

NLP 102

Transformer

Тарас Хахулин

Skoltech, MIPT

Deep Learning Engineer, Samsung AI Center

Tg: @vitaminotar

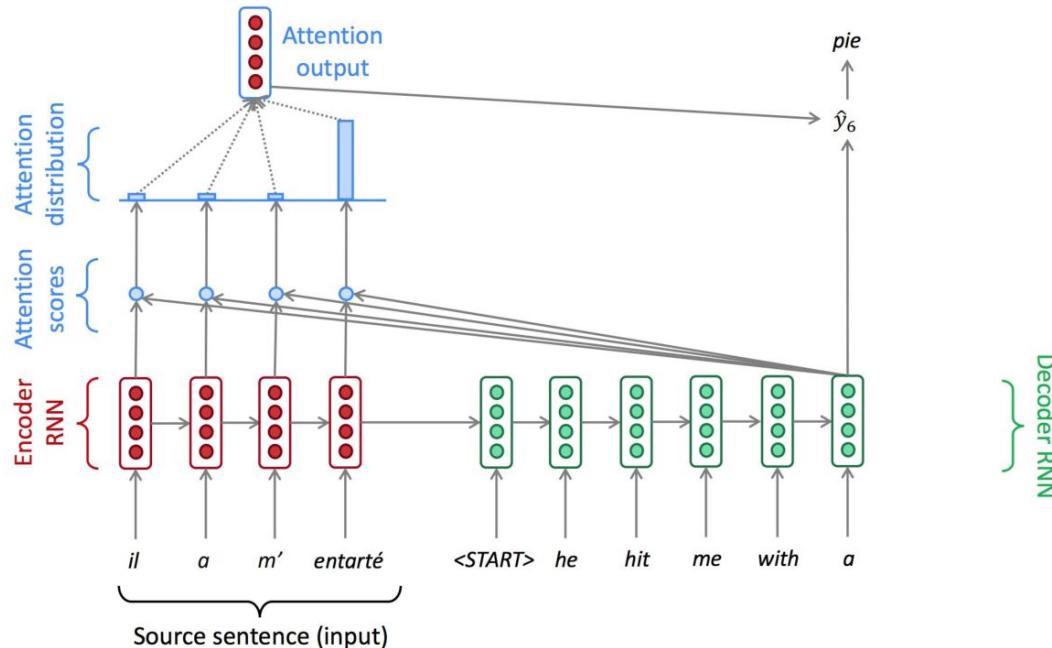
<https://github.com/khakhulin/>

Overview

- Attention
- Global view on the transformer
- Self-Attention
- Positional Encoding
- Layer Normalization
- Questions



recap: Attention in seq2seq

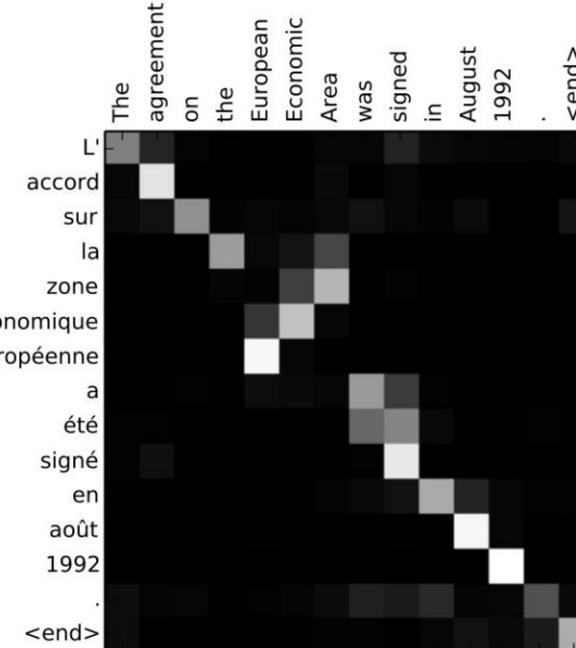
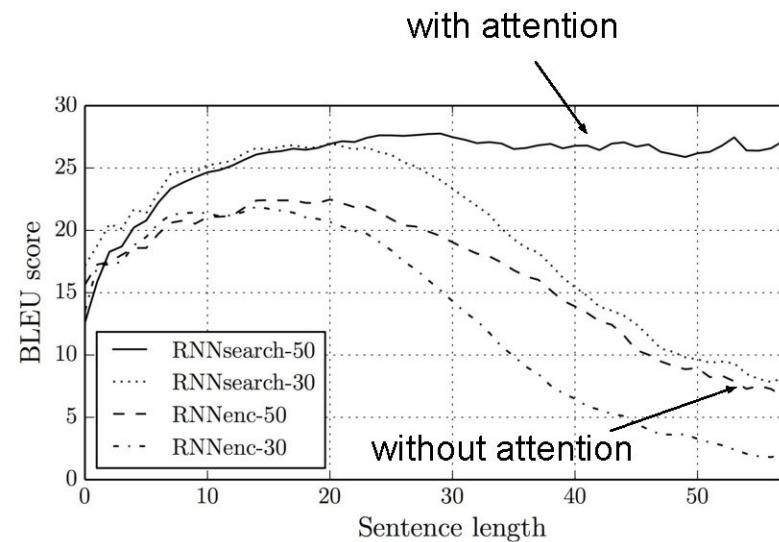


Attention variants

- 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

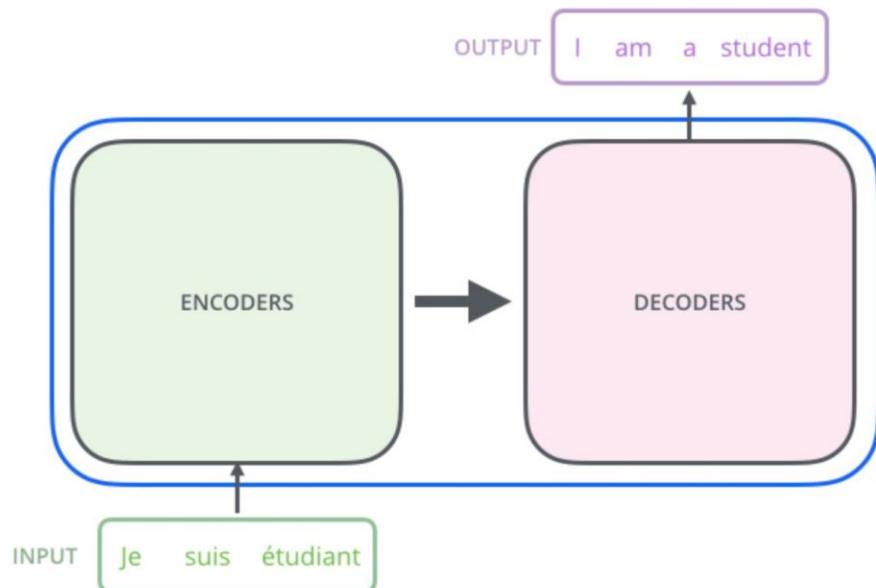


The Transformer

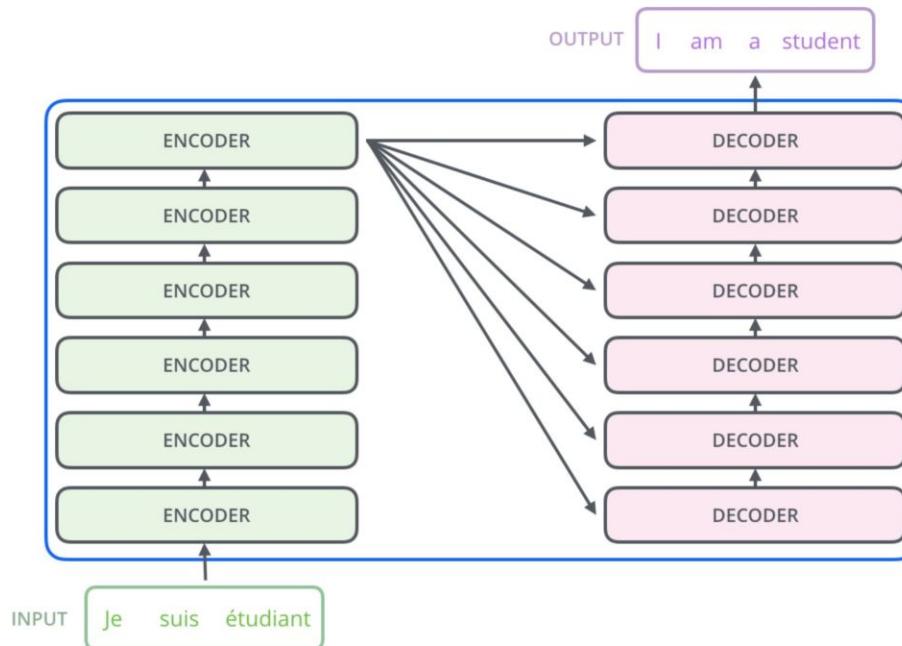
The Transformer



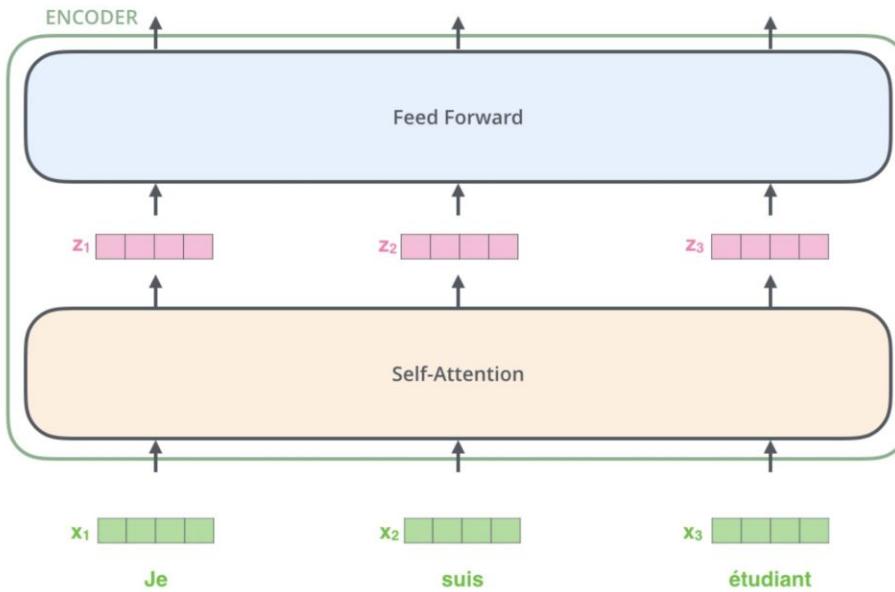
The Transformer



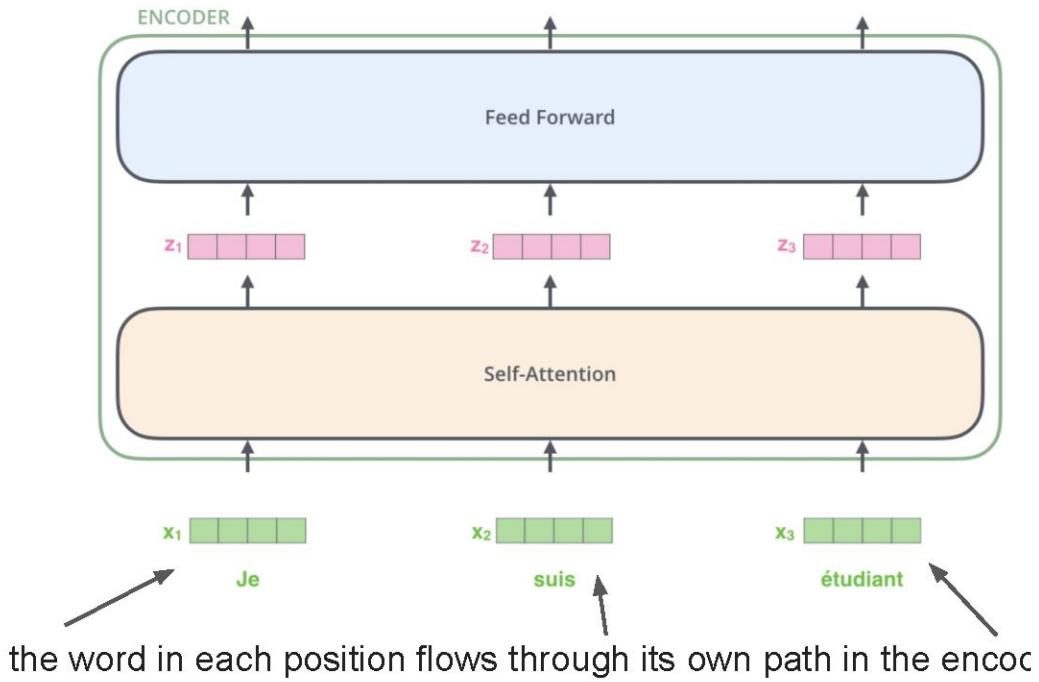
The Transformer



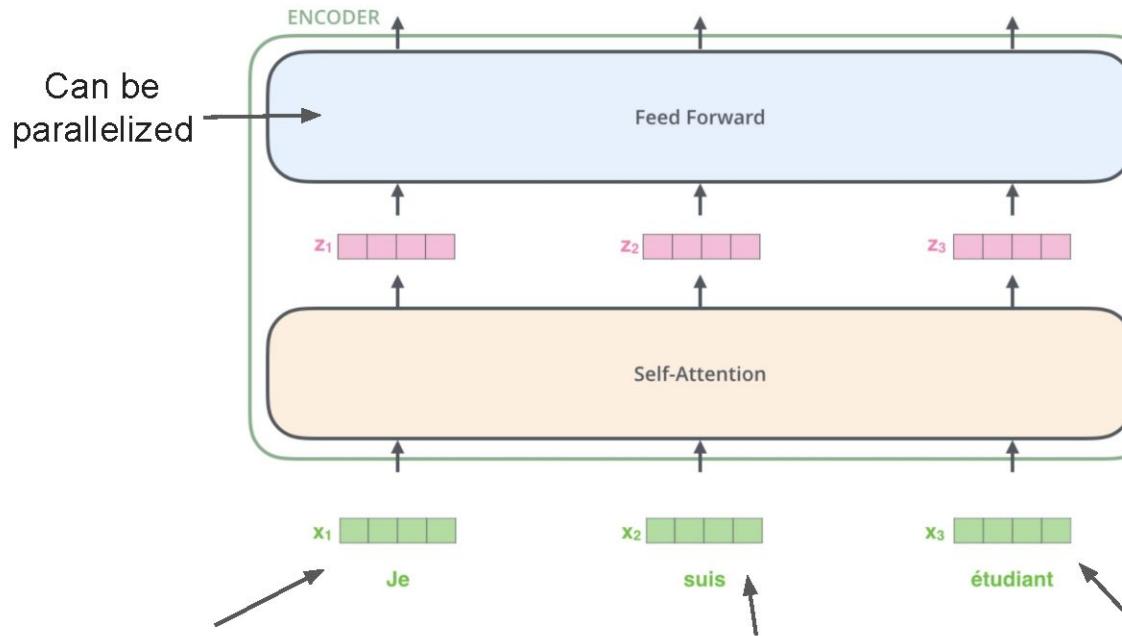
The Encoder Side



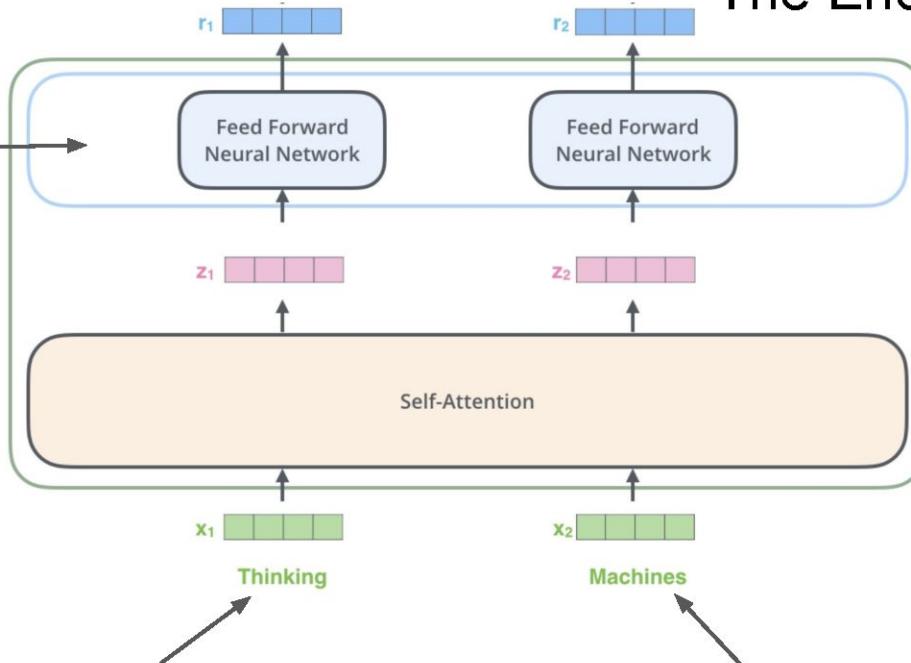
The Encoder Side



The Encoder Side



The Encoder Side



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

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

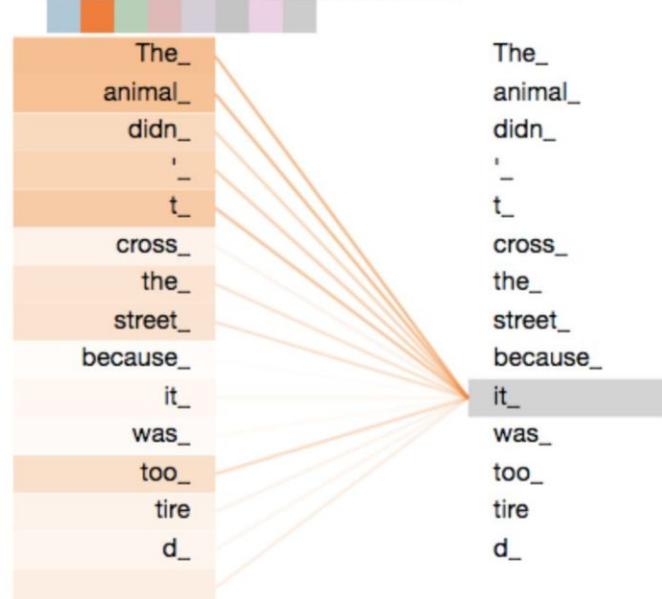
"The animal didn't cross the street because it was too tired"

- 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

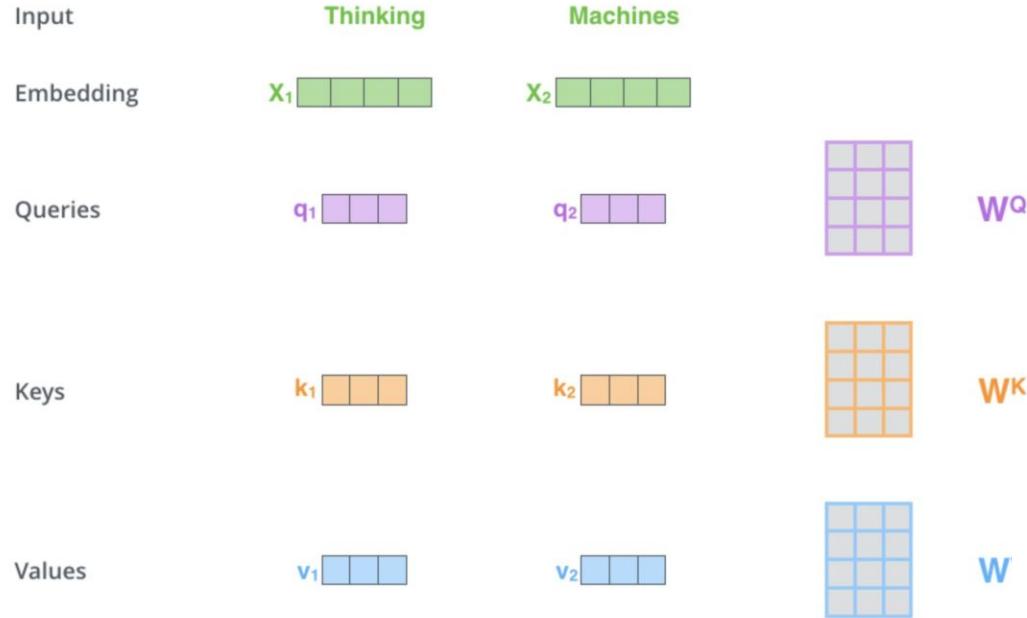
Layer: 5 Attention: Input - Input



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

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

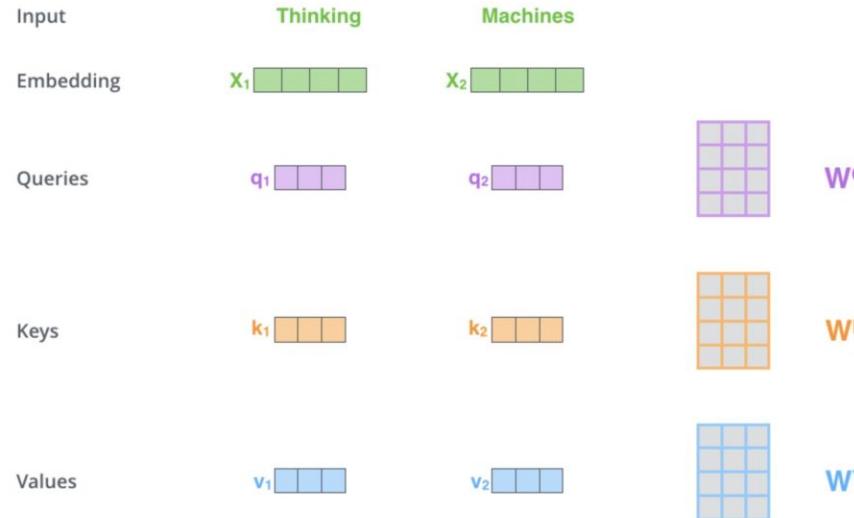
Self-Attention: detailed explanation



Self-Attention: detailed explanation

STEP 1:

create 3 vectors
(Query, Key, Value)
from each of the encoder's
input vectors



Self-Attention: detailed explanation

What are the “query”, “key”, and “value” vectors?

Self-Attention: detailed explanation

What are the “query”, “key”, and “value” vectors?

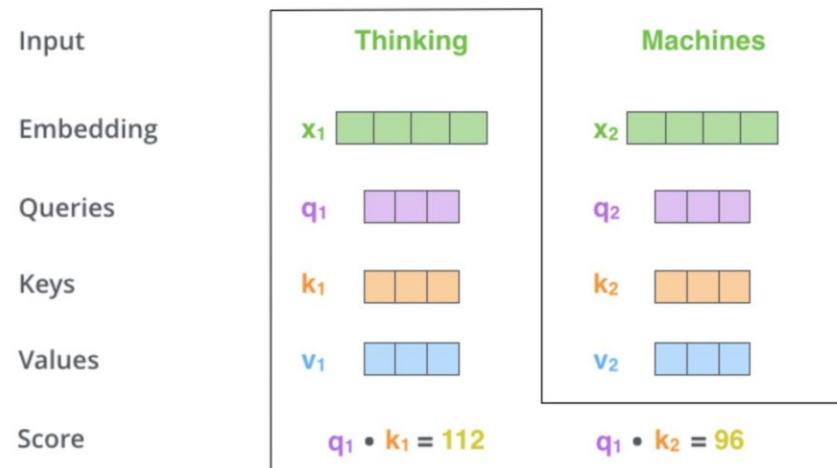
They’re abstractions that are useful for
calculating and thinking about attention.

Self-Attention: detailed explanation

STEP 2:

calculate a score

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



Self-Attention: detailed explanation

STEP 3:

divide the scores by 8

Queries
Keys
(the square root of the
dimension of the key vectors) Values

STEP 4:

softmax

Input

Embedding

Queries

Keys

(the square root of the
dimension of the key vectors)

Score

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

Softmax

Thinking

x_1

q_1

k_1

v_1

Machines

x_2

q_2

k_2

v_2

$$q_1 \cdot k_1 = 112$$

$$14$$

$$0.88$$

$$q_1 \cdot k_2 = 96$$

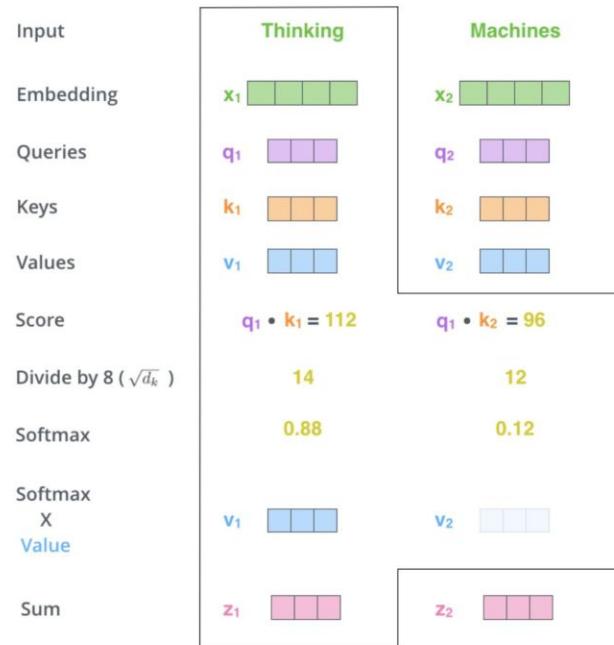
$$12$$

$$0.12$$

Self-Attention: detailed explanation

STEP 5:

multiply each value vector by the softmax score

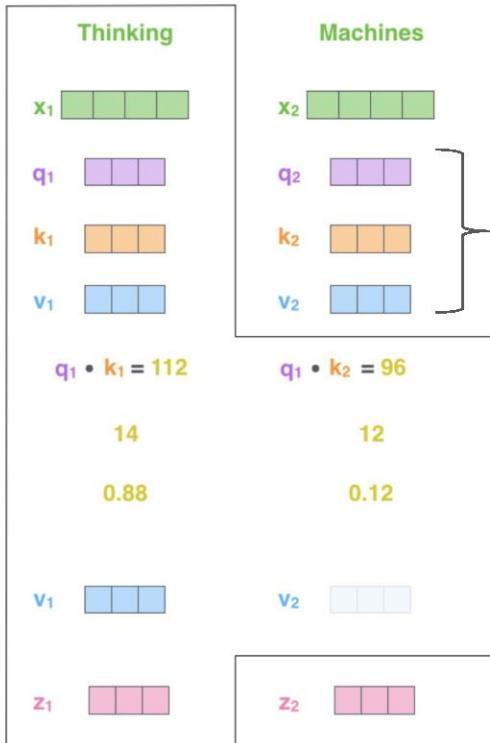


STEP 6:

sum up the weighted value vectors

Self-Attention

Input
Embedding
Queries
Keys
Values
Score
Divide by 8 ($\sqrt{d_k}$)
Softmax
Softmax X Value
Sum



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 \mathbf{X}

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

Multiply \mathbf{X} by weight matrices we've trained (\mathbf{W}_k , \mathbf{W}_q , \mathbf{W}_v)

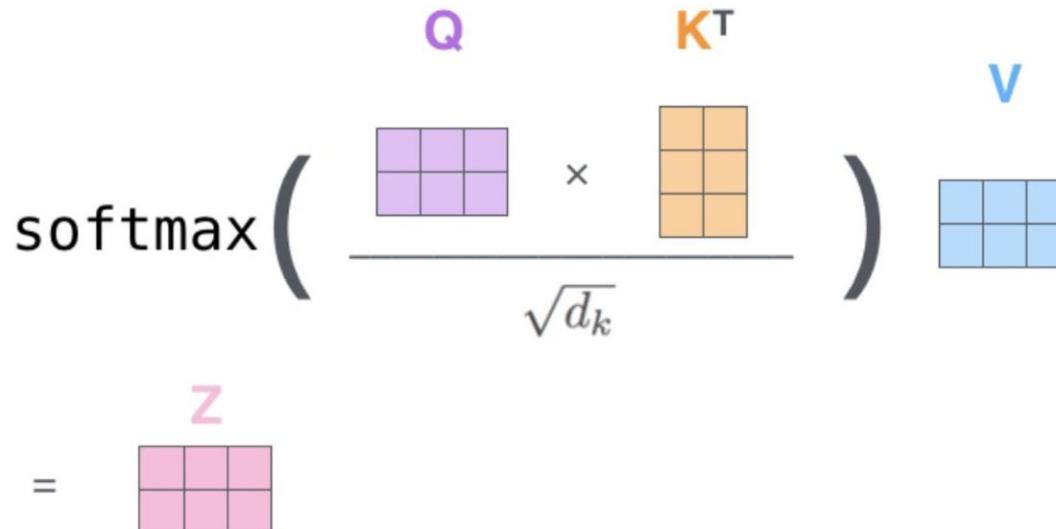
$$\mathbf{X} \quad \times \quad \mathbf{W}^K \quad = \quad \mathbf{K}$$

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

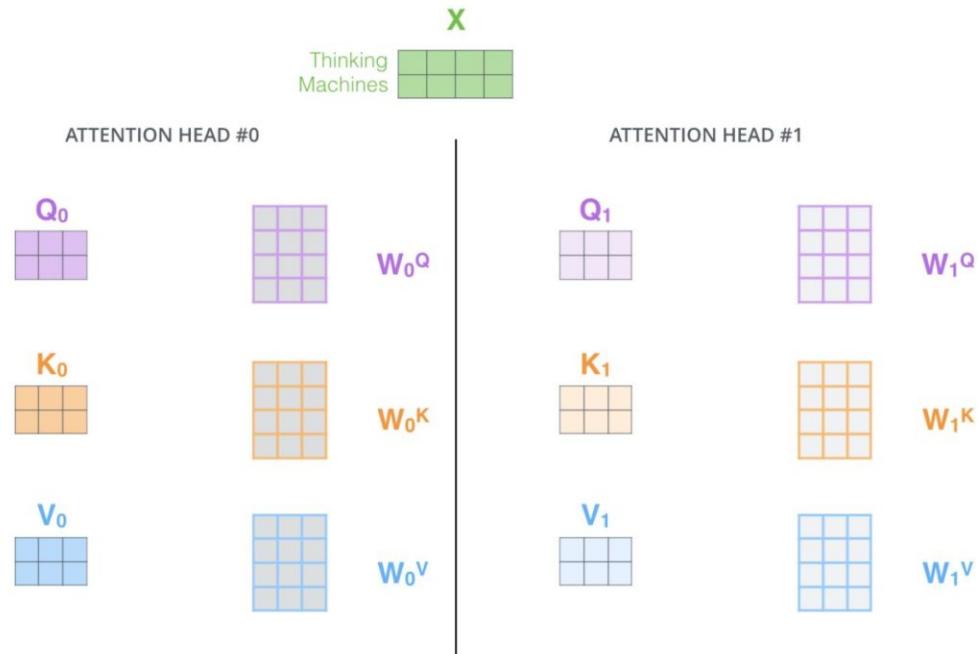
Self-Attention: Matrix Calculation

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

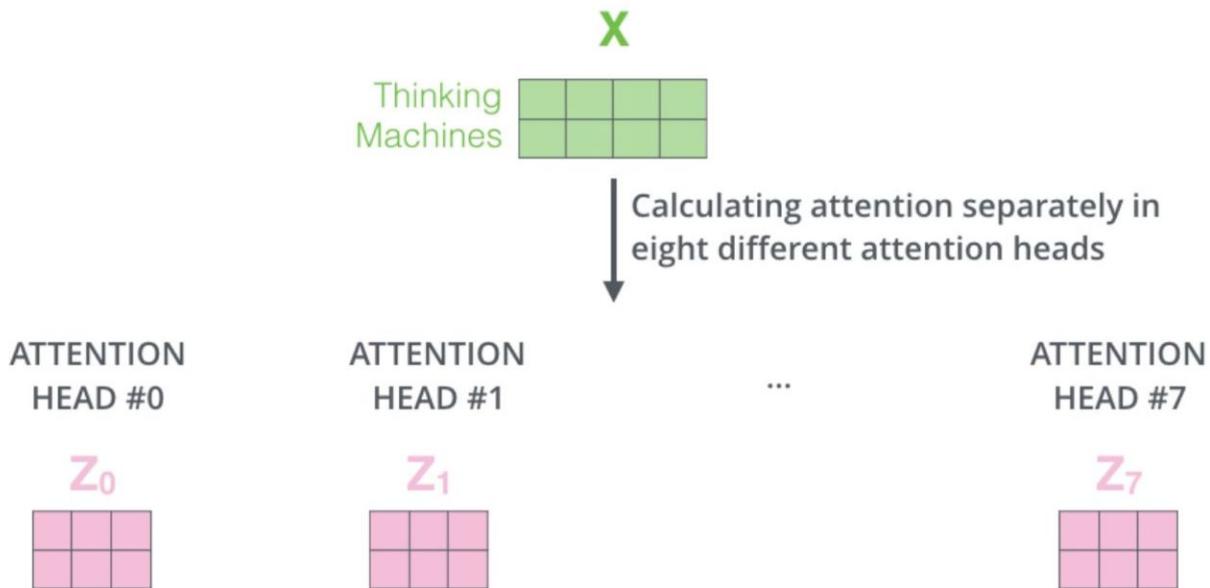
$=$ \mathbf{Z}



Multi-Head Attention



Multi-Head Attention



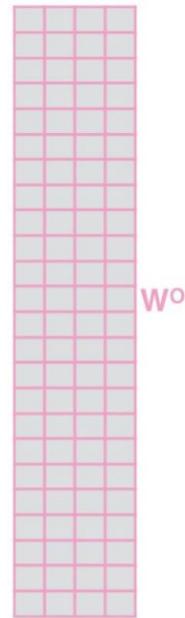
Multi-Head Attention

1) Concatenate all the attention heads



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

X

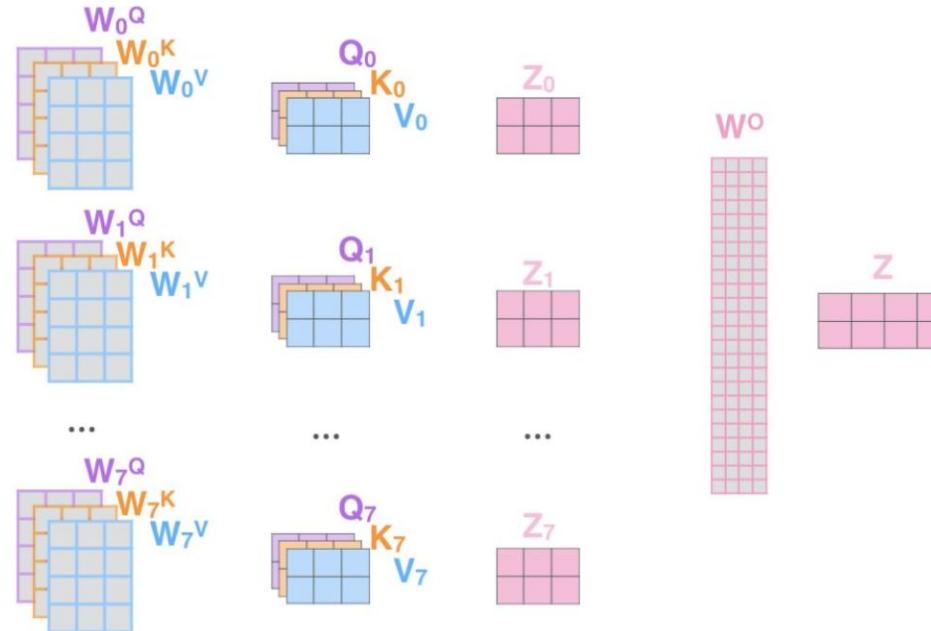


3) The result would be the Z matrix that captures information from all the attention heads. We can send this forward to the FFNN

$$= \begin{matrix} Z \\ \hline \end{matrix}$$

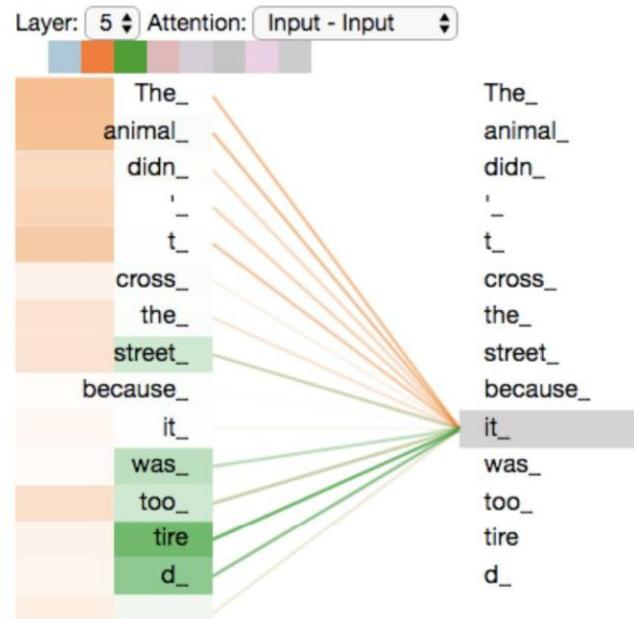
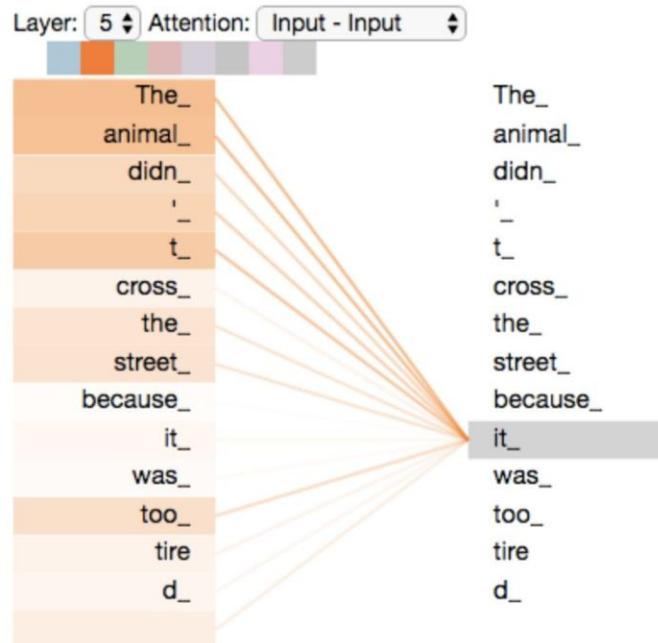
- 1) This is our input sentence*
- 2) We embed each word*
- 3) Split into 8 heads. We multiply X or R with weight matrices
- 4) Calculate attention using the resulting $Q/K/V$ matrices
- 5) Concatenate the resulting Z matrices, then multiply with weight matrix W^o to produce the output of the layer

Thinking
Machines

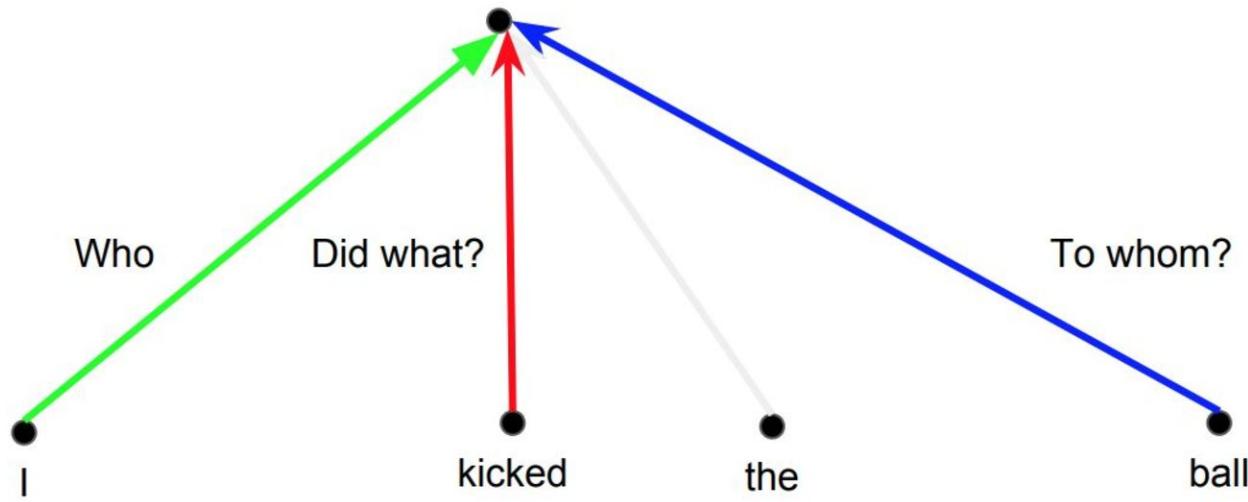


* In all encoders other than #0, we don't need embedding. We start directly with the output of the encoder right below this one

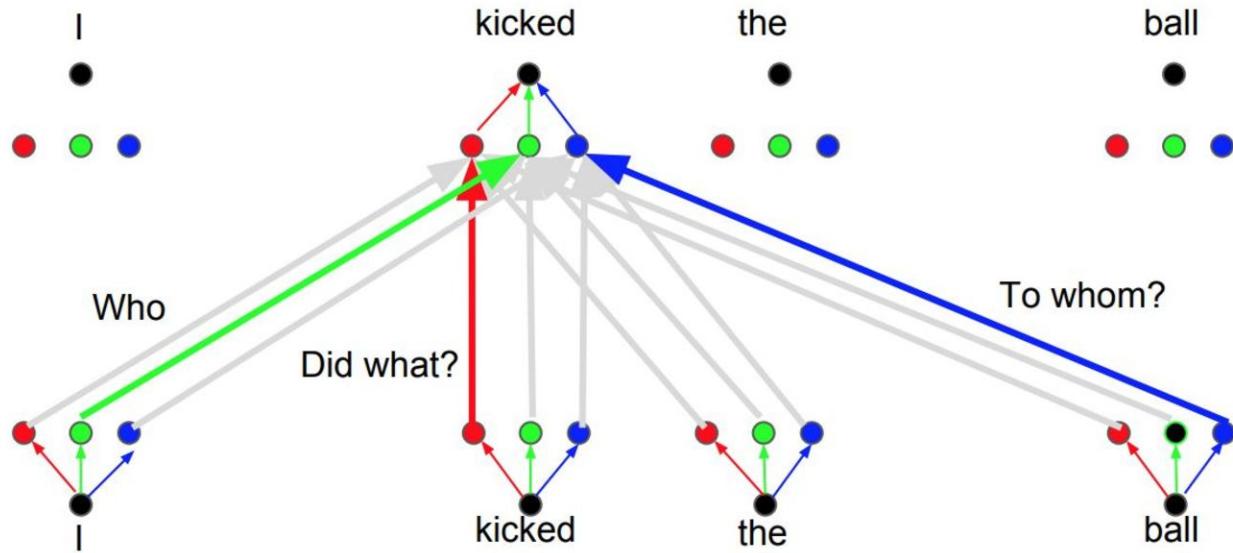
Multi-Head Attention



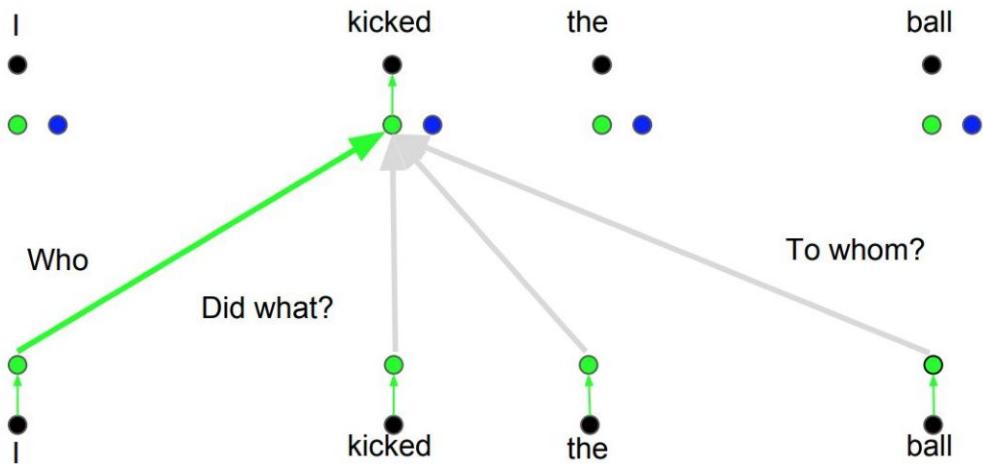
Why Multi-Head Attention?



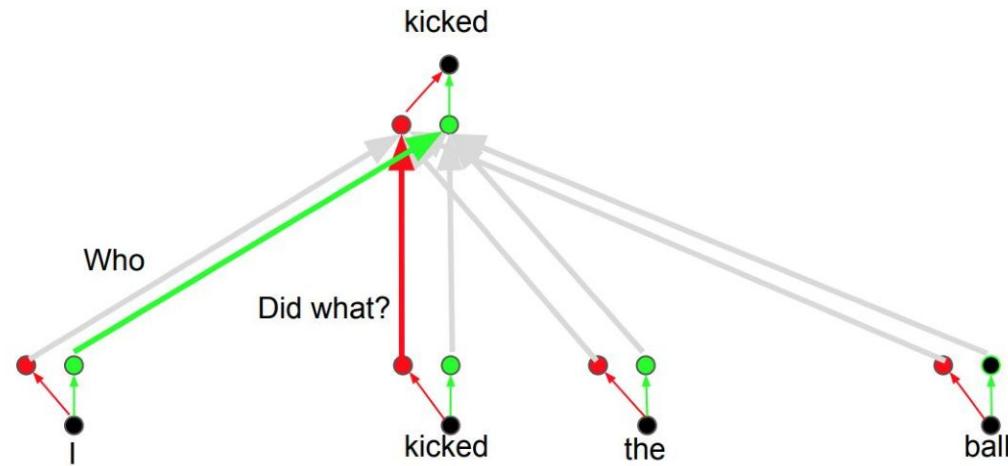
Why Multi-Head Attention?



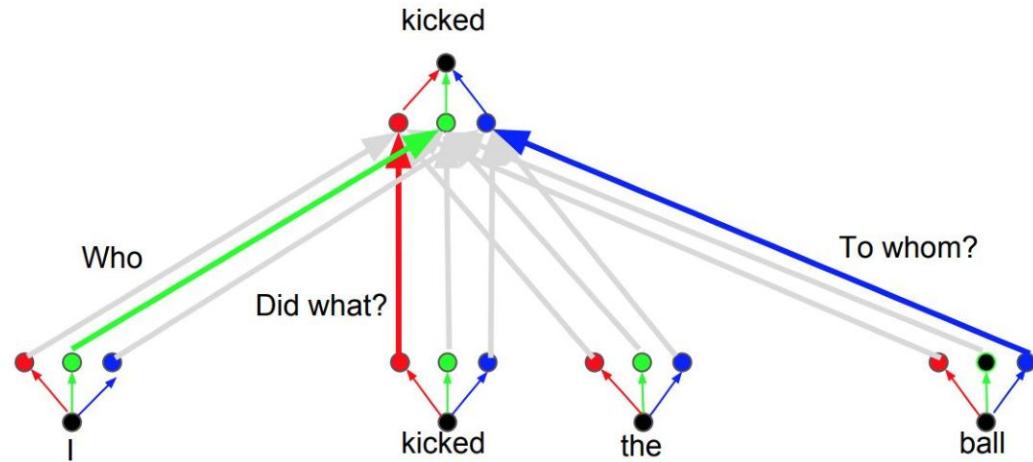
Attention head: Who



Attention head: Did What?

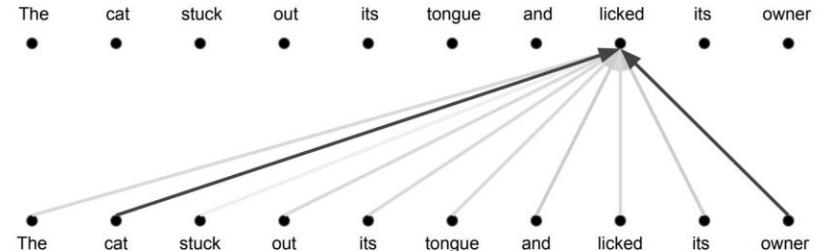


Attention head: To Whom?



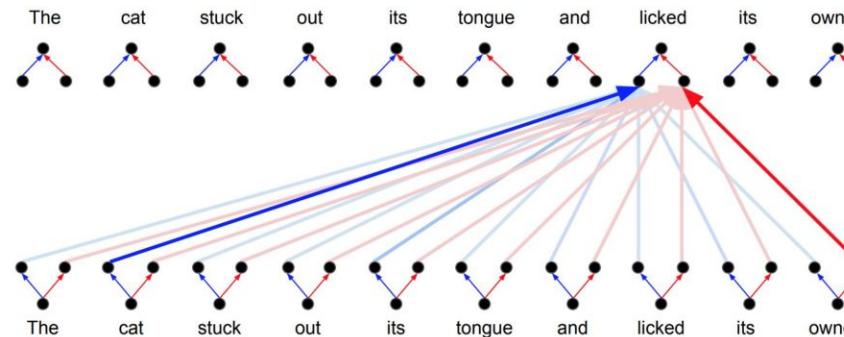
Attention vs. Multi-Head Attention

Attention: a weighted average



Multi-Head Attention:

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



Performance: WMT 2014 BLEU

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

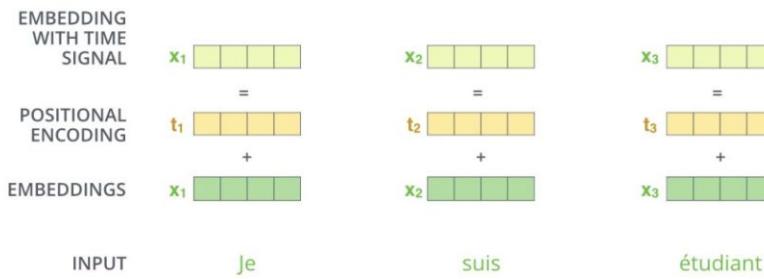
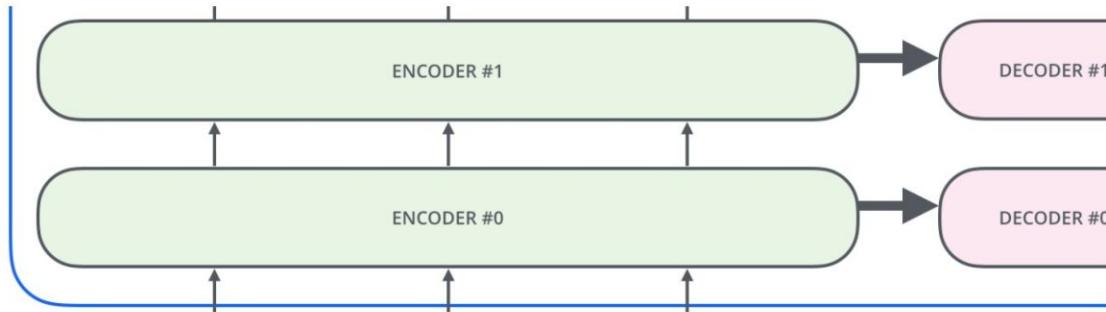
*Transformer models trained >3x faster than the others.

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

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

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

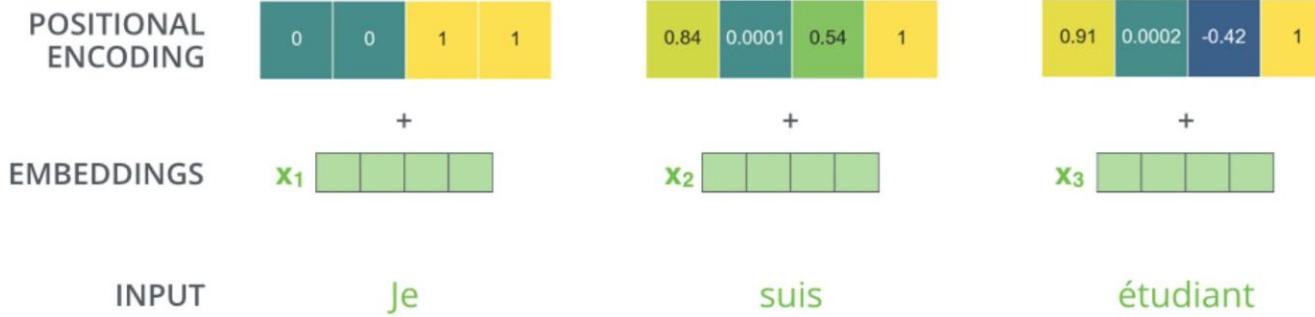
$$PE_{(pos, \ 2i + 1)} = \cos(pos / 10000^{2i / d_{\text{model}}})$$

- pos is the position
- i is the dimension.

Each dimension of the positional encoding corresponds to a sinusoid.

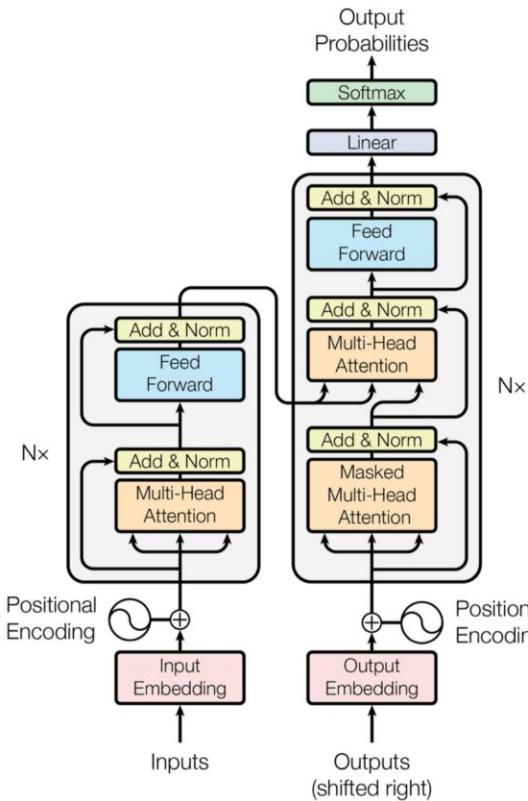
The wavelengths form a geometric progression from 2π to $2\pi * 10000$

Positional Encoding

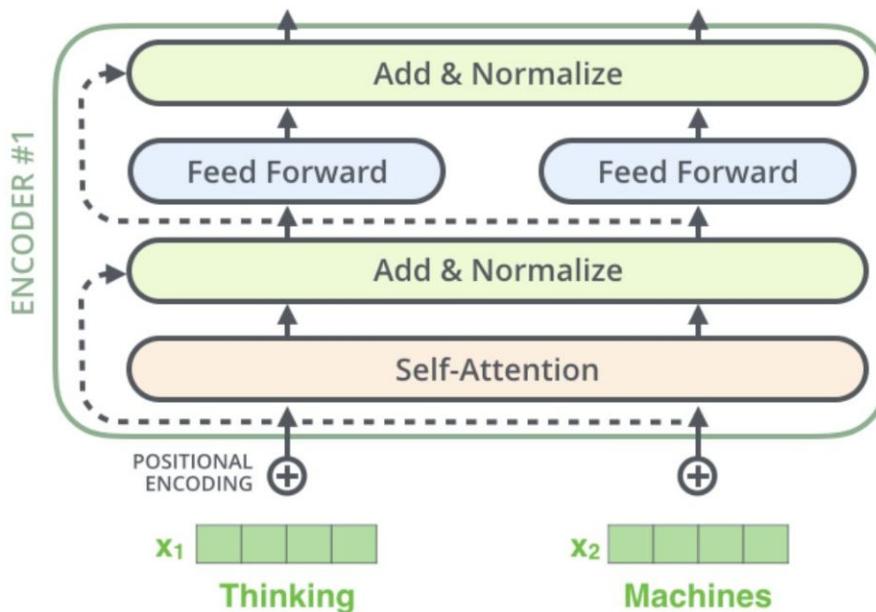


Layer Normalization

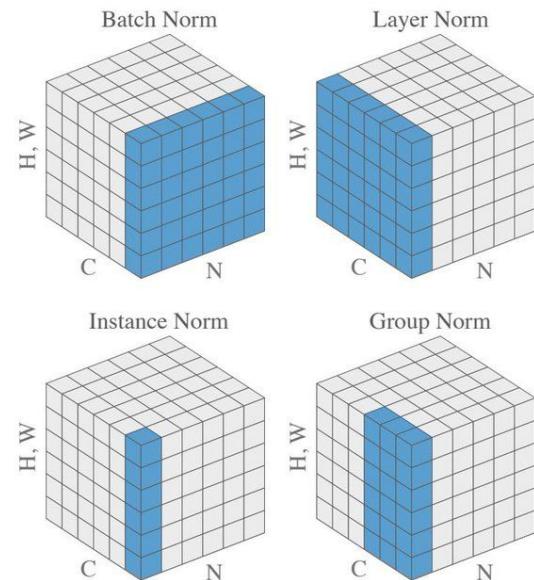
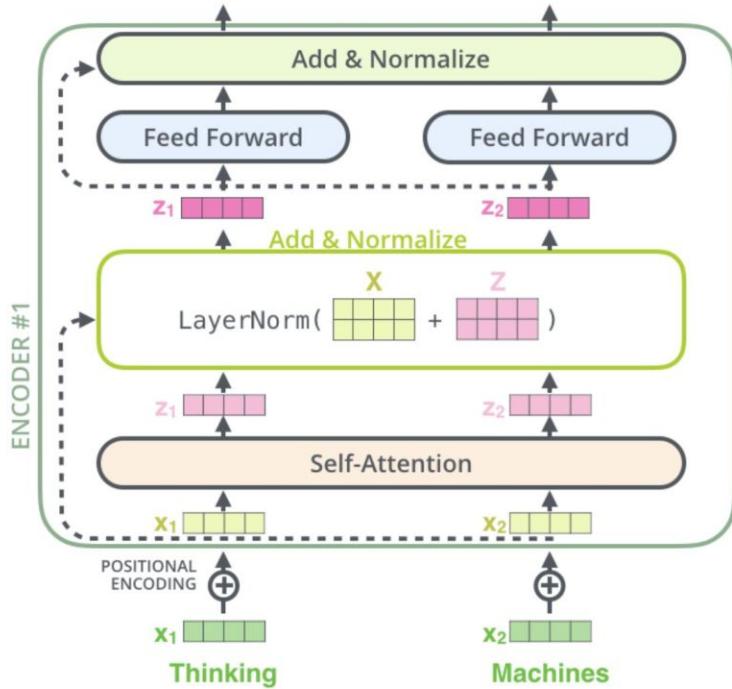
The Transformer: recap



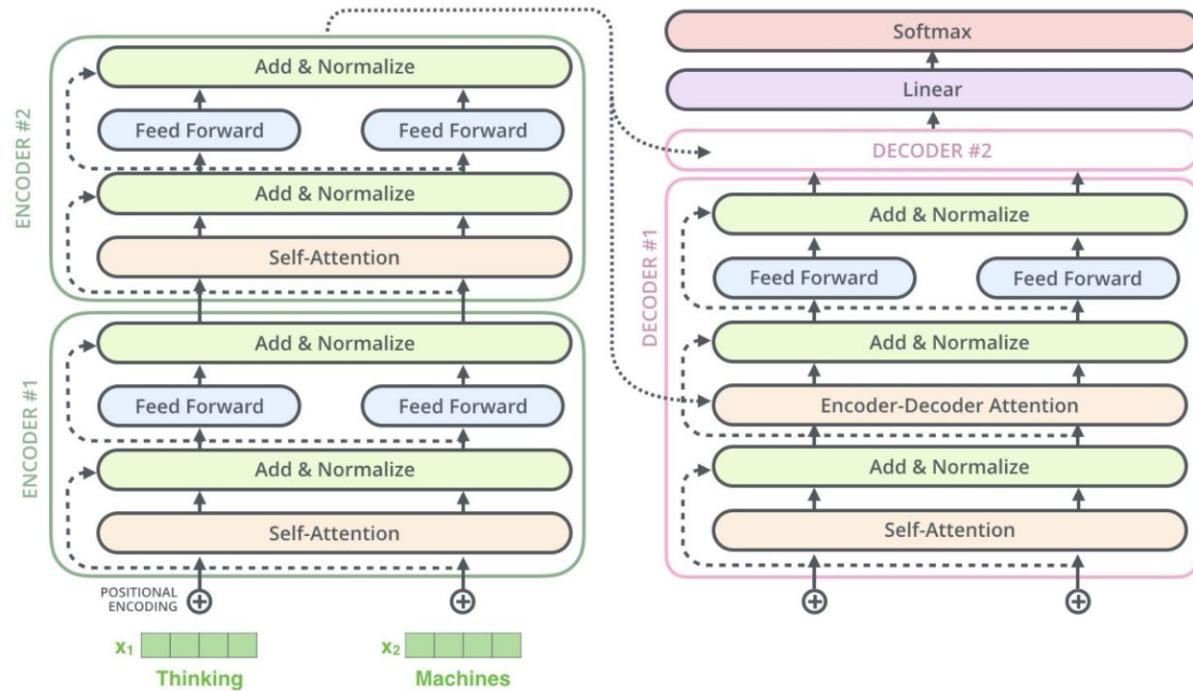
Layer Normalization



Layer Normalization

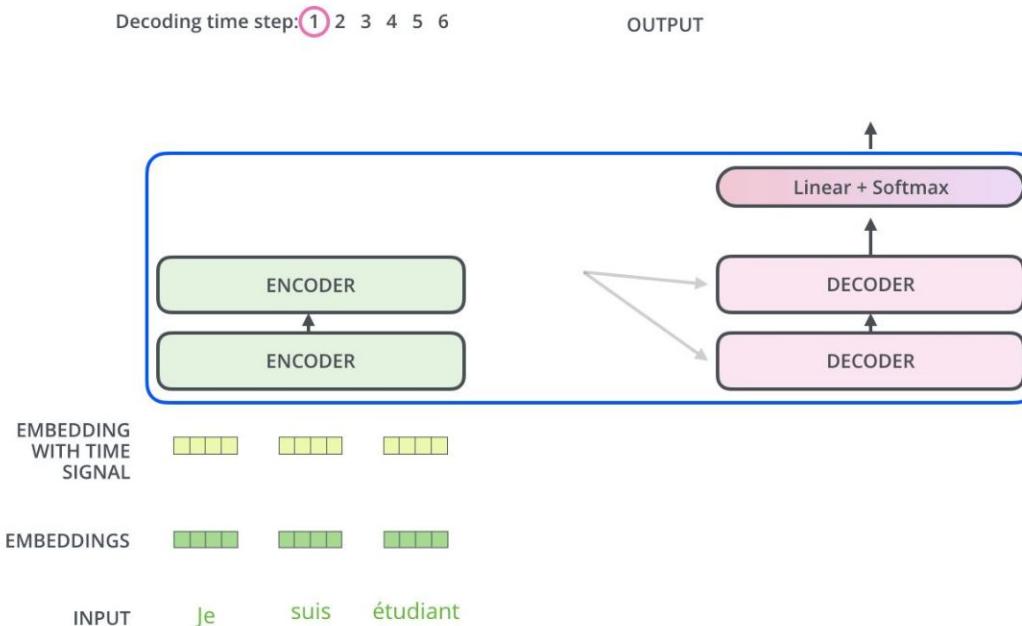


Layer Normalization

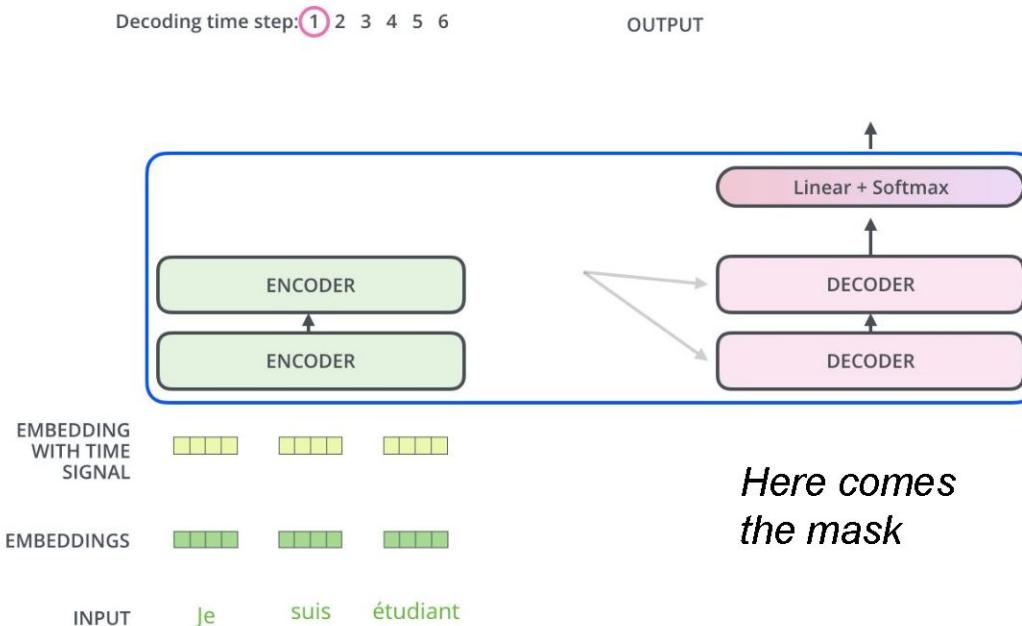


The Decoder

The Decoder Side



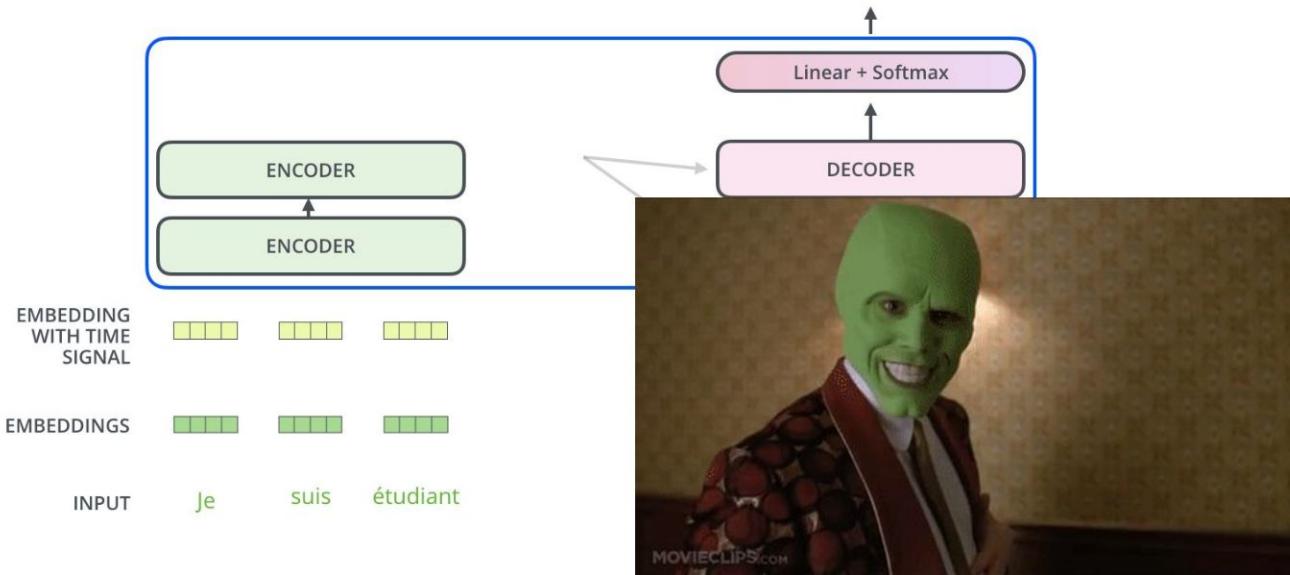
The Decoder Side



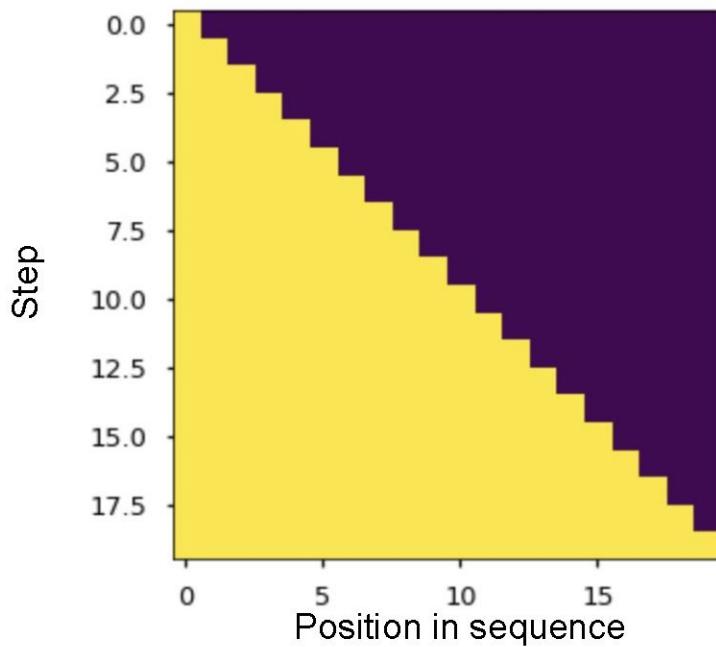
The Decoder Side

Decoding time step: ① 2 3 4 5 6

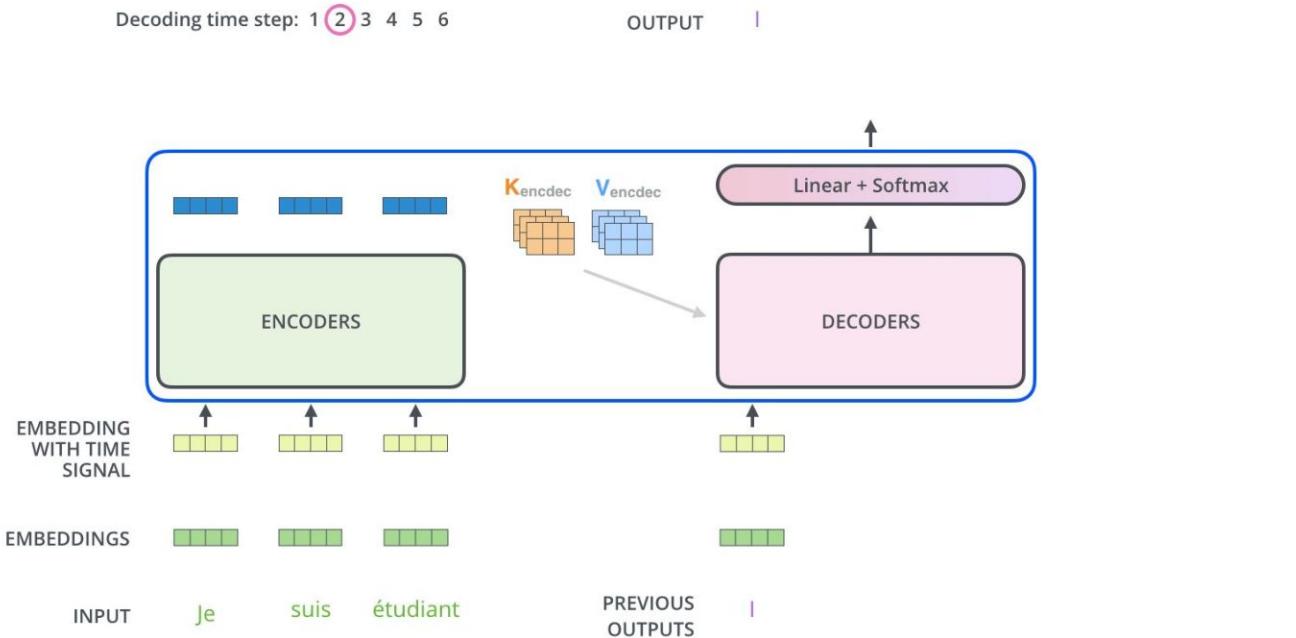
OUTPUT



The masked decoder input



The Decoder Side



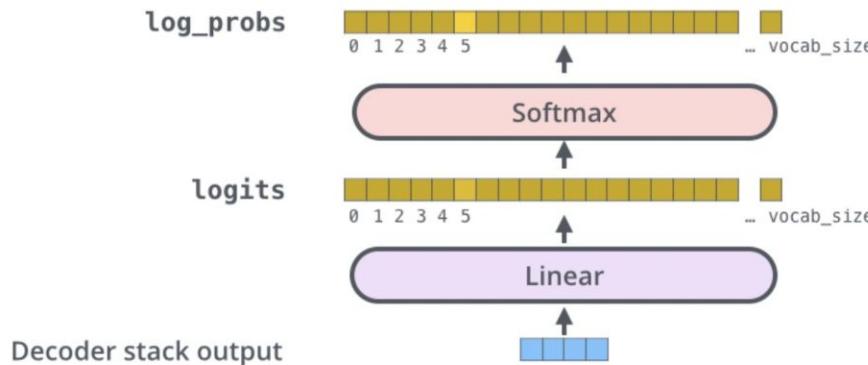
Final Linear and Softmax Layer

Which word in our vocabulary
is associated with this index?

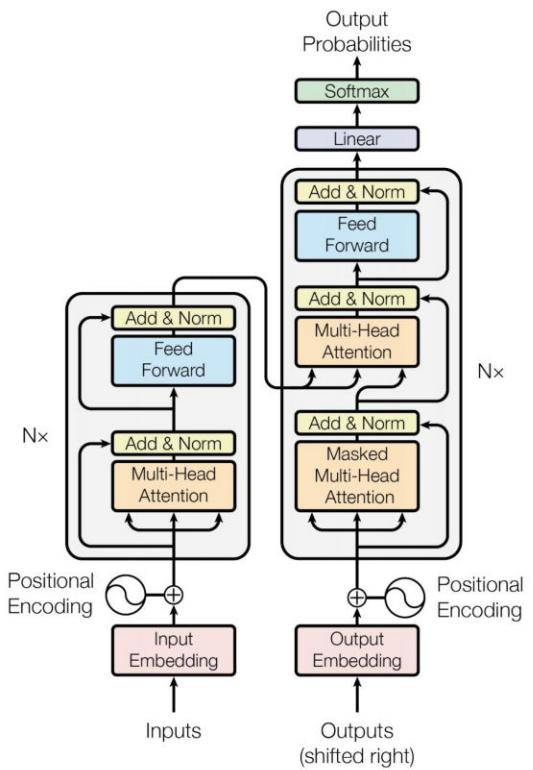
am

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

5



The Transformer

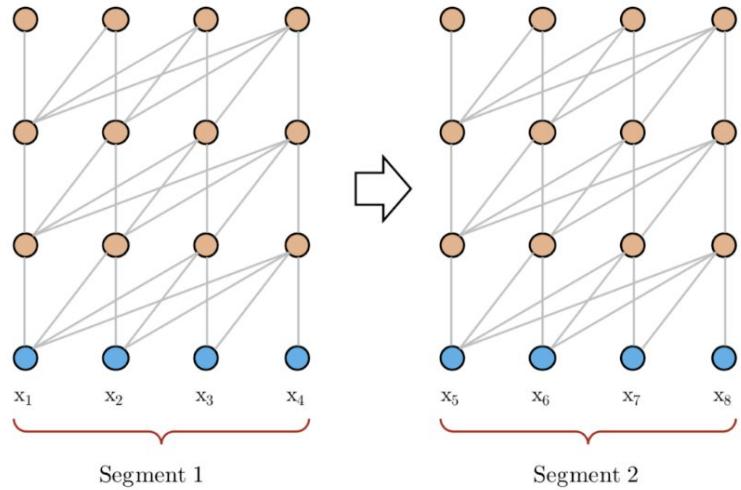


Improved Transformers

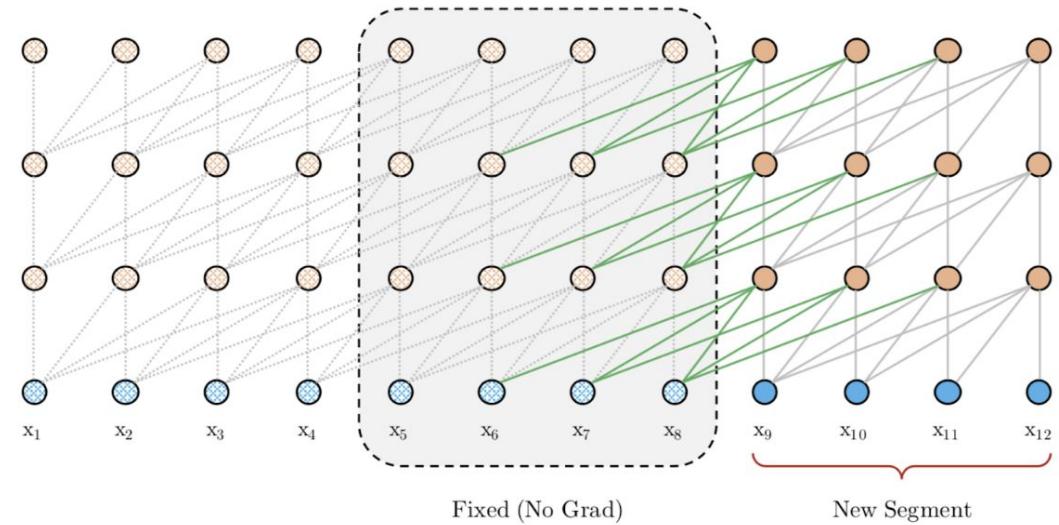


Transformer-XL

Transformer (Training)



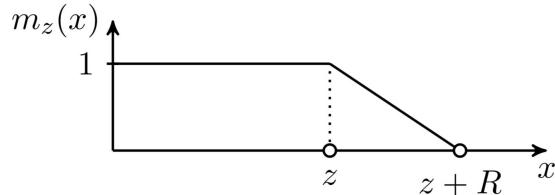
Transformer-XL (Training)



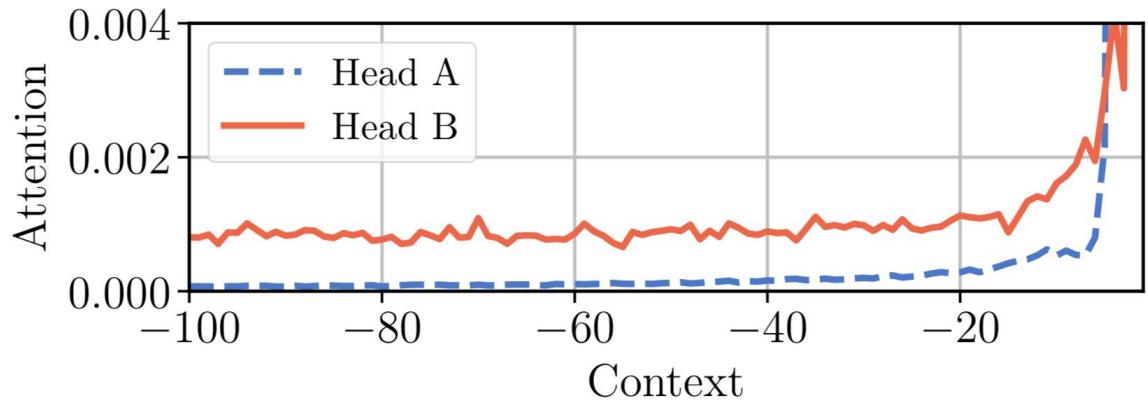
A comparison between the training phrase of vanilla Transformer & Transformer-XL
[Transformer-XL paper](#)

Adaptive Attention Span

$$m_z(x) = \text{clamp}\left(\frac{1}{R}(R + z - x), 0, 1\right)$$



$$a_{ij} = \frac{m_z(i - j) \exp(s_{ij})}{\sum_{r=i-s}^{i-1} m_z(i - r) \exp(s_{ir})}$$



Two attention heads in the same model, A & B, assign attention differently within the same context window. Head A attends more to the recent tokens, while head B look further back into the past uniformly.

[Adaptive attention Span paper](#)

Reformer

[Reformer Paper](#)

Problems

- Memory in a model with N layers is N -times larger than in a single-layer model because we need to store activations for back-propagation.
- The intermediate FF layers are often quite large.
- The attention matrix on sequences of length L often requires $O(L^2)$ in both memory and time.

Sorting by hash

	q_1	q_2	q_3	q_4	q_5	q_6
k_1	•	•		•		
k_2		•				•
k_3				•		
k_4				•		
k_5				•		
k_6		•				•

(a) Normal

	q_1	q_2	q_4	q_3	q_6	q_5
k_1	•	•	•			
k_2				•	•	
k_6				•	•	
k_3				•		
k_4					•	
k_5						•
k_6						•

(b) Bucketed

shared-QK

	q_1	q_2	q_4	q_3	q_6	q_5
q_1	•		•	•		
q_2		•		•		
q_4			•	•		
q_3				•	•	
q_6					•	
q_5						•

(c) $Q = K$

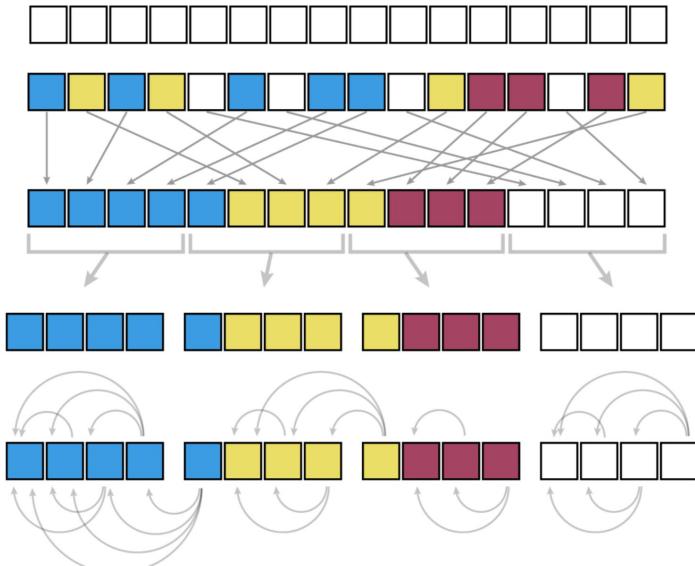
	q_1	q_2	q_4	q_3	q_6	q_5
q_1	•		•	•		
q_2		•		•		
q_4			•	•		
q_3				•	•	
q_6					•	
q_5						•

(d) Chunked

Reformer

[Reformer Paper](#)

Sequence of queries=keys
LSH bucketing
Sort by LSH bucket
Chunk sorted sequence to parallelize
Attend within same bucket in own chunk and previous chunk



Solution

- Replace the dot-product attention with locality-sensitive hashing (LSH) attention, reducing the complexity from $O(L^2)$ to $O(L \log L)$.
- Replace the standard residual blocks with reversible residual layers, which allows storing activations only once during training instead of N times (i.e. proportional to the number of layers).

Reversible Residual Network?

$$y_1 = x_1 + F(x_2), \quad y_2 = x_2 + G(y_1)$$
$$x_2 = y_2 - G(y_1), \quad x_1 = y_1 - F(x_2)$$

$$Y_1 = X_1 + \text{Attention}(X_2), \quad Y_2 = X_2 + \text{FeedForward}(Y_1)$$

The resulting reversible Transformer does not need to store activation in every layer.

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