

Deep Learning in Applications

Lecture 5: Transformer overview

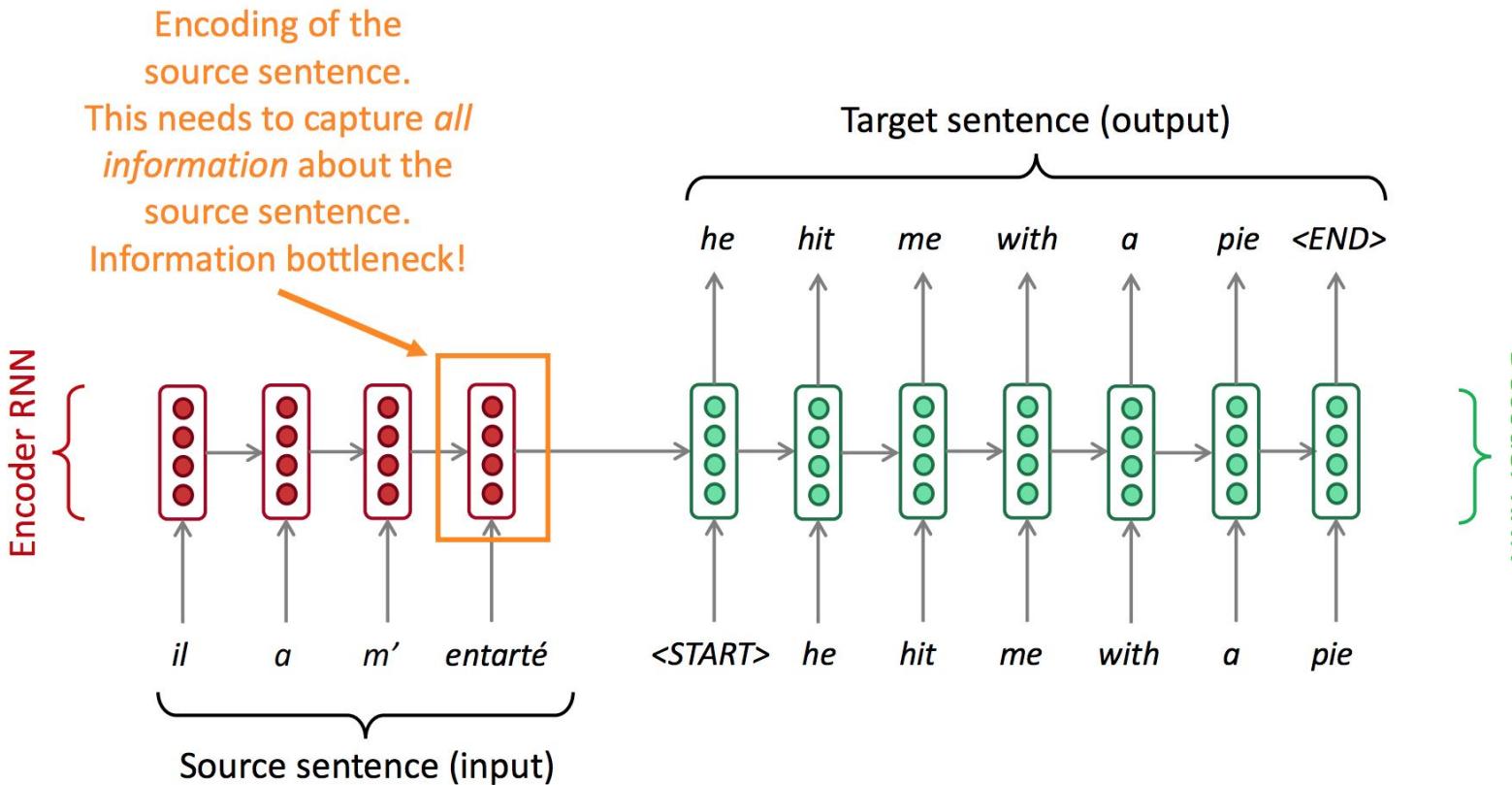
Radoslav Neychev

Harbour.Space University
12.06.2020, Barcelona, Spain

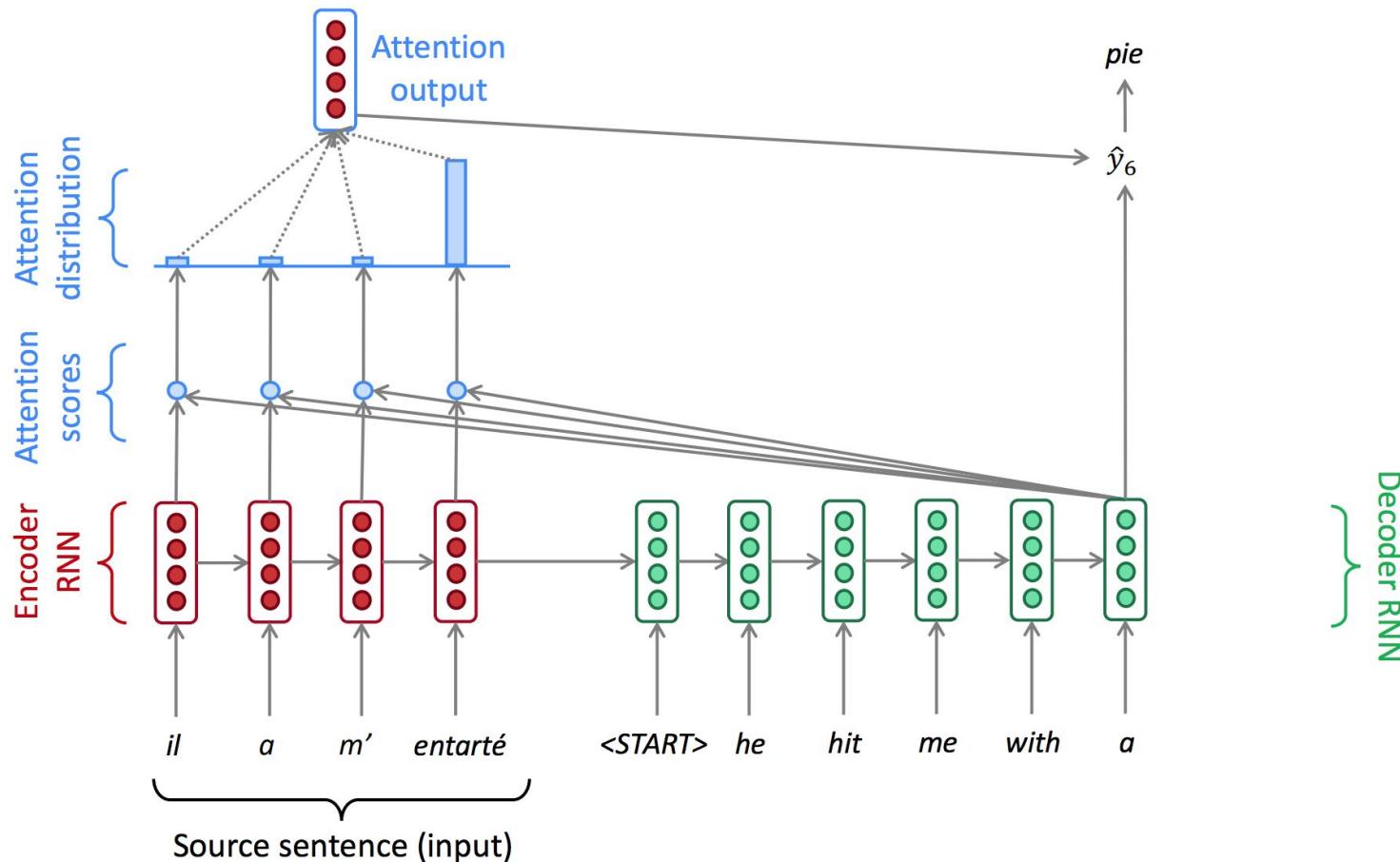
Outline

1. recap: Attention in seq2seq
2. Transformer architecture
3. Self-Attention
4. Positional encoding
5. Layer normalization
6. Q & A

recap: Attention in seq2seq



recap: Attention in seq2seq



recap: Attention in seq2seq

We have encoder hidden states $h_1, \dots, h_N \in \mathbb{R}^h$

On timestep t , we have decoder hidden state $s_t \in \mathbb{R}^h$

We get the attention scores e^t for this step:

$$e^t = [s_t^T h_1, \dots, s_t^T h_N] \in \mathbb{R}^N$$

We take softmax to get the attention distribution α^t for this step (this is a probability distribution and sums to 1)

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

We use α^t to take a weighted sum of the encoder hidden states to get the attention output a_t

$$a_t = \sum_{i=1}^N \alpha_i^t h_i \in \mathbb{R}^h$$

Finally we concatenate the attention output a_t with the decoder hidden state s_t and proceed as in the non-attention seq2seq model

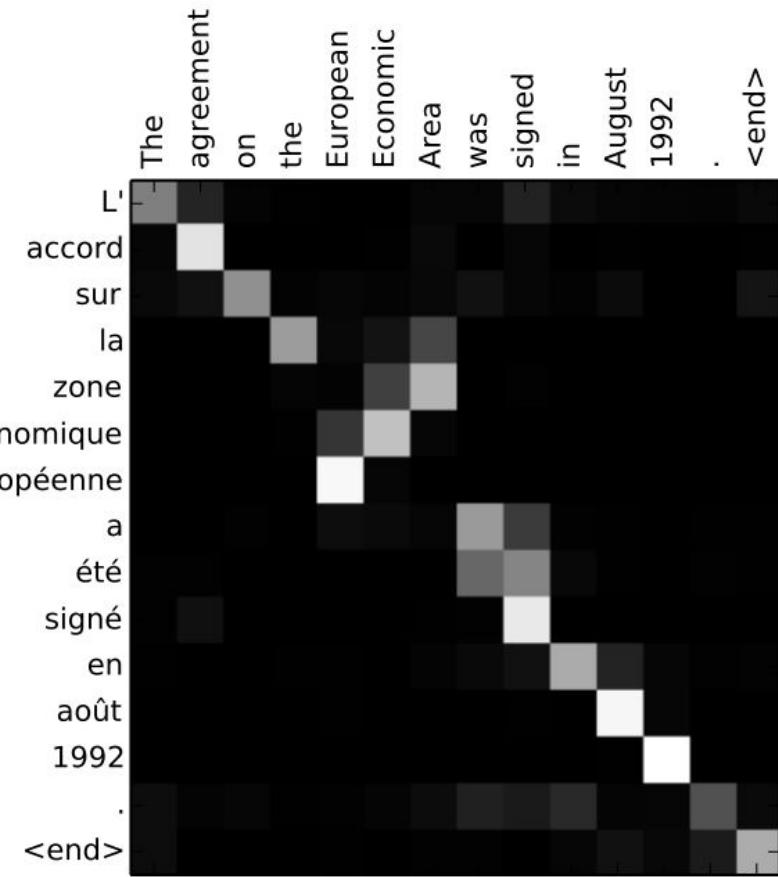
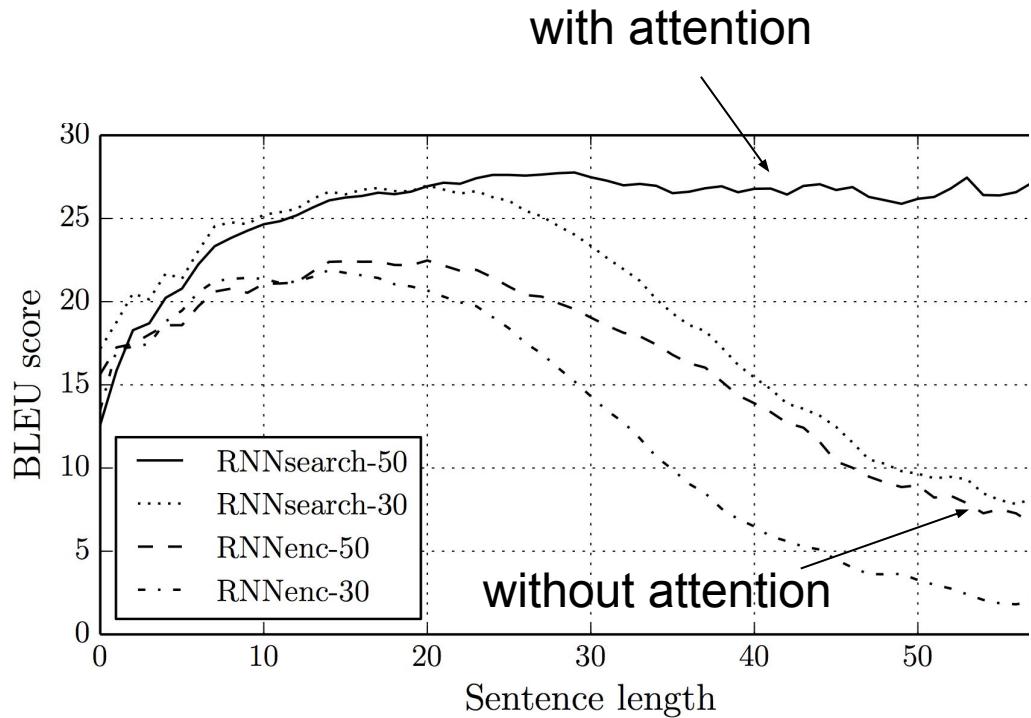
$$[a_t; s_t] \in \mathbb{R}^{2h}$$

Attention variants

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

Attention advantages

- “Free” word alignment
- Better results on long sequences

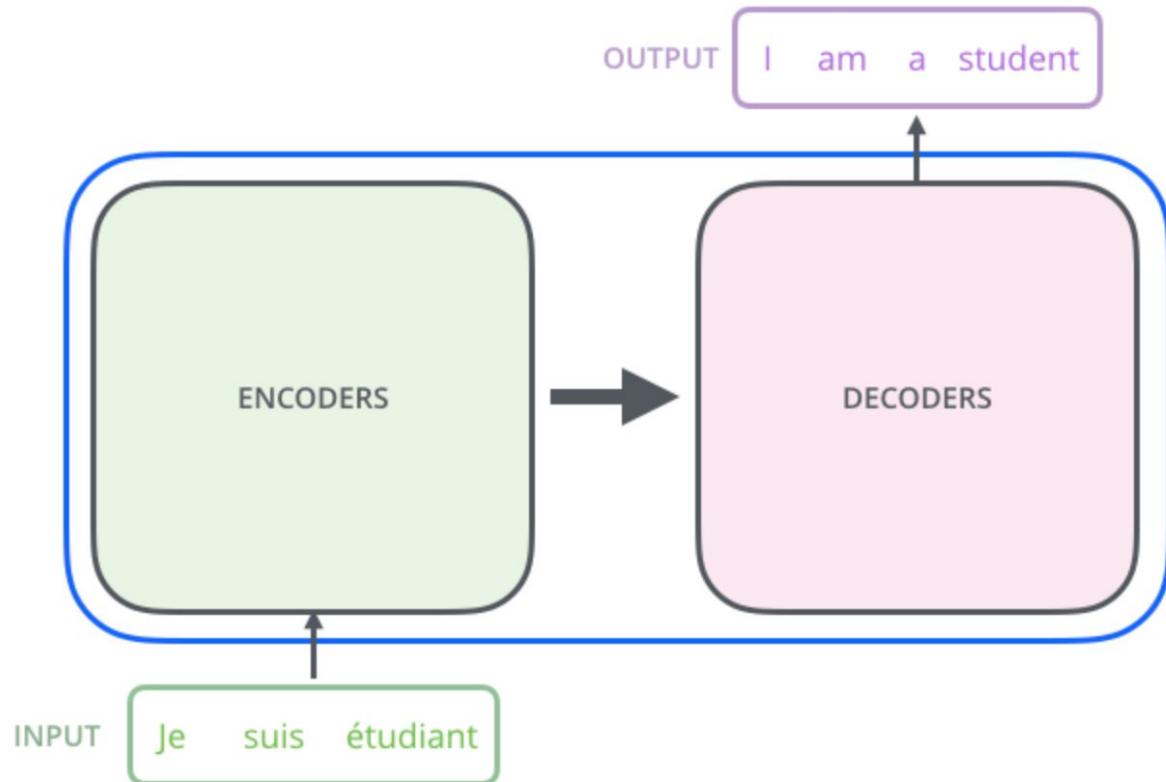


The Transformer

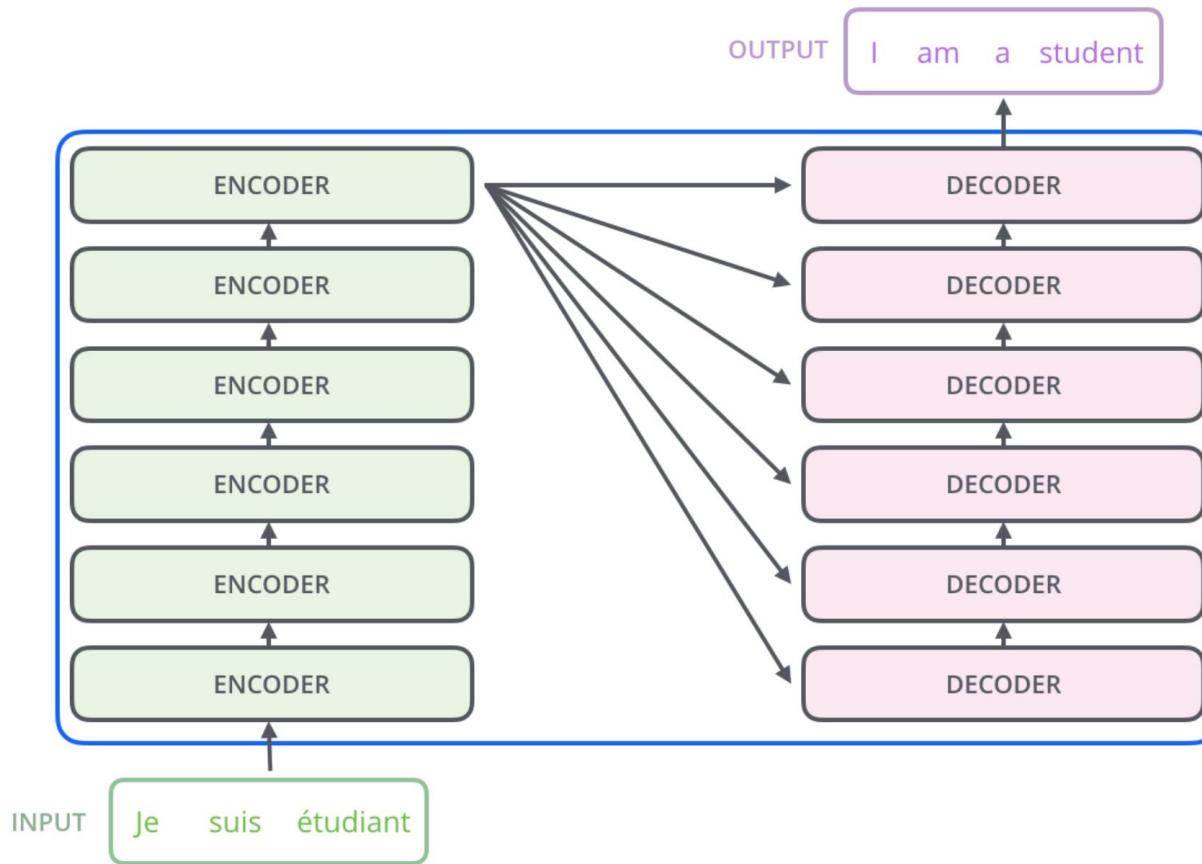
The Transformer



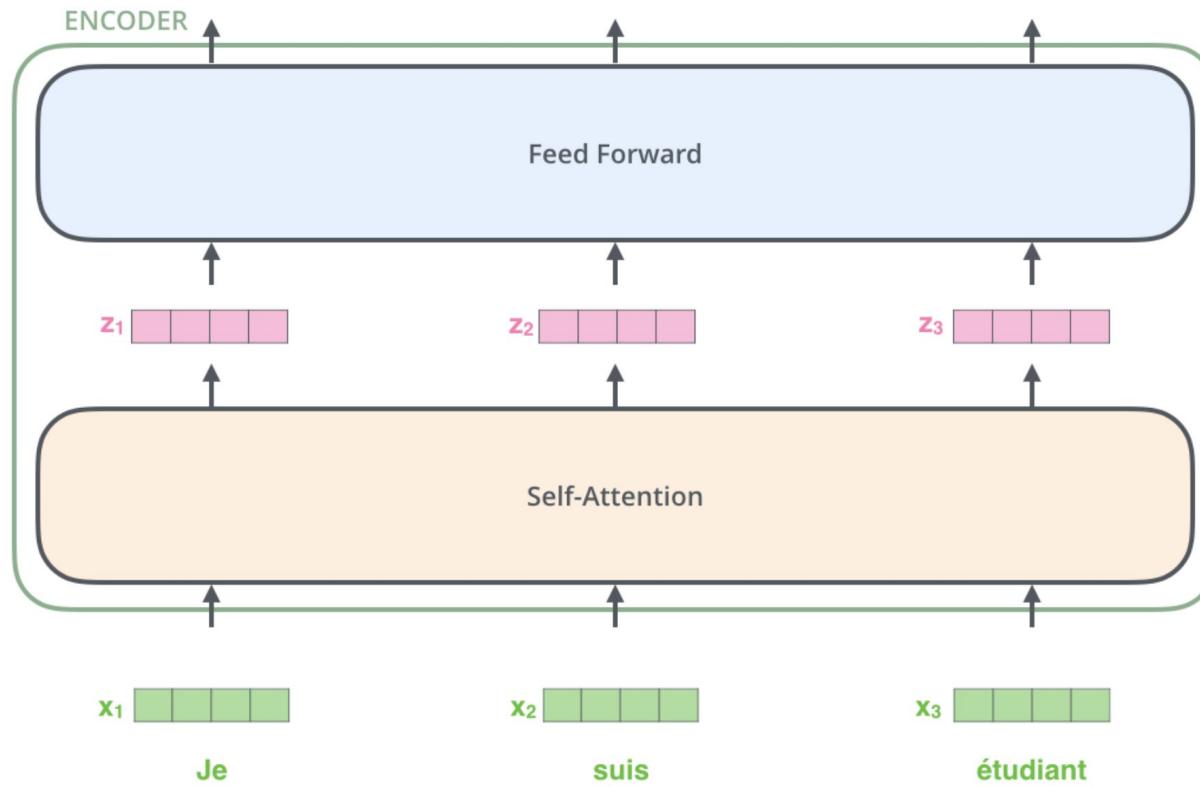
The Transformer



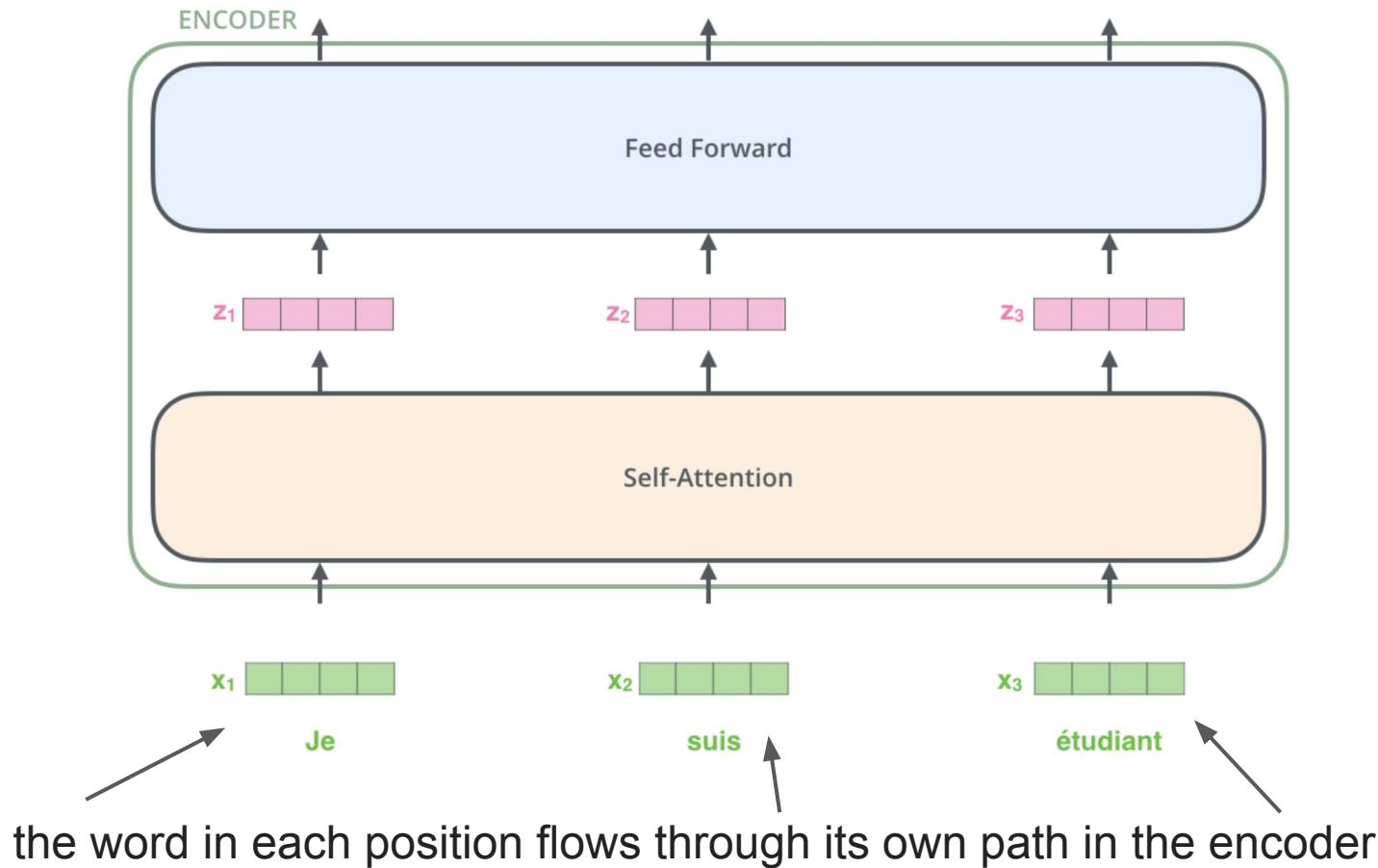
The Transformer



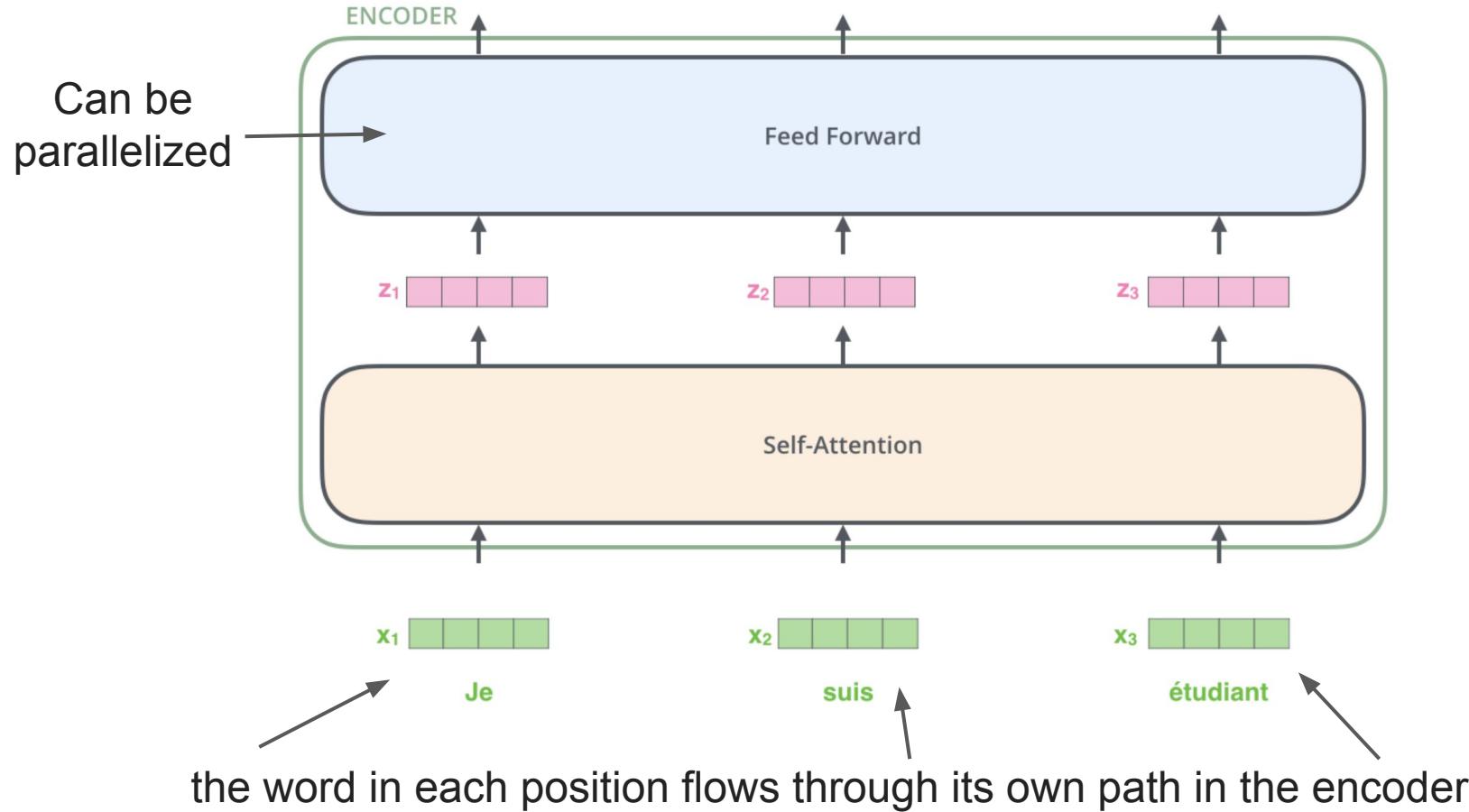
The Encoder Side



The Encoder Side

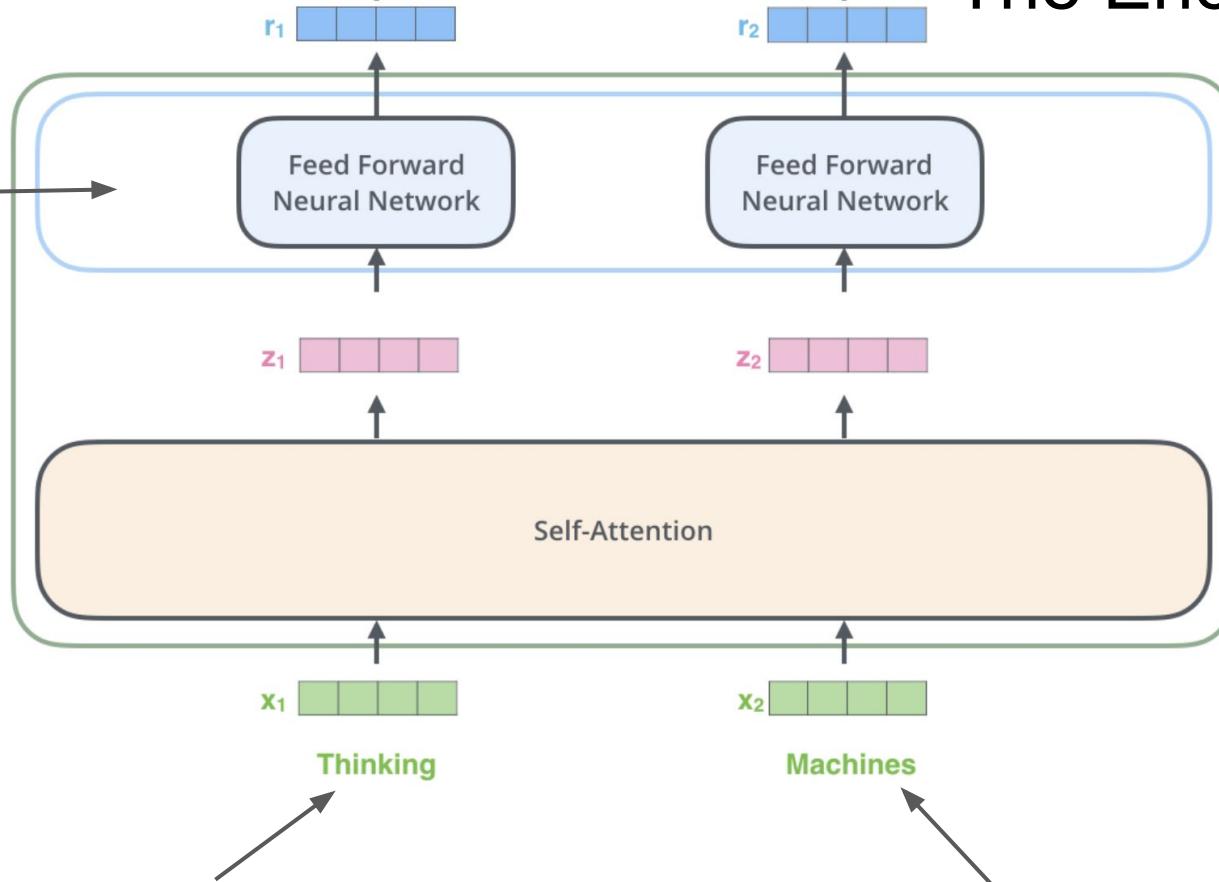


The Encoder Side



The Encoder Side

Can be parallelized



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

The Transformer: quick overview

- Proposed in the paper “Attention is All You Need” (Ashish Vaswani et al.)
- No recurrent or convolutional neural networks -> just 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?

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 at a High Level

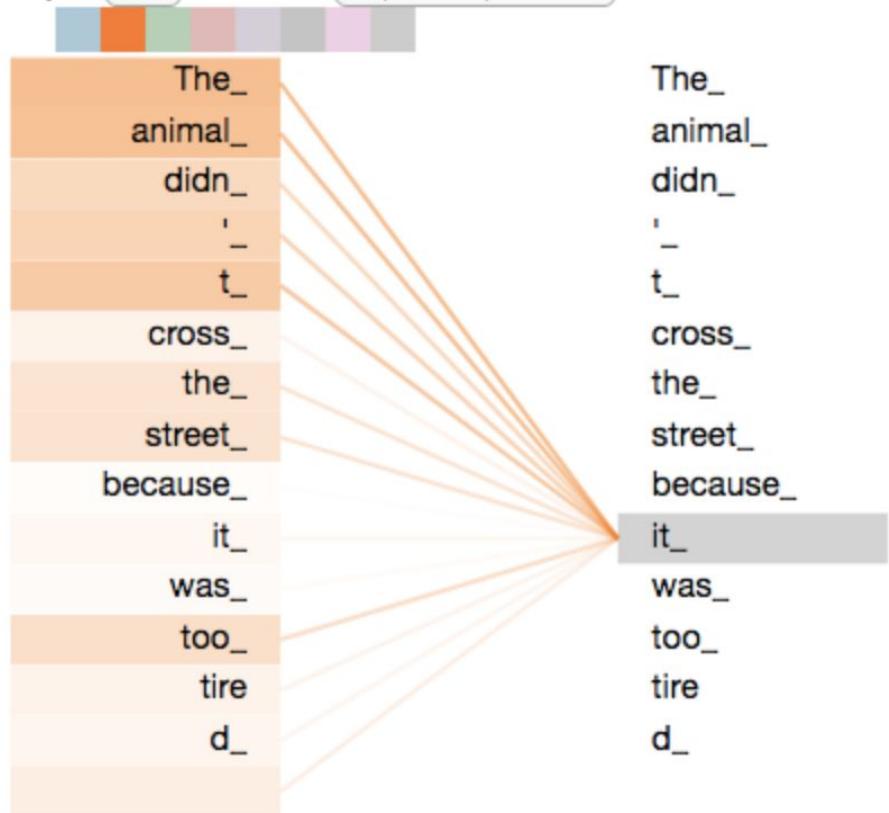
”The animal didn't cross the street because it was too tired”

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

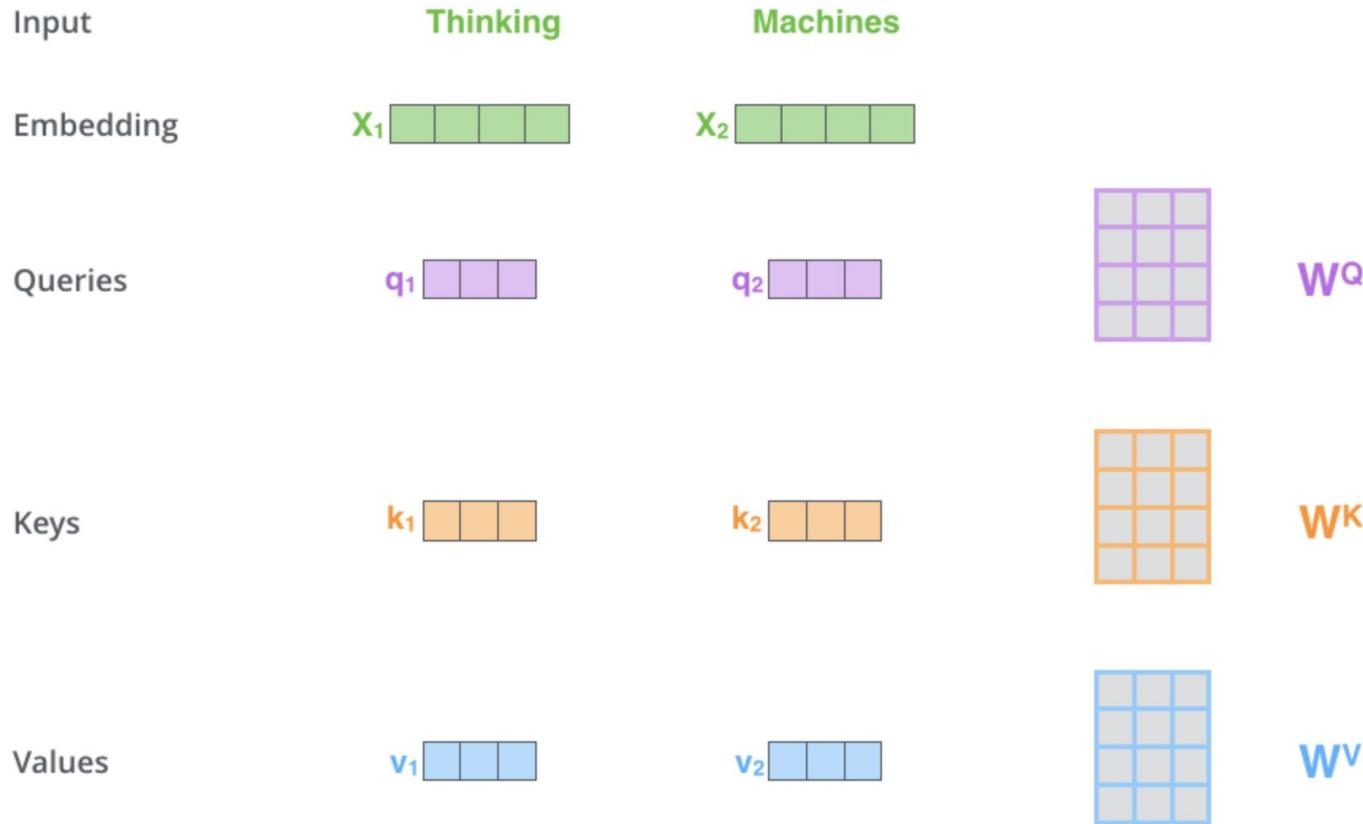
- Self-attention is the method the Transformer uses to bake the “understanding” of other relevant words into the one we’re currently processing

Self-Attention at a High Level

Layer: 5 ⬆️ Attention: Input - Input ⬆️



Self-Attention: detailed explanation

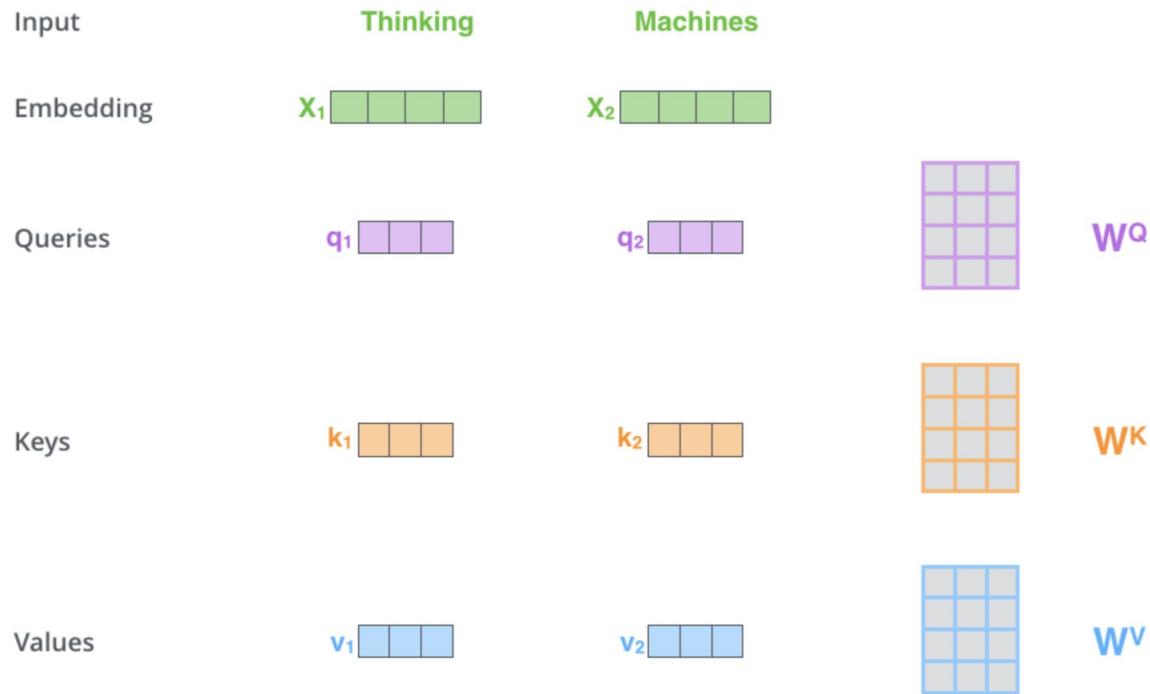


Self-Attention: detailed explanation

STEP 1:

create 3 vectors
(Query, Key, Value)

from each of the encoder's
input vectors



Self-Attention

Input

Embedding

Queries

Keys

Values

Score

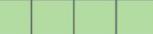
Divide by 8 ($\sqrt{d_k}$)

Softmax

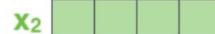
Softmax
X
Value

Sum

Thinking

x_1 

Machines

x_2 

q_1 

q_2 

k_1 

k_2 

v_1 

v_2 

$$q_1 \cdot k_1 = 112$$

$$q_1 \cdot k_2 = 96$$

14

12

0.88

0.12

v_1 

v_2 

z_1 

z_2 

STEP 1: create Query, Key, Value

STEP 2: calculate scores

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

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 (**W_k**, **W_q**, **W_v**)

$$\begin{matrix} \mathbf{X} \\ \begin{array}{|c|c|c|c|} \hline & & & \\ \hline \end{array} \end{matrix} \times \begin{matrix} \mathbf{W^Q} \\ \begin{array}{|c|c|c|c|} \hline & & & \\ \hline \end{array} \end{matrix} = \begin{matrix} \mathbf{Q} \\ \begin{array}{|c|c|c|} \hline & & \\ \hline & & \\ \hline \end{array} \end{matrix}$$

$$\begin{matrix} \mathbf{X} \\ \begin{array}{|c|c|c|c|} \hline & & & \\ \hline \end{array} \end{matrix} \times \begin{matrix} \mathbf{W^K} \\ \begin{array}{|c|c|c|c|} \hline & & & \\ \hline \end{array} \end{matrix} = \begin{matrix} \mathbf{K} \\ \begin{array}{|c|c|c|} \hline & & \\ \hline & & \\ \hline \end{array} \end{matrix}$$

$$\begin{matrix} \mathbf{X} \\ \begin{array}{|c|c|c|c|} \hline & & & \\ \hline \end{array} \end{matrix} \times \begin{matrix} \mathbf{W^V} \\ \begin{array}{|c|c|c|c|} \hline & & & \\ \hline \end{array} \end{matrix} = \begin{matrix} \mathbf{V} \\ \begin{array}{|c|c|c|} \hline & & \\ \hline & & \\ \hline \end{array} \end{matrix}$$

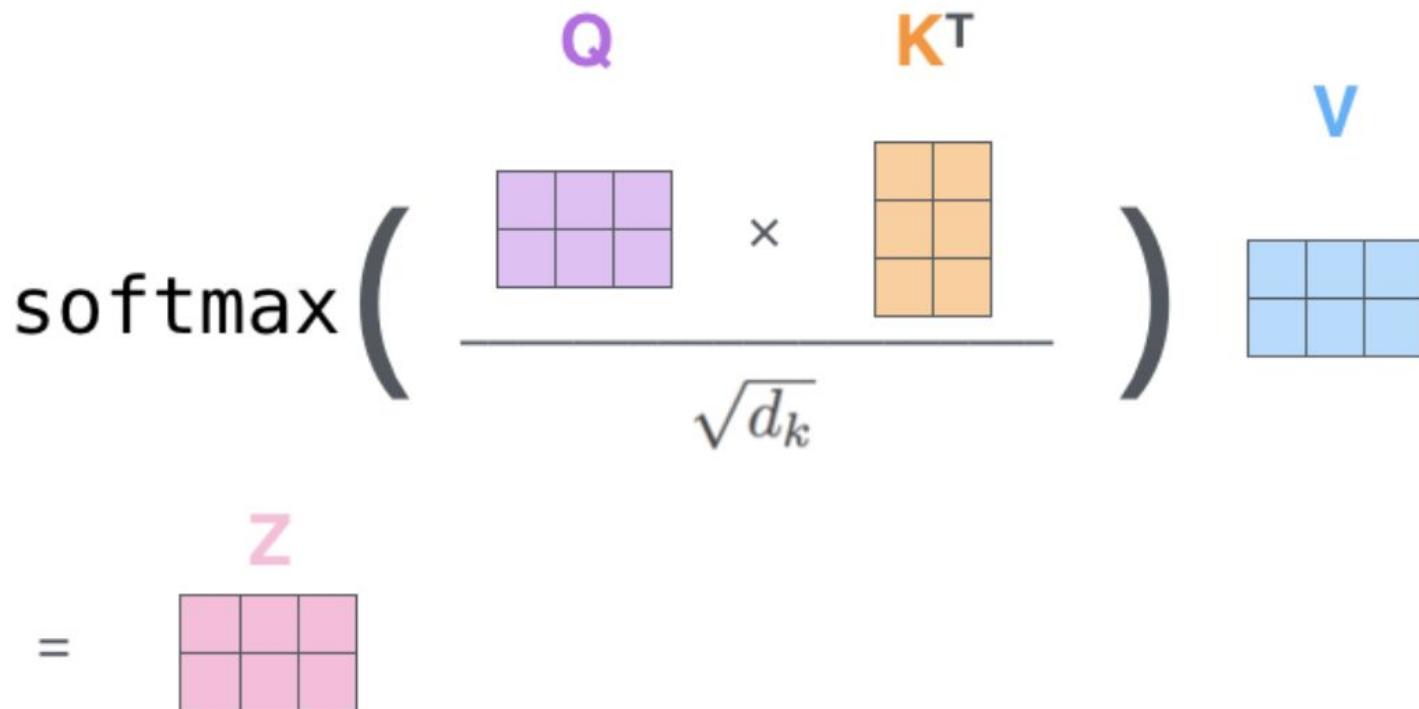
Self-Attention: Matrix Calculation

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

Q K^T V

Z

=

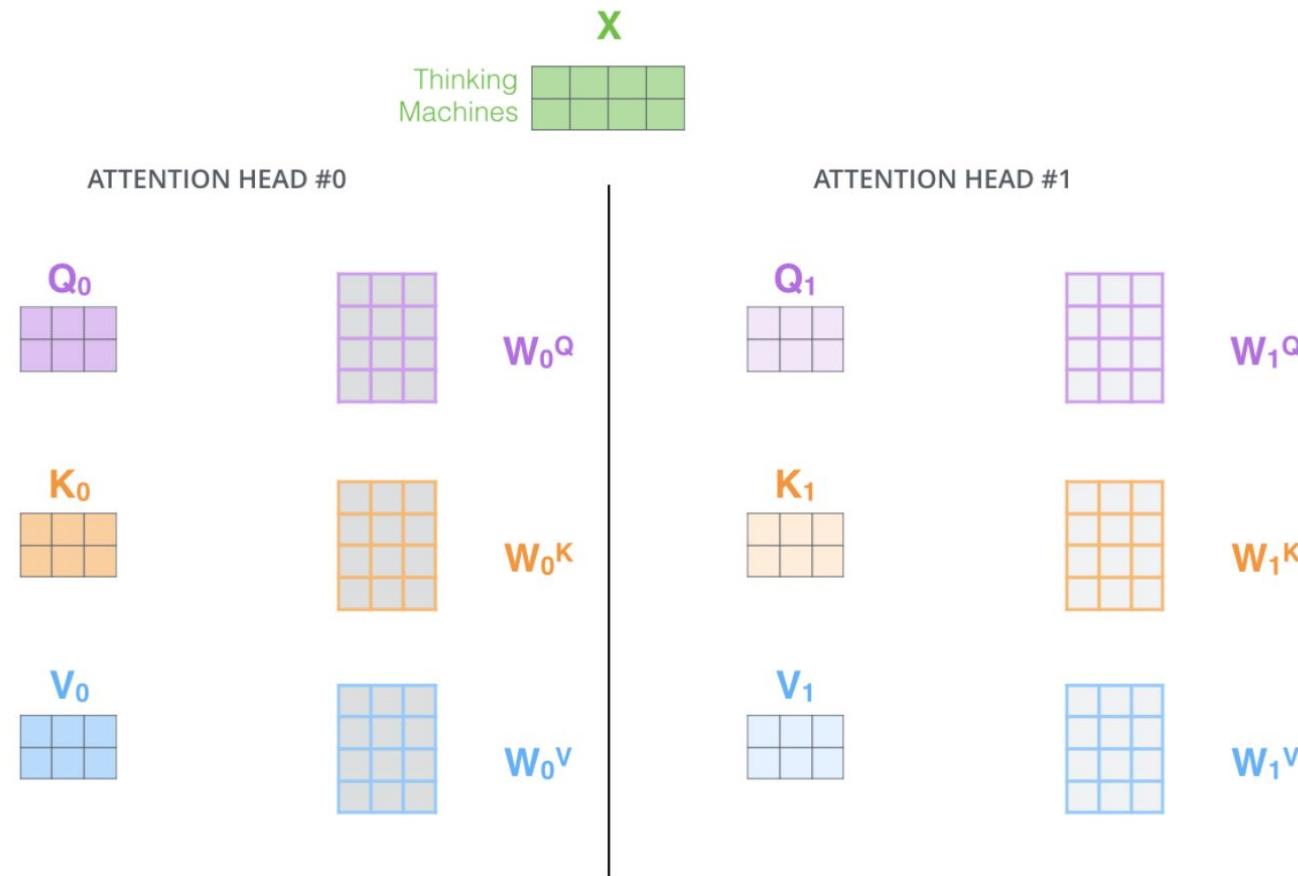


The diagram illustrates the computation of self-attention scores. It shows the softmax function applied to the product of the Query matrix (Q) and the transpose of the Key matrix (K^T), scaled by the square root of the dimension d_k, and then multiplied by the Value matrix (V). The result is labeled Z.

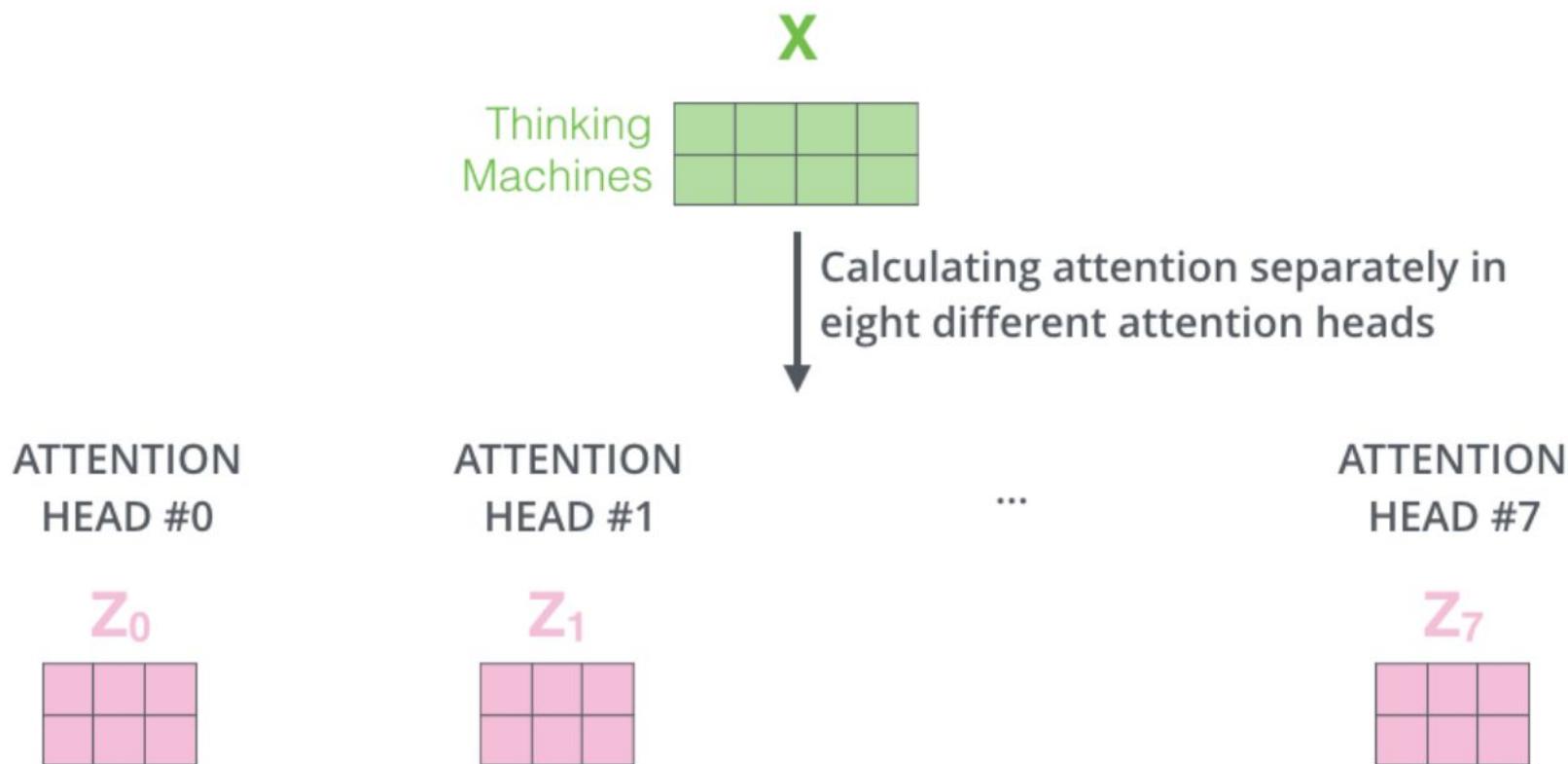
Matrices shown:

- Q (Query): A 3x3 grid of purple squares.
- K^T (Key Transpose): A 3x3 grid of orange squares.
- V (Value): A 3x3 grid of blue squares.
- Z (Output): A 3x3 grid of pink squares.

Multi-Head Attention



Multi-Head Attention



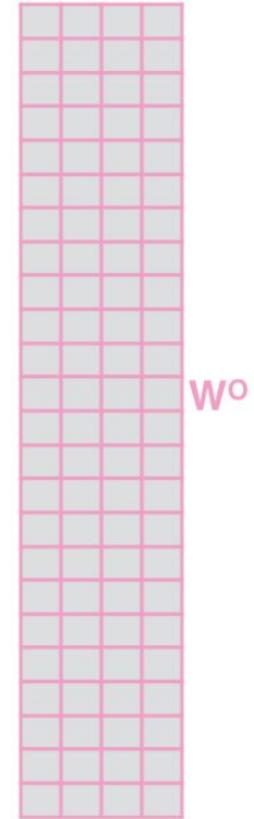
Multi-Head Attention

1) Concatenate all the attention heads



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

X



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

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

1) This is our input sentence*

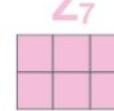
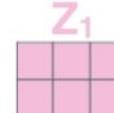
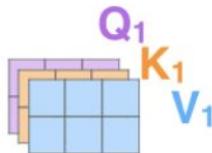
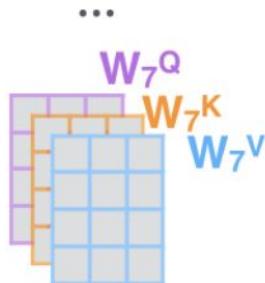
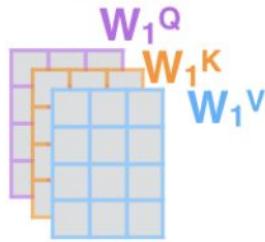
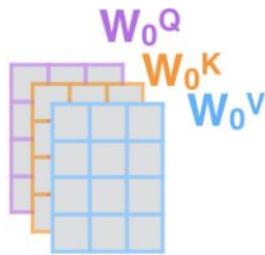
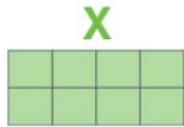
2) We embed each word*

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

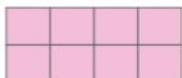
Thinking
Machines



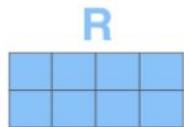
W^o



Z

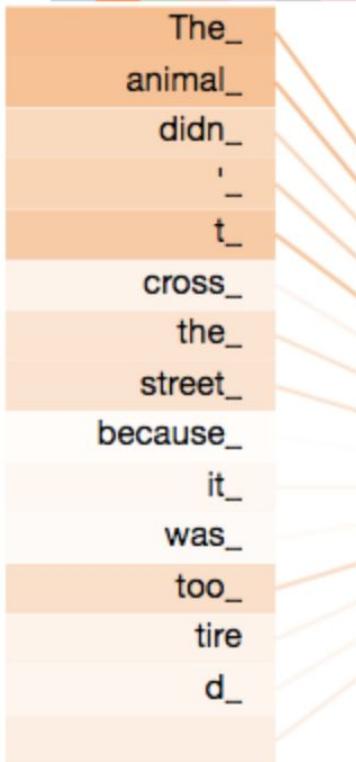


* In all encoders other than #0, we don't need embedding.
We start directly with the output of the encoder right below this one



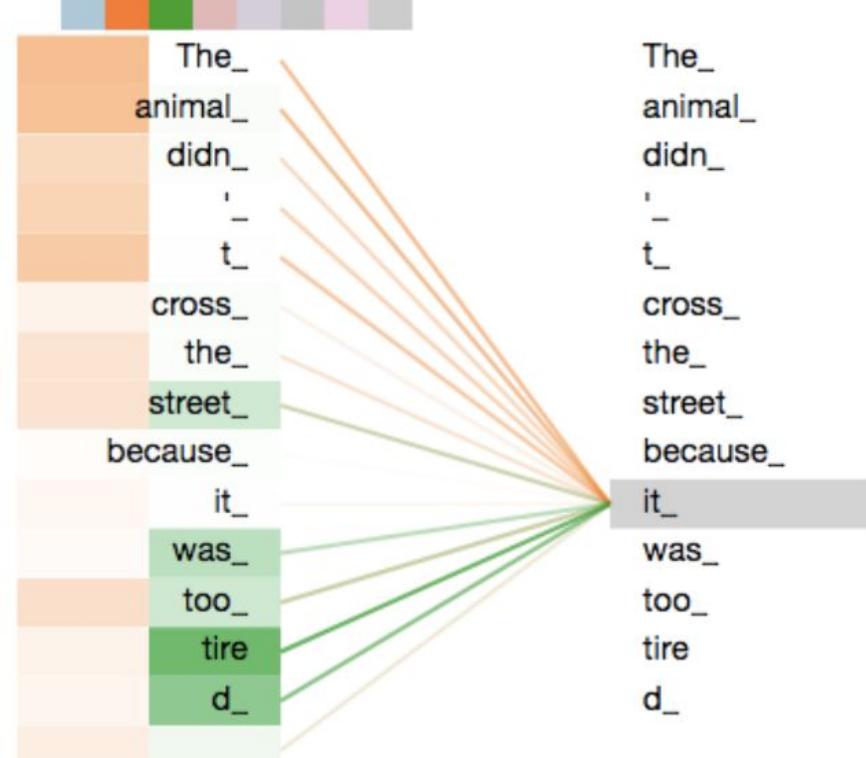
Multi-Head Attention

Layer: 5 ⬆️ Attention: Input - Input ⬆️



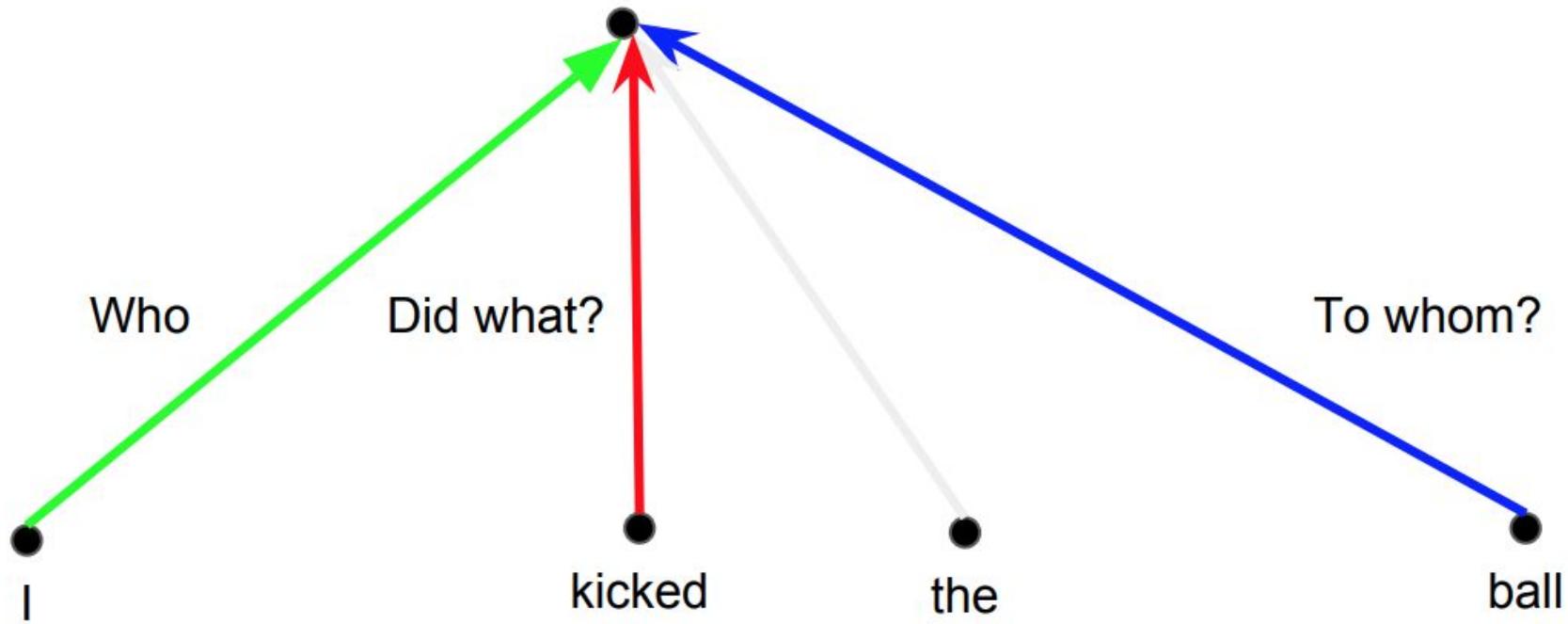
The_ animal_ didn_ ' t_ cross_ the_ street_ because_ it_ was_ too_ tire d_

Layer: 5 ⬆️ Attention: Input - Input ⬆️

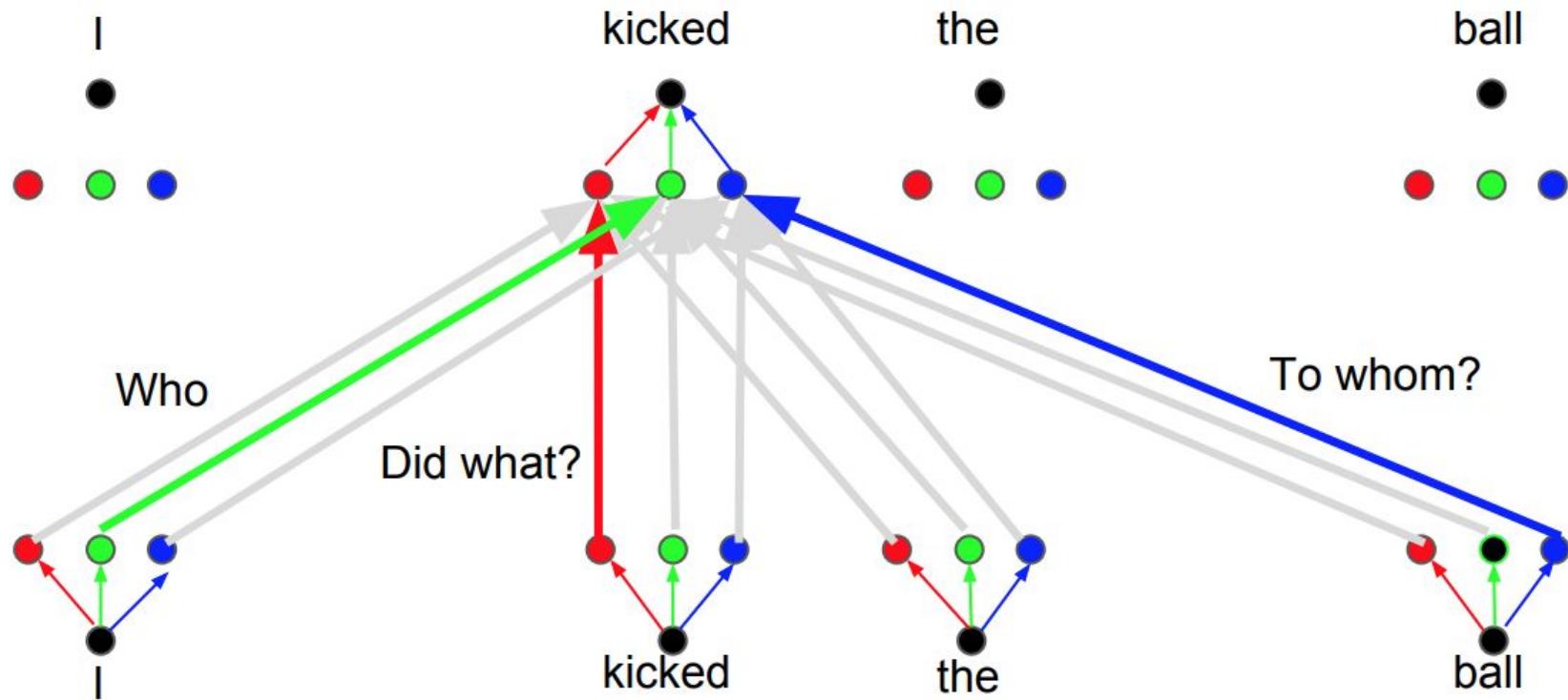


The_ animal_ didn_ ' t_ cross_ the_ street_ because_ it_ was_ too_ tire d_

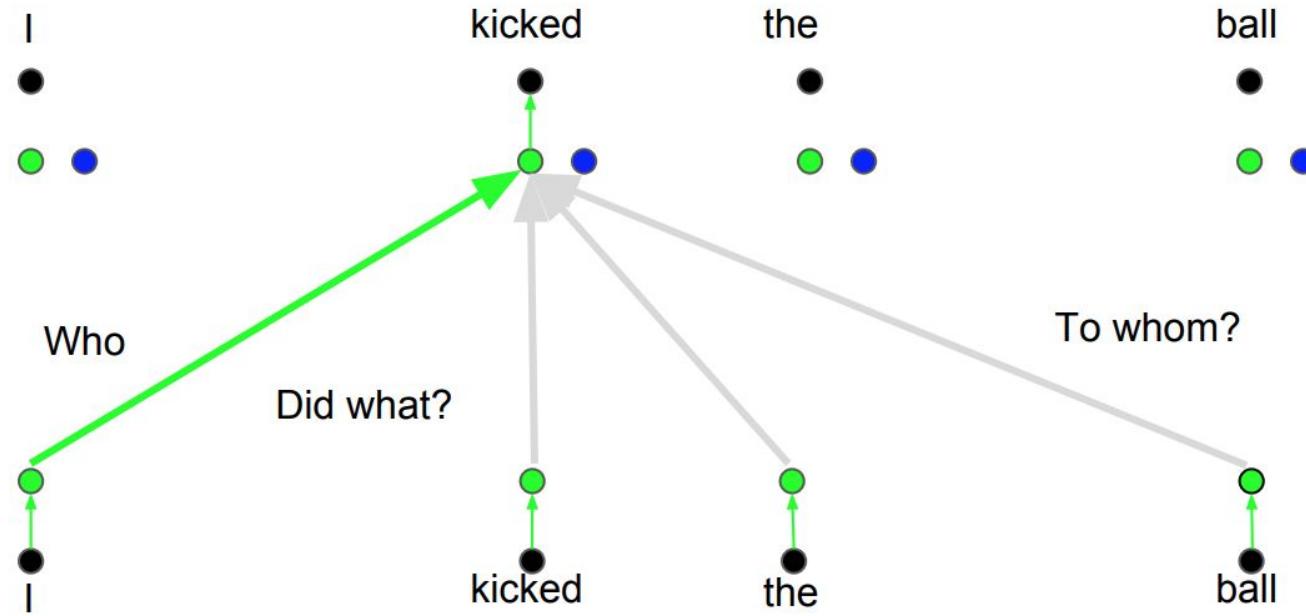
Why Multi-Head Attention?



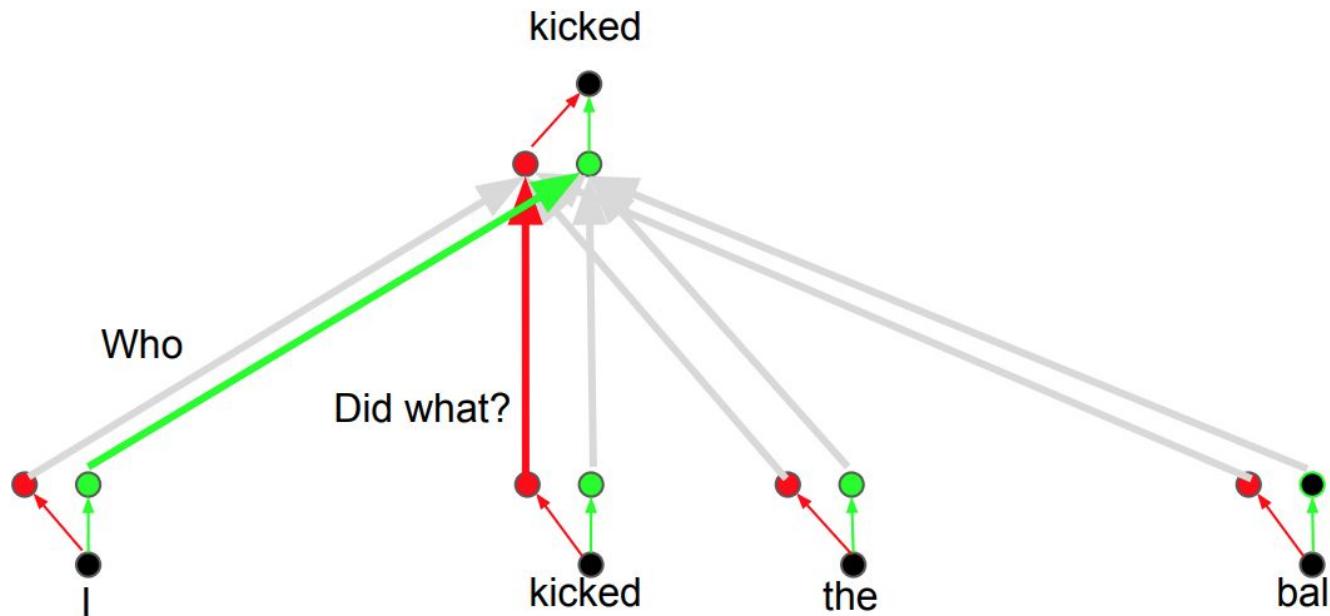
Why Multi-Head Attention?



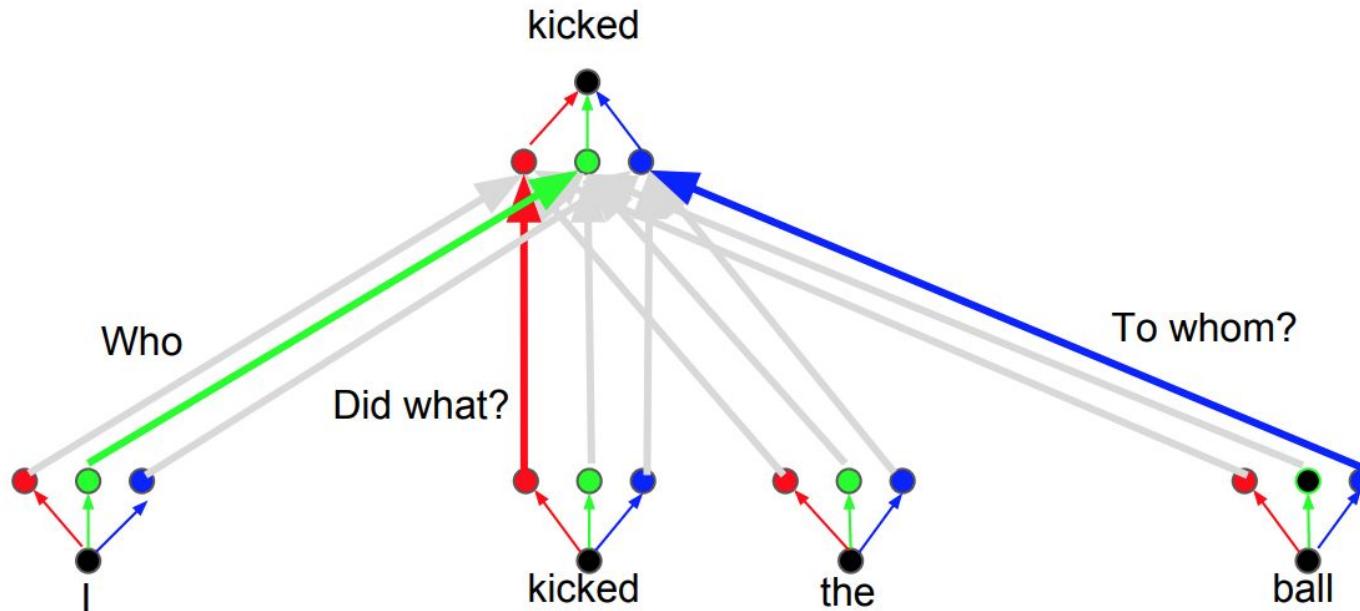
Attention head: Who



Attention head: Did What?

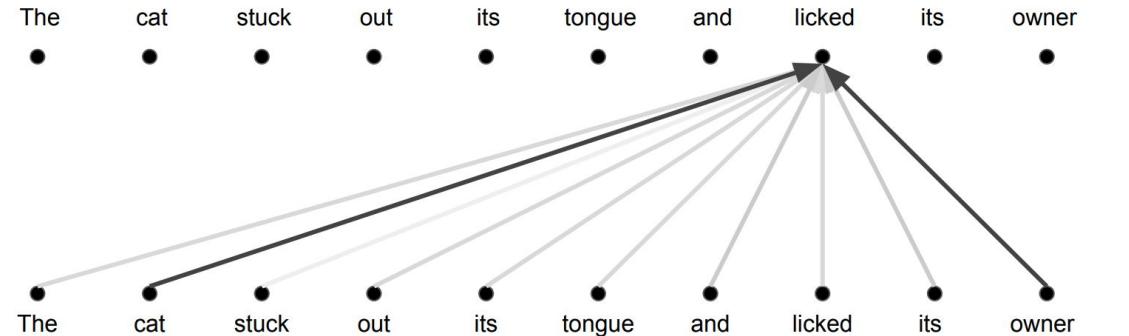


Attention head: To Whom?



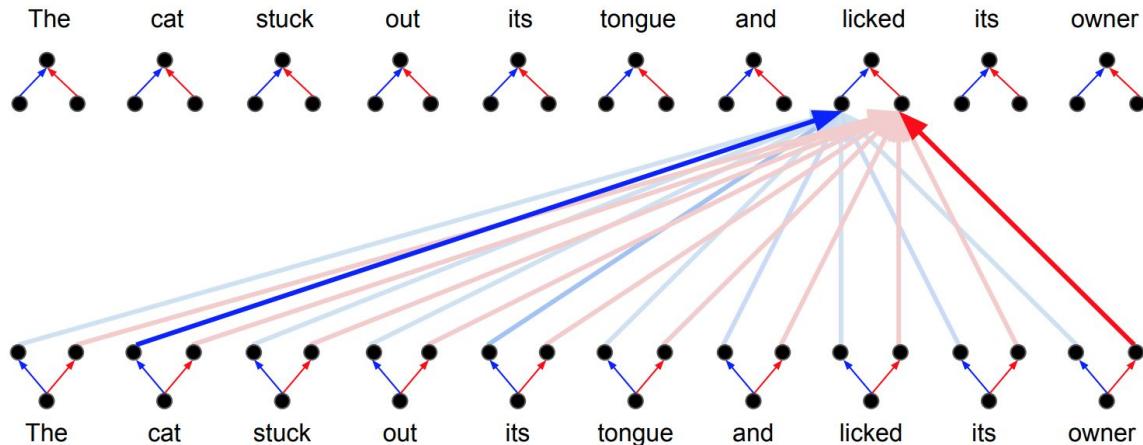
Attention vs. Multi-Head Attention

Attention: a weighted average



Multi-Head Attention:

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

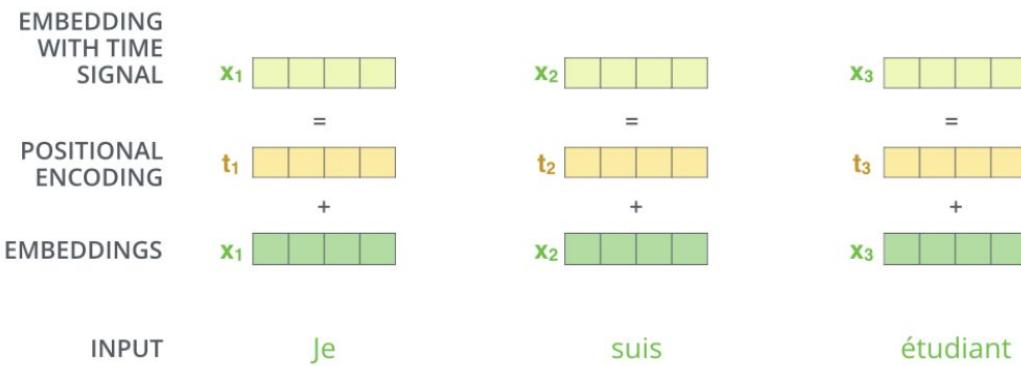
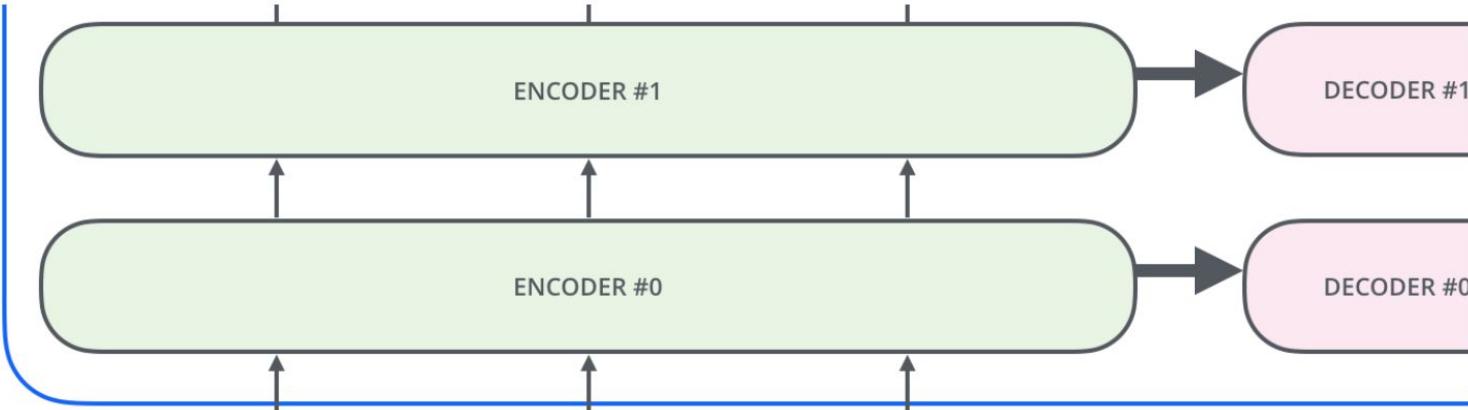


Research Challenges

- 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

Positional Encoding



It provides meaningful distances between the embedding vectors once they're projected into Q/K/V vectors and during dot-product attention

Positional Encoding

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

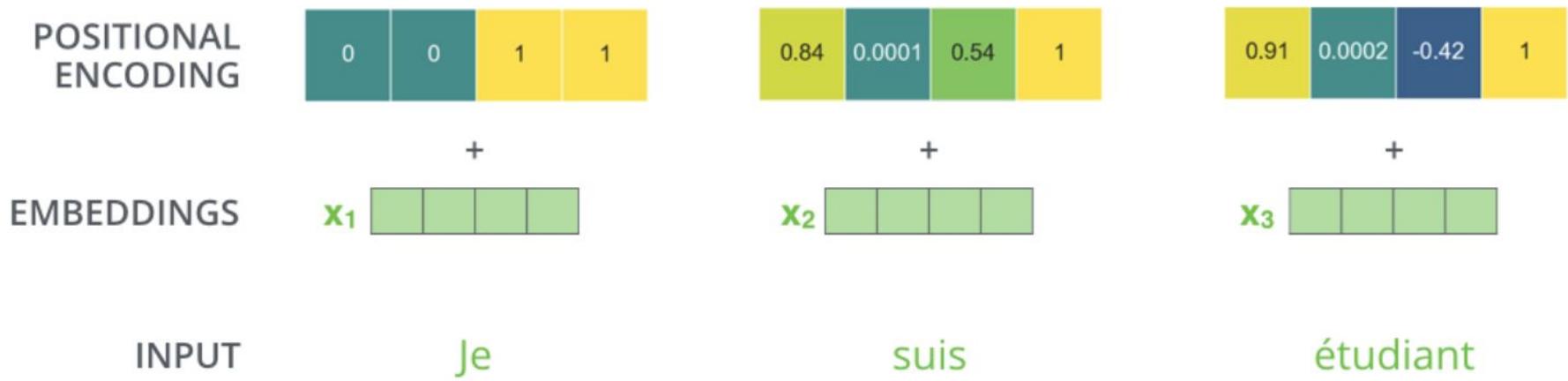
$$PE_{(pos, 2i + 1)} = \cos(pos / 10000^{2i / d_{\text{model}}})$$

- pos is the position
- i is the dimension.

Each dimension of the positional encoding corresponds to a sinusoid.

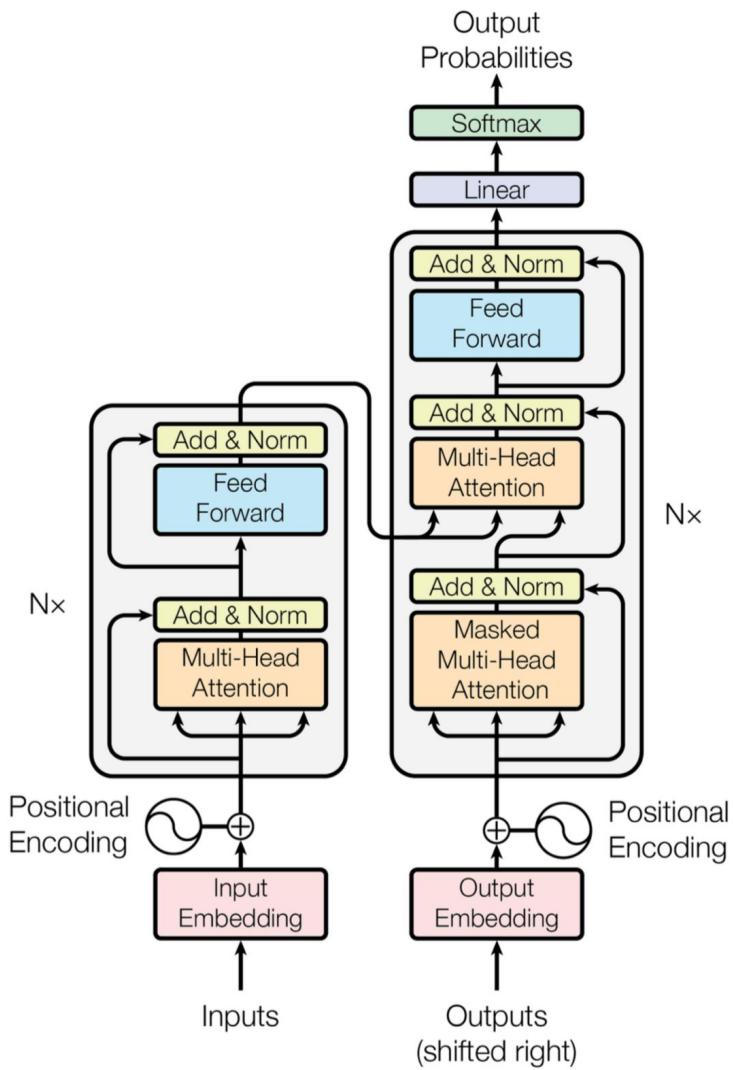
The wavelengths form a geometric progression from 2π to $2\pi * 10000$

Positional Encoding

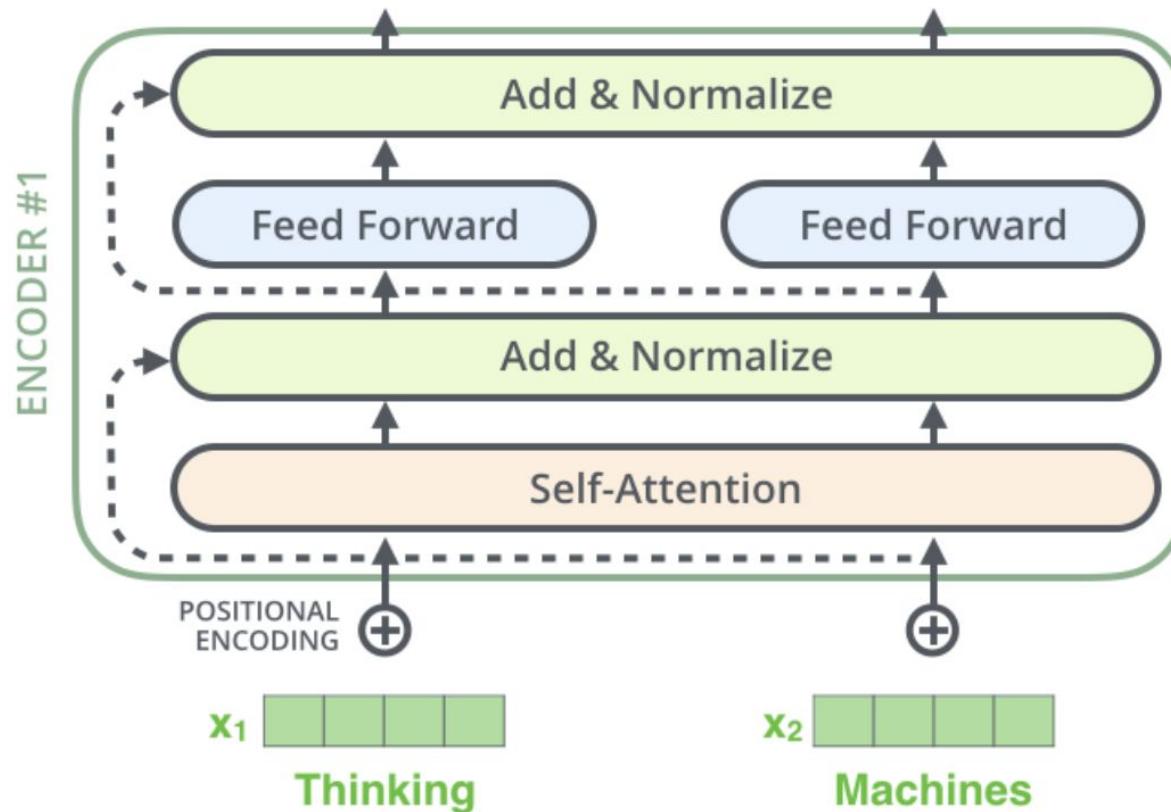


Layer Normalization

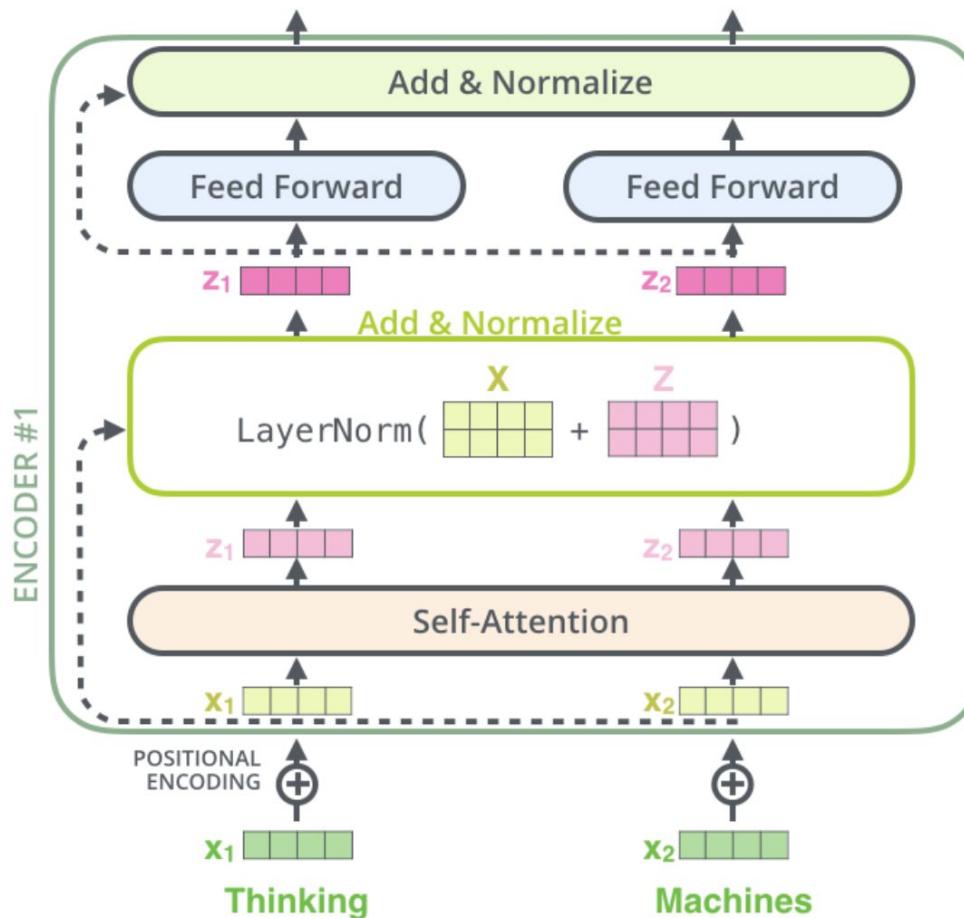
The Transformer: recap



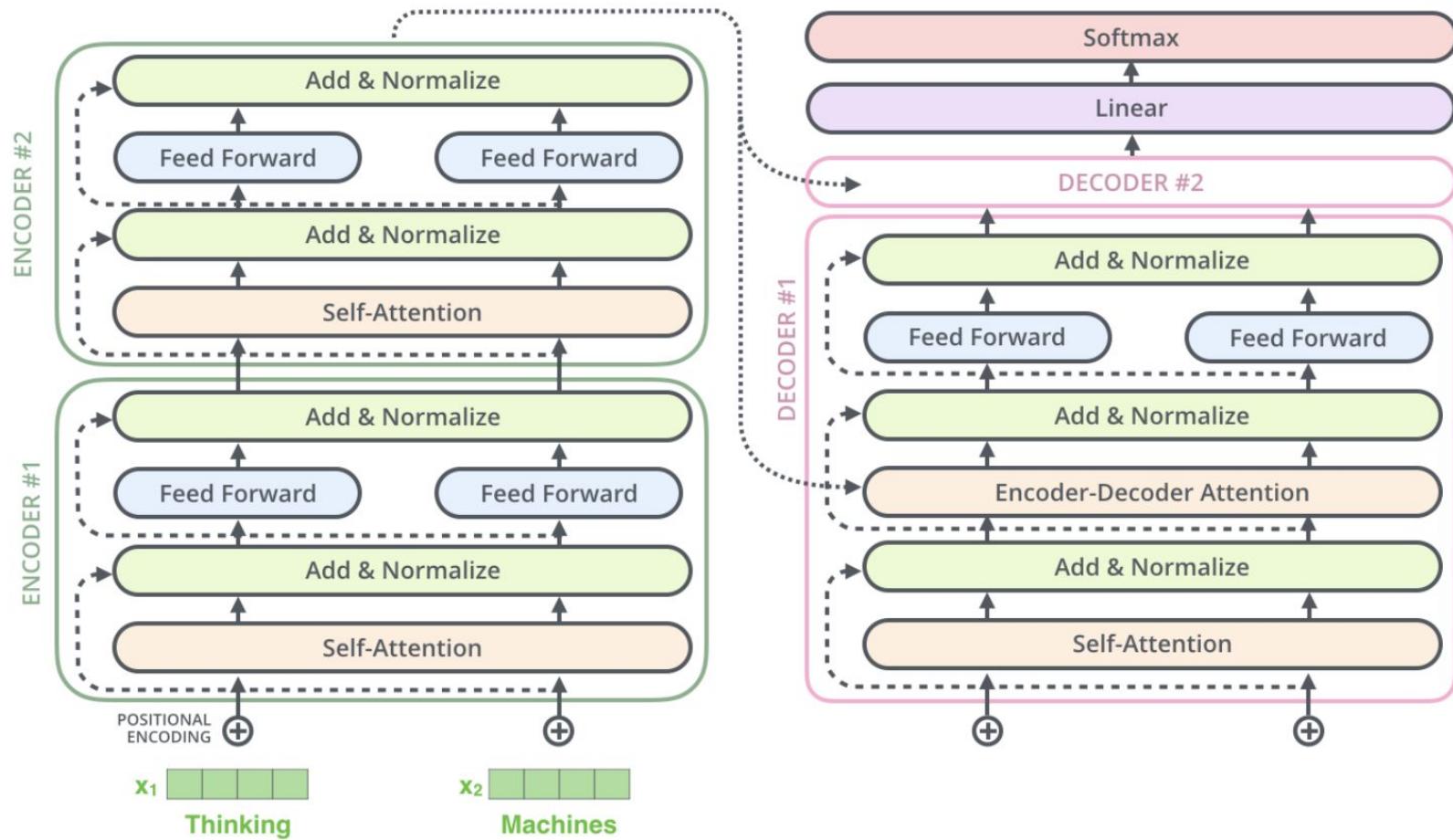
Layer Normalization



Layer Normalization

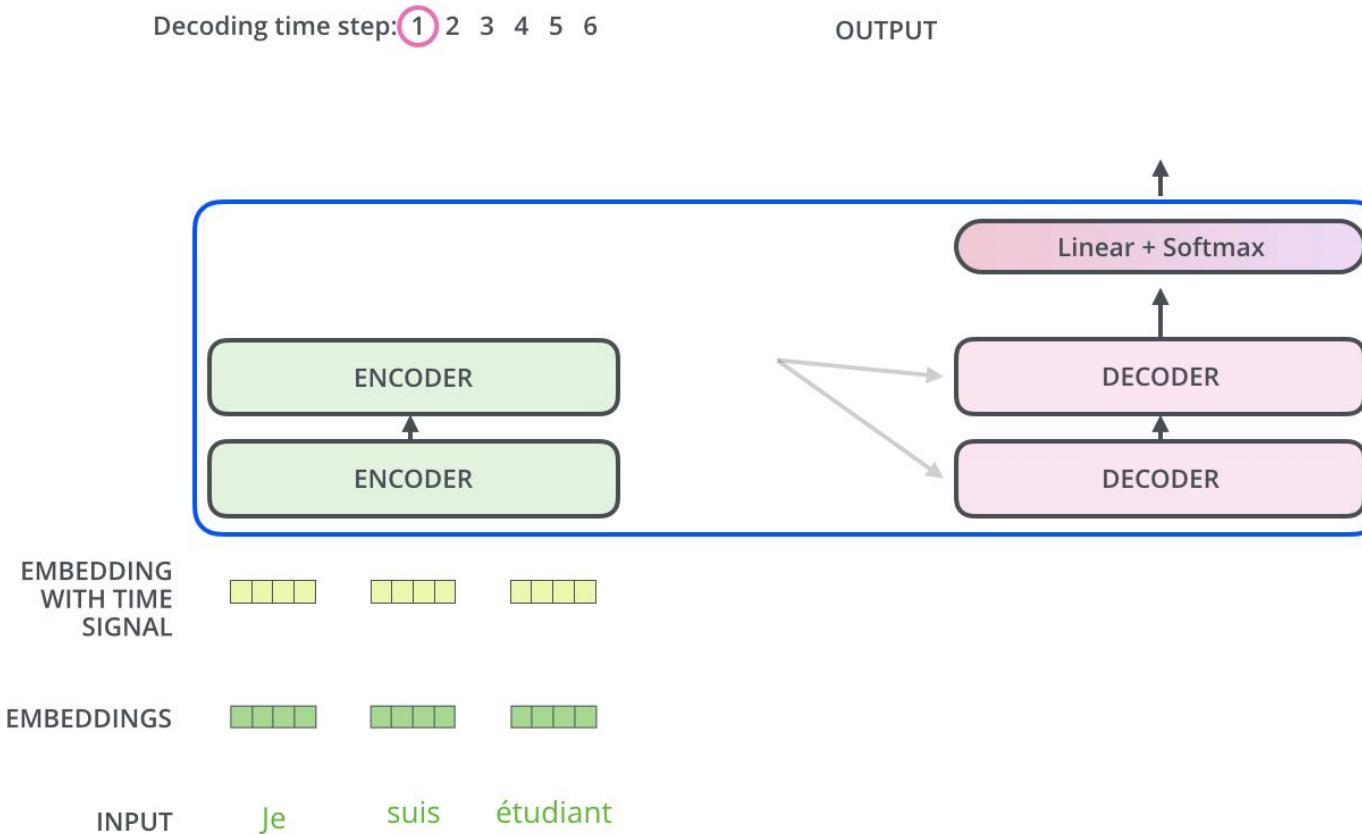


Layer Normalization

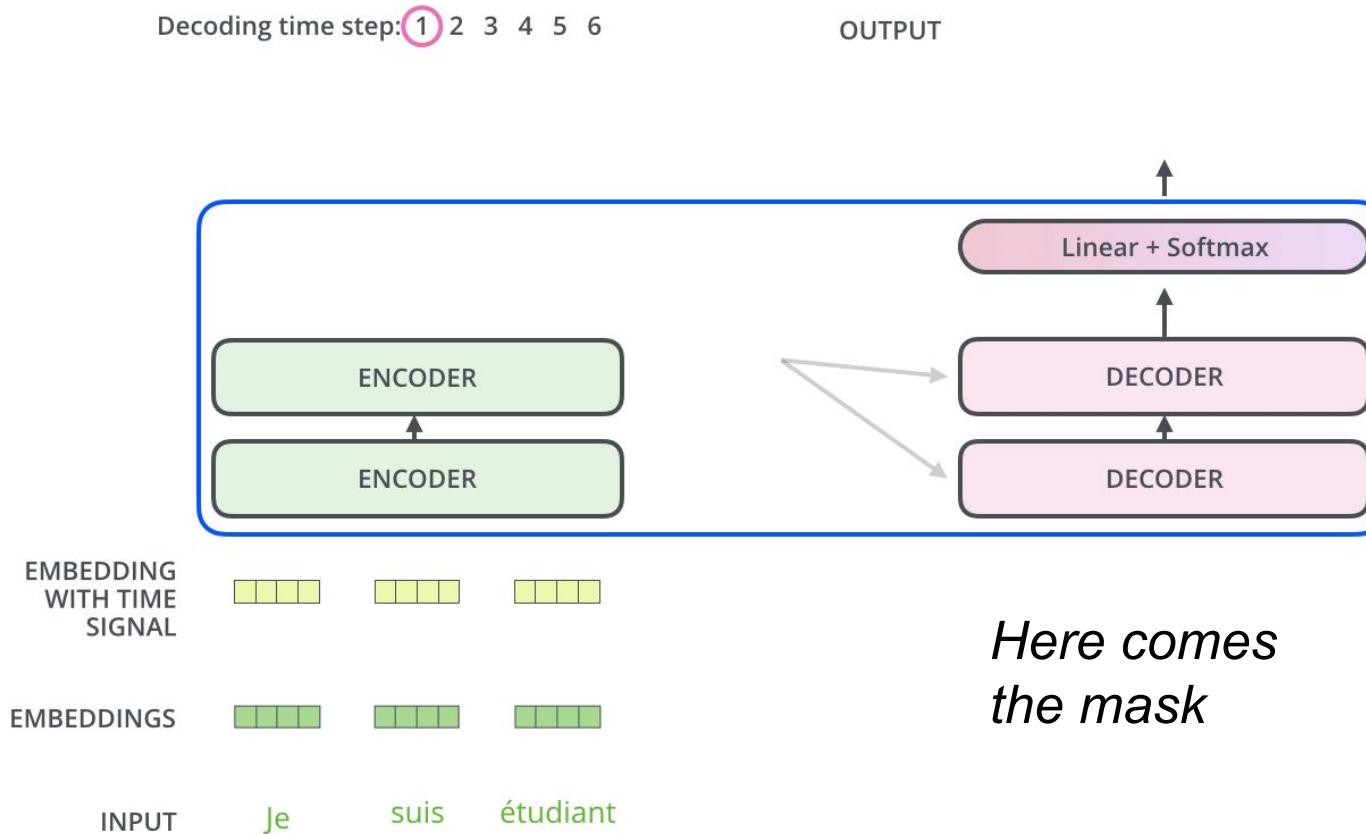


The Decoder

The Decoder Side



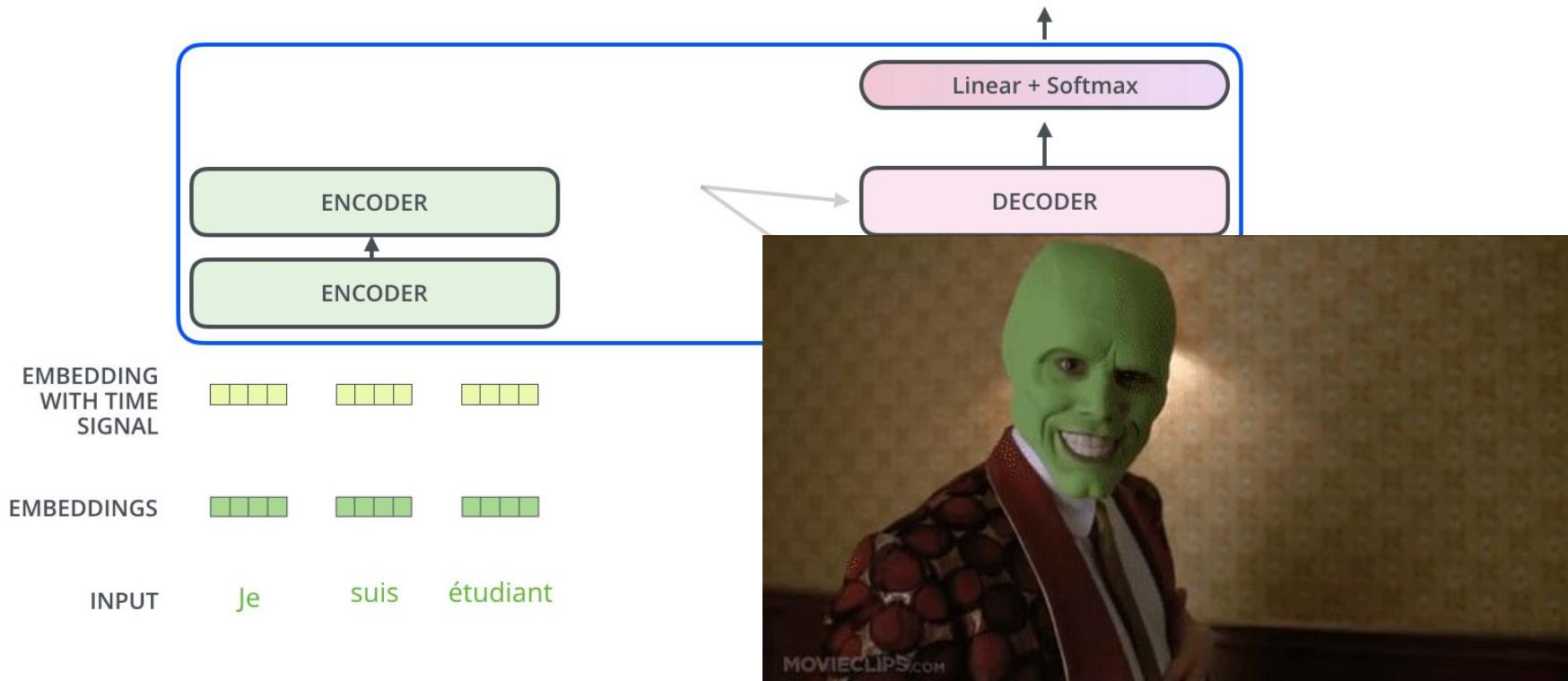
The Decoder Side



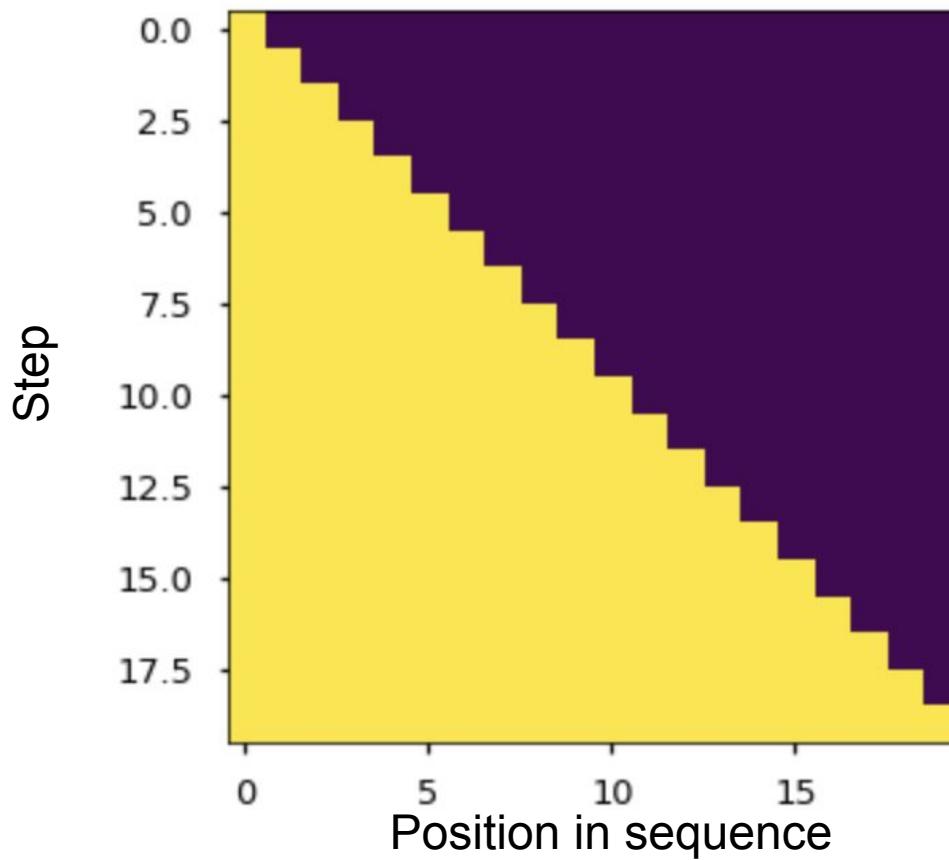
The Decoder Side

Decoding time step: 1 2 3 4 5 6

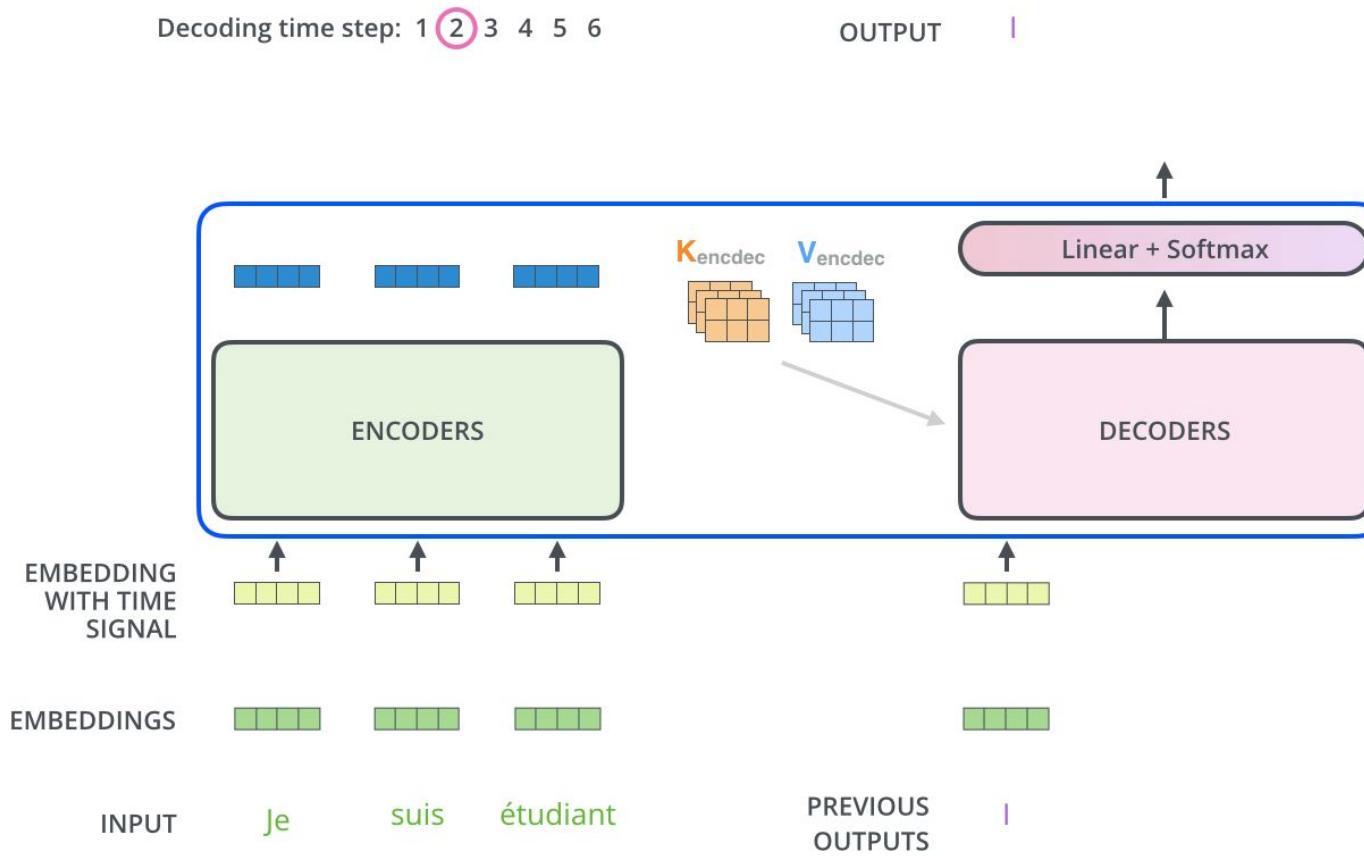
OUTPUT



The masked decoder input



The Decoder Side



Final Linear and Softmax Layer

Which word in our vocabulary
is associated with this index?

am

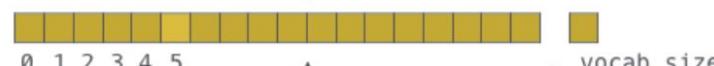
Get the index of the cell
with the highest value
(argmax)

5

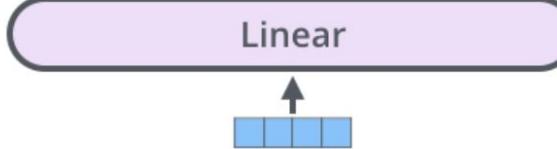
log_probs



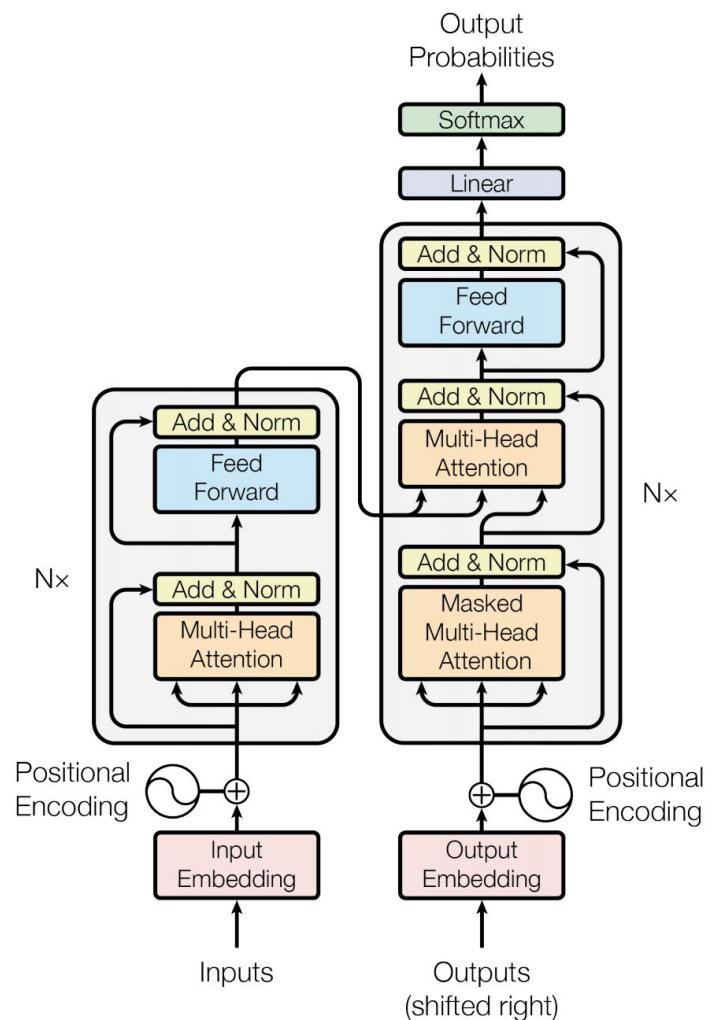
logits

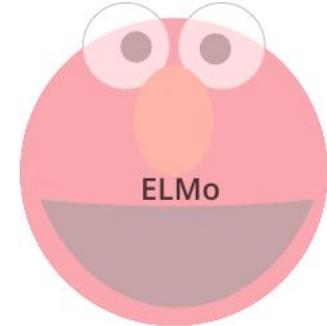


Decoder stack output



The Transformer





OpenAI Transformer: Pre-training Decoder for Language Modeling

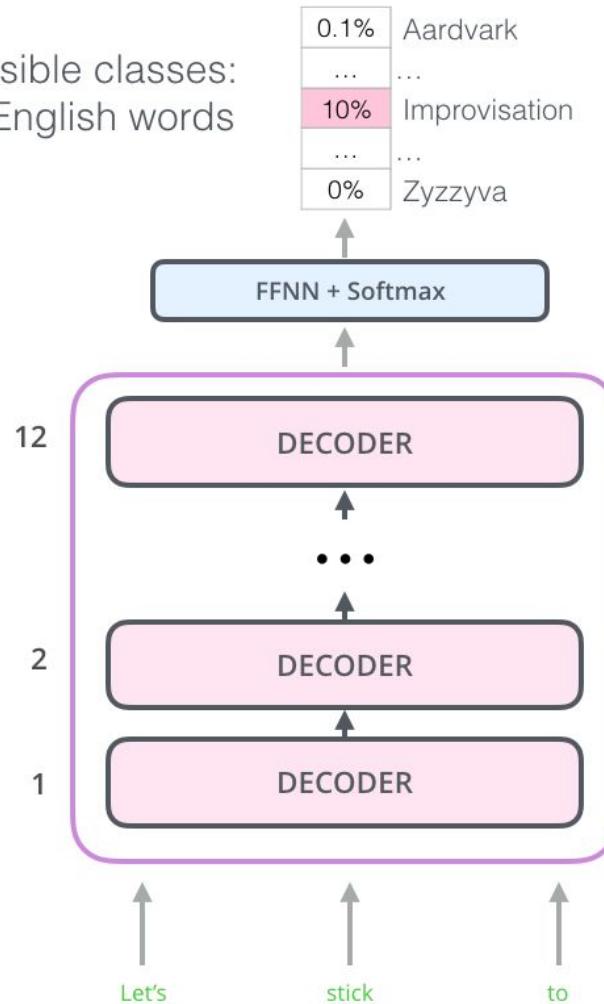
OpenAI Transformer

- The Encoder-Decoder structure of the transformer made it perfect for machine translation
- But what about sentence classification?
- **Main goal: pre-train a language model that can be fine-tuned for other tasks**



OpenAI Transformer

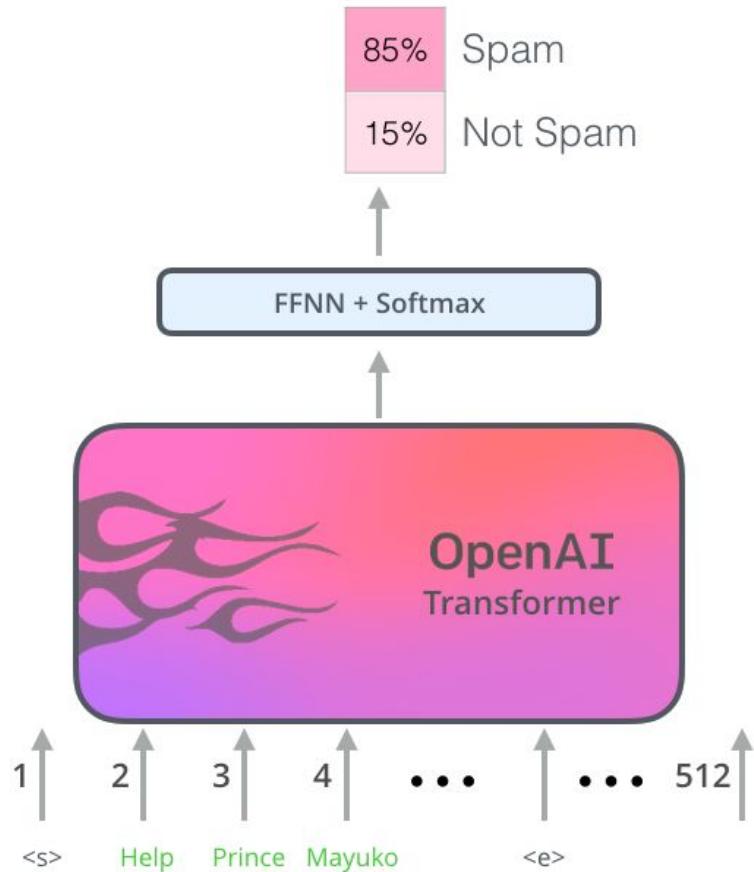
Possible classes:
All English words



Differences from vanilla Transformer:

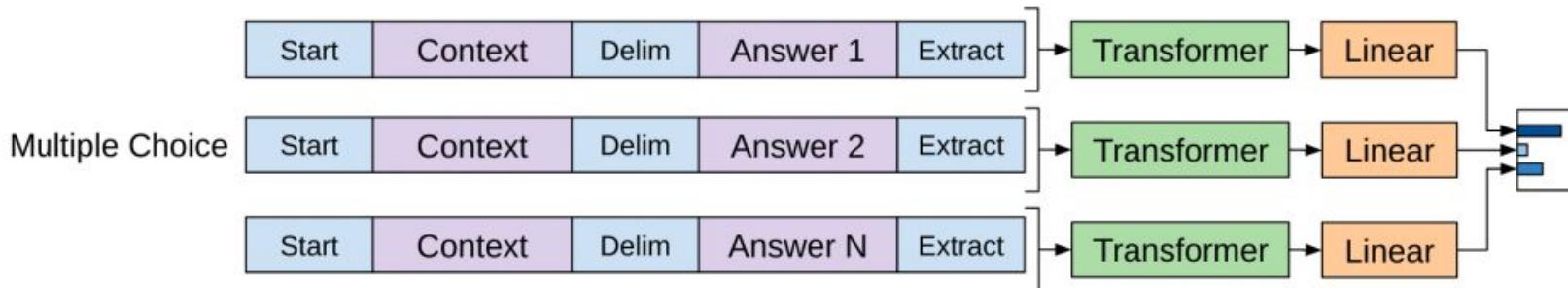
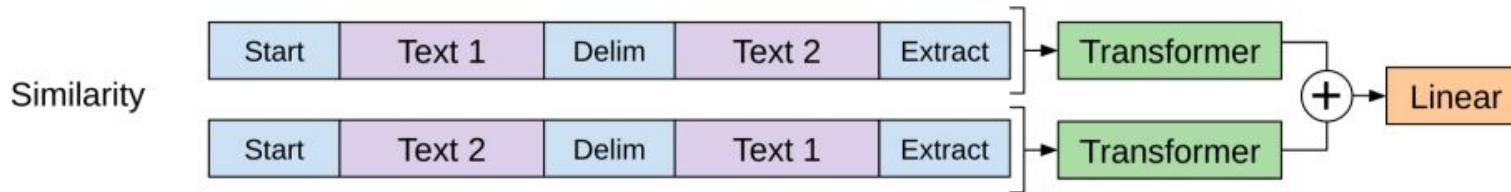
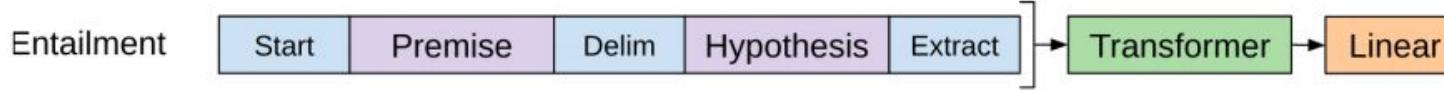
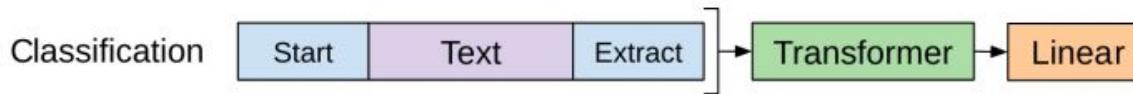
- no encoder
- decoder layers would not have the encoder-decoder attention sublayer
- Pre-train the model on predicting the next word using massive (unlabeled) datasets

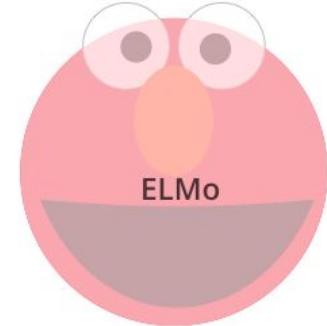
OpenAI Transformer



- During pre-training phase layers have been tuned to reasonably handle language
- Now let's use it for downstream tasks (e.g. sentence classification)

Input transformations for different tasks





GPT-2

- Transformer-based architecture
- trained to predict the **next** word
- 1.5 billion parameters
- Trained on 8 million web-pages



- Transformer-based architecture
- trained to predict the **next** word
- 1.5 billion parameters
- Trained on 8 million web-pages

On language tasks (question answering, reading comprehension, summarization, translation) works well **WITHOUT** fine-tuning

GPT-2: question answering

EXAMPLES

Who wrote the book the origin of species?

Correct answer: *Charles Darwin*

Model answer: Charles Darwin

What is the largest state in the U.S. by land mass?

Correct answer: *Alaska*

Model answer: California

GPT-2: language modeling

EXAMPLE

Both its sun-speckled shade and the cool grass beneath were a welcome respite after the stifling kitchen, and I was glad to relax against the tree's rough, brittle bark and begin my breakfast of buttery, toasted bread and fresh fruit. Even the water was tasty, it was so clean and cold. It almost made up for the lack of...

Correct answer: coffee

Model answer: food

GPT-2: machine translation

EXAMPLE

French sentence:

Un homme a expliqué que l'opération gratuite qu'il avait subie pour soigner une hernie lui permettrait de travailler à nouveau.

Reference translation:

One man explained that the free hernia surgery he'd received will allow him to work again.

Model translation:

A man told me that the operation gratuity he had been promised would not allow him to travel.

New AI fake text generator may be too dangerous to ... - The Guardian

<https://www.theguardian.com/.../elon-musk-backed-ai-writes-convincing-news-fiction>

4 days ago - The Elon Musk-backed nonprofit company OpenAI declines to release research publicly for fear of misuse. The creators of a revolutionary AI system that can write news stories and works of fiction – dubbed "deepfakes for text" – have taken the unusual step of not releasing ...

OpenAI built a text generator so good, it's considered too dangerous to ...

<https://techcrunch.com/2019/02/17/openai-text-generator-dangerous/> ▾

12 hours ago - A storm is brewing over a new language model, built by non-profit artificial intelligence research company OpenAI, which it says is so good at ...

The AI Text Generator That's Too Dangerous to Make Public | WIRED

<https://www.wired.com/story/ai-text-generator-too-dangerous-to-make-public/> ▾

4 days ago - In 2015, car-and-rocket man Elon Musk joined with influential startup backer Sam Altman to put artificial intelligence on a new, more open ...

Elon Musk-backed AI Company Claims It Made a Text Generator ...

<https://gizmodo.com/elon-musk-backed-ai-company-claims-it-made-a-text-gener-183...> ▾

Elon Musk-backed AI Company Claims It Made a Text Generator That's Too Dangerous to Release · Rhett Jones · Friday 12:15pm · Filed to: OpenAI Filed to: ...

Scientists have made an AI that they think is too dangerous to ...

<https://www.weforum.org/.../amazing-new-ai-churns-out-coherent-paragraphs-of-text/> ▾

3 days ago - Sample outputs suggest that the AI system is an extraordinary step forward, producing text rich with context, nuance and even something ...

New AI Fake Text Generator May Be Too Dangerous To ... - Slashdot

<https://news.slashdot.org/.../new-ai-fake-text-generator-may-be-too-dangerous-to-rele...> ▾

3 days ago - An anonymous reader shares a report: The creators of a revolutionary AI system that can write news stories and works of fiction – dubbed ...

GPT-2: fake news and hype

Top stories



OpenAI built a text generator so good, it's considered too dangerous to release

TechCrunch

11 hours ago



Elon Musk's AI company created a fake news generator it's too scared to make public

BGR.com

9 hours ago



The AI That Can Write A Fake News Story From A Handful Of Words

NDTV.com

2 hours ago

When Is Technology Too Dangerous to Release to the Public?

Slate • 2 days ago



Scientists Developed an AI So Advanced They Say It's Too Dangerous to Release

ScienceAlert • 6 days ago





- 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