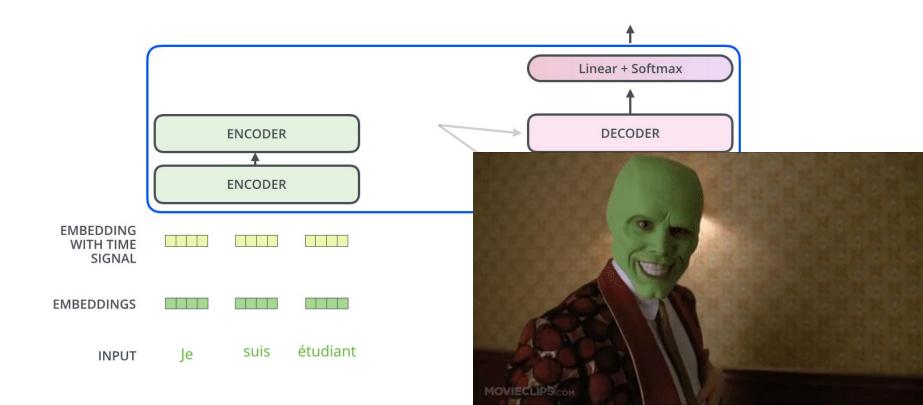
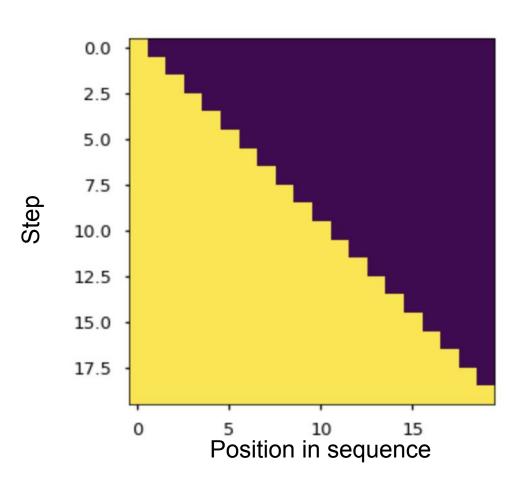


Image source: https://jalammar.github.io/illustrated-transformer/

Decoding time step: 1 2 3 4 5 6 OUTPUT



The masked decoder input



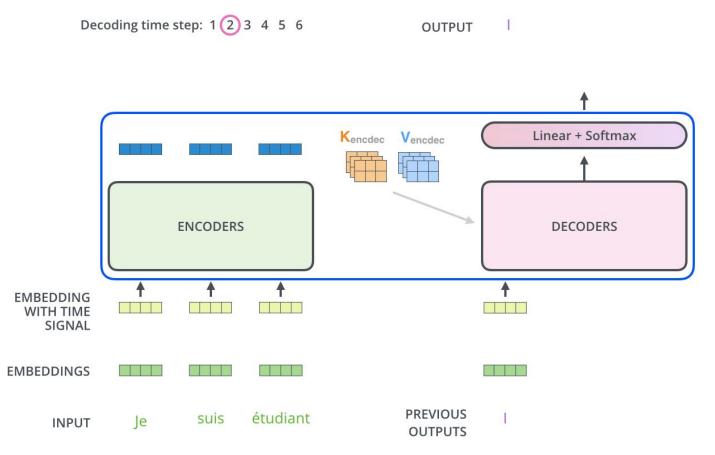


Image source: https://jalammar.github.io/illustrated-transformer/

Final Linear and Softmax Layer

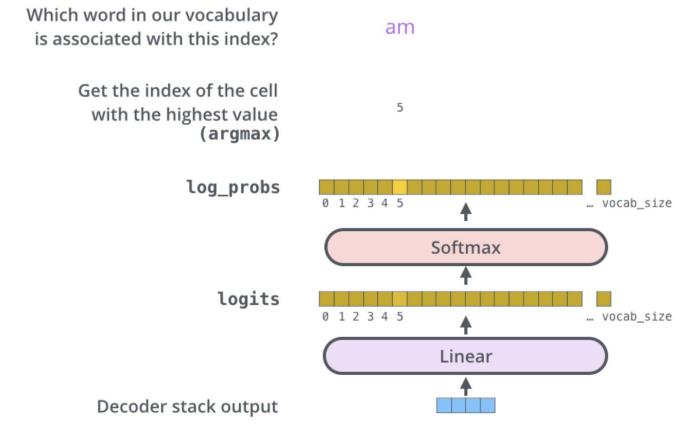
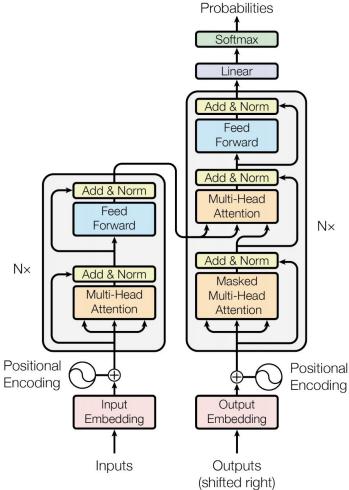


Image source: https://jalammar.github.io/illustrated-transformer/

The Transformer



Output

Image source: Attention Is All You Need, Neural Information Processing Systems 2017

Outro and Q&A

- Transformer is novel and very powerful architecture
- It is worth it to understand how Self-Attention works
- Physical analogues can help you

Further readings are available in the repo