

These slides are almost the exact copy of [ML-MIPT course](#). Special thanks to ML-MIPT team.

Attention is All You Need

Based on:

https://github.com/qirafe-ai/ml-mipt/blob/master/week1_04_Transformer/week04_Transformer.pdf

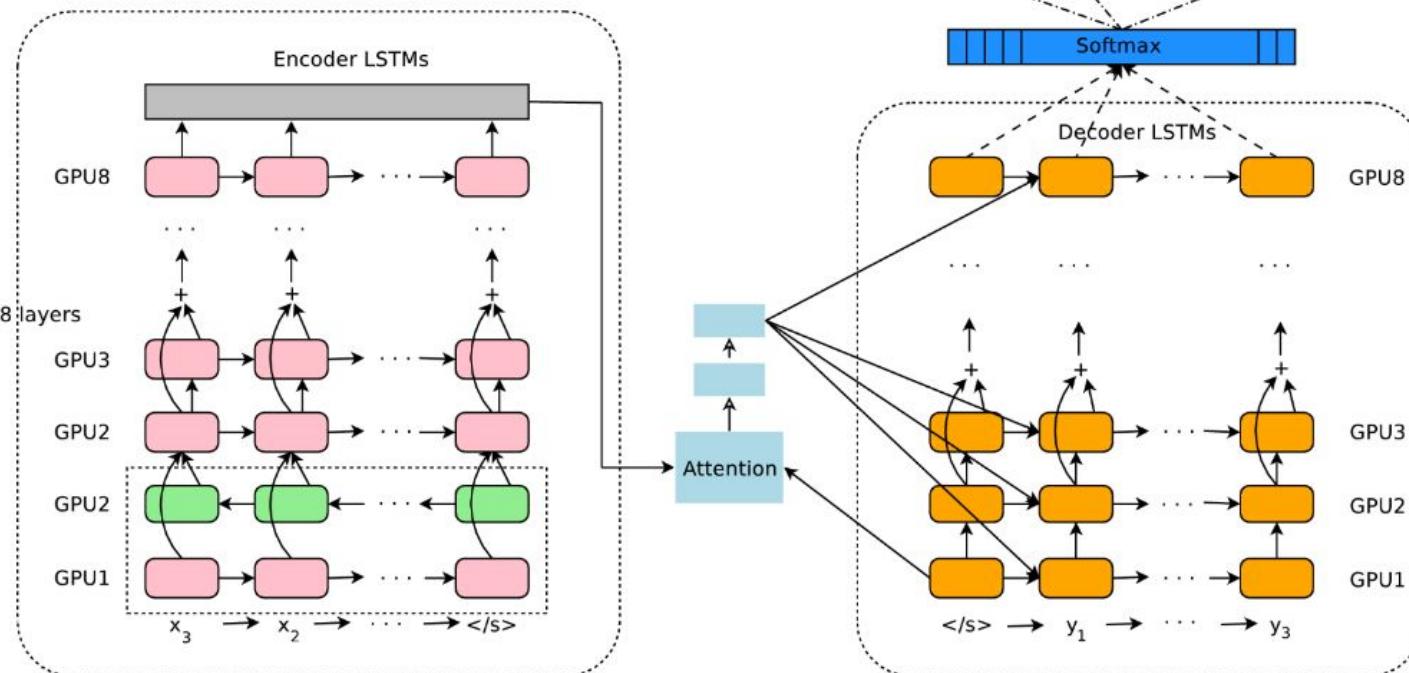
<http://web.stanford.edu/class/cs224n/slides/cs224n-2019-lecture08-nmt.pdf>

<https://jalamar.github.io/illustrated-transformer/>

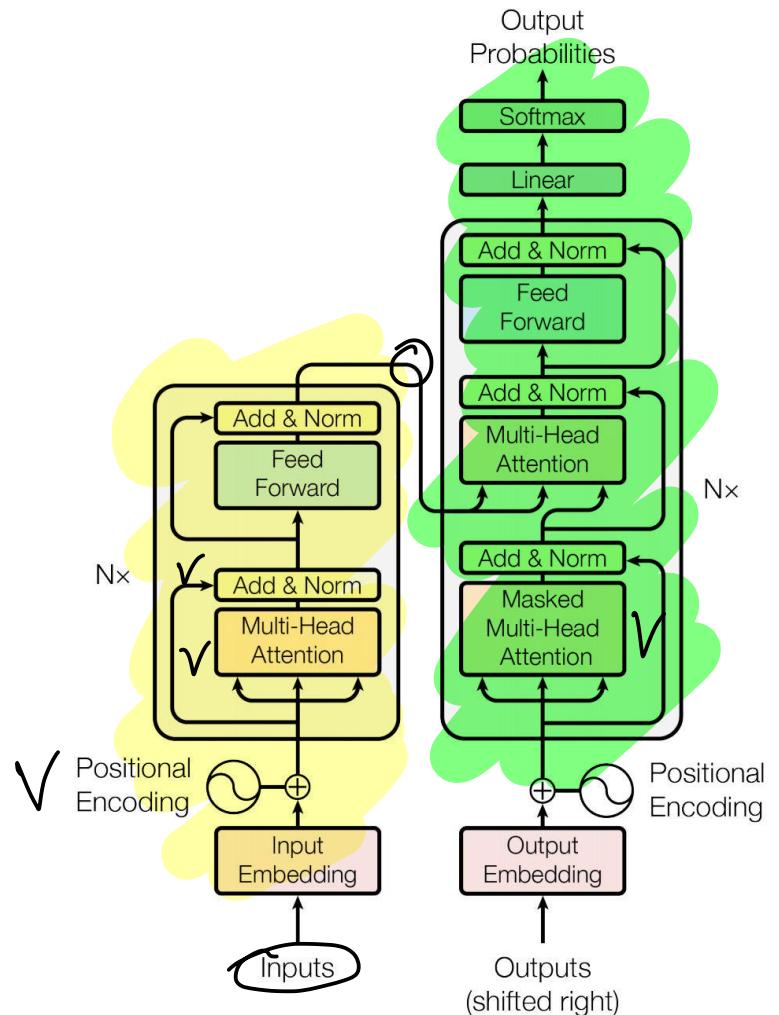
https://github.com/yandexdataschool/nlp_course

Deep Encoder-Decoder Models (GNMT)

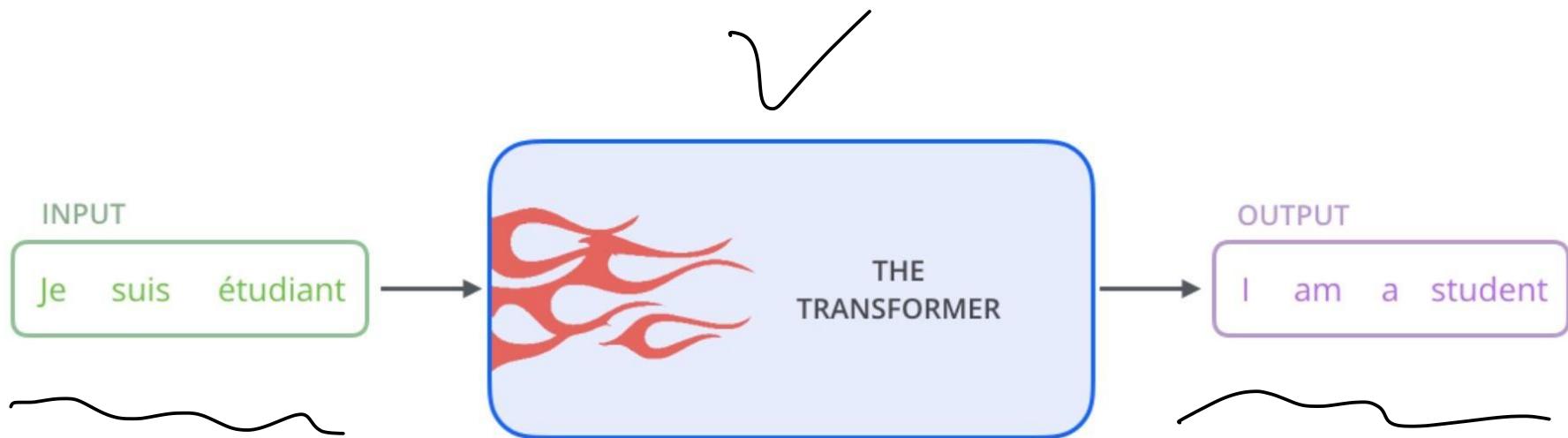
Wu et al. 2016



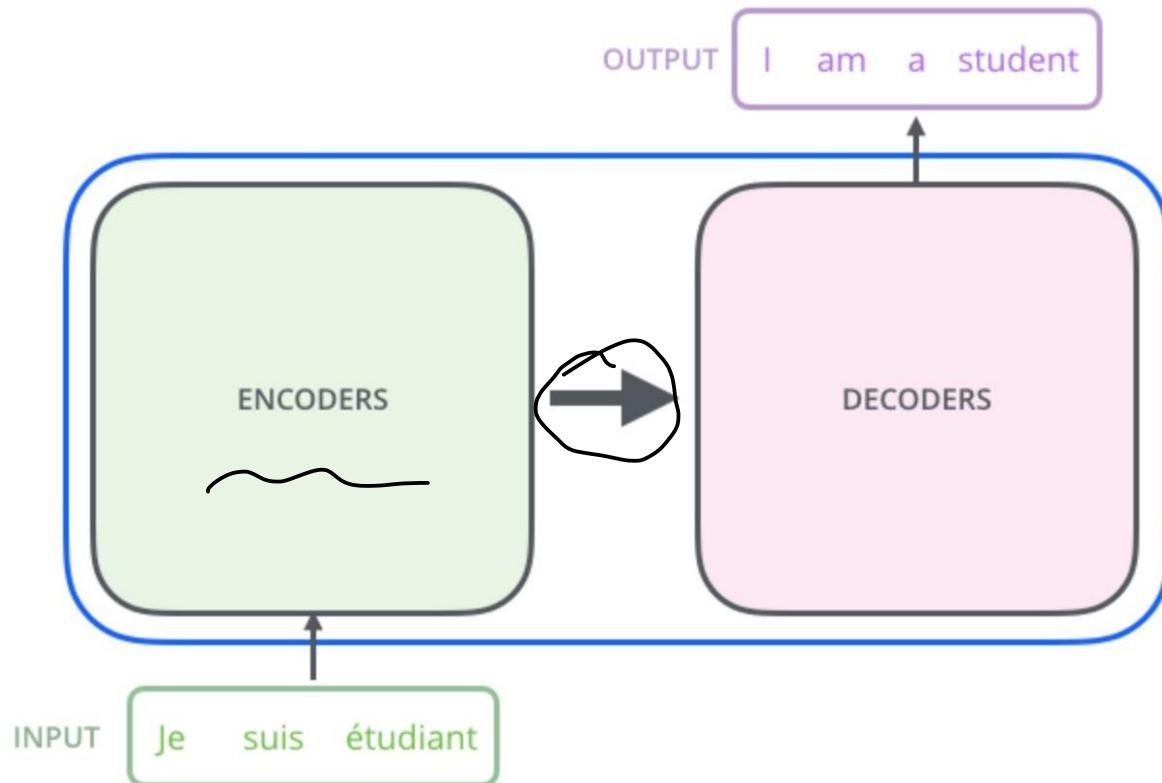
The Transformer



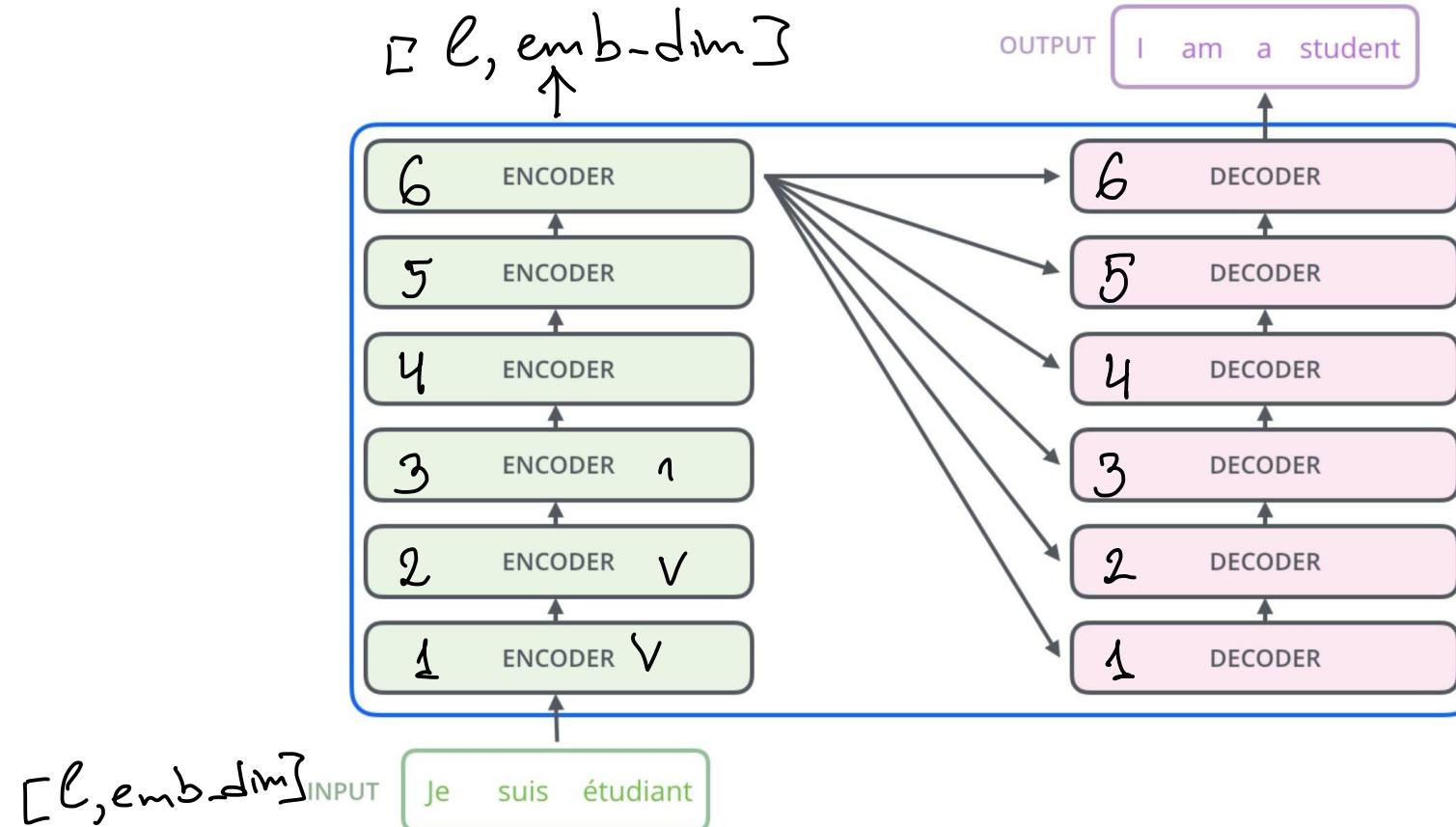
The Transformer



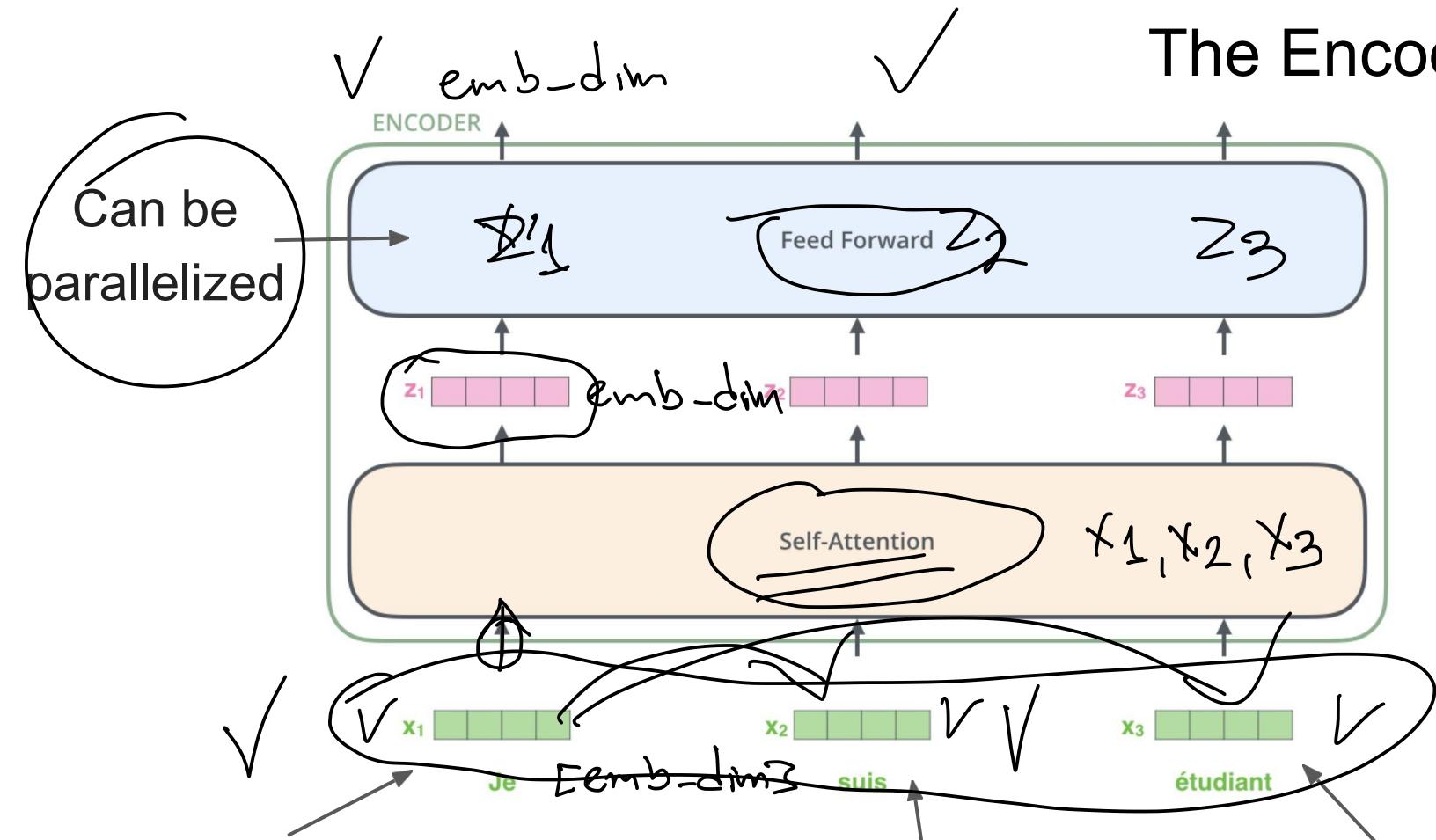
The Transformer



The Transformer



The Encoder Side



the word in each position flows through its own path in the encoder 7

The Transformer: quick overview

- Proposed in 2017 in paper [Attention is All You Need 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

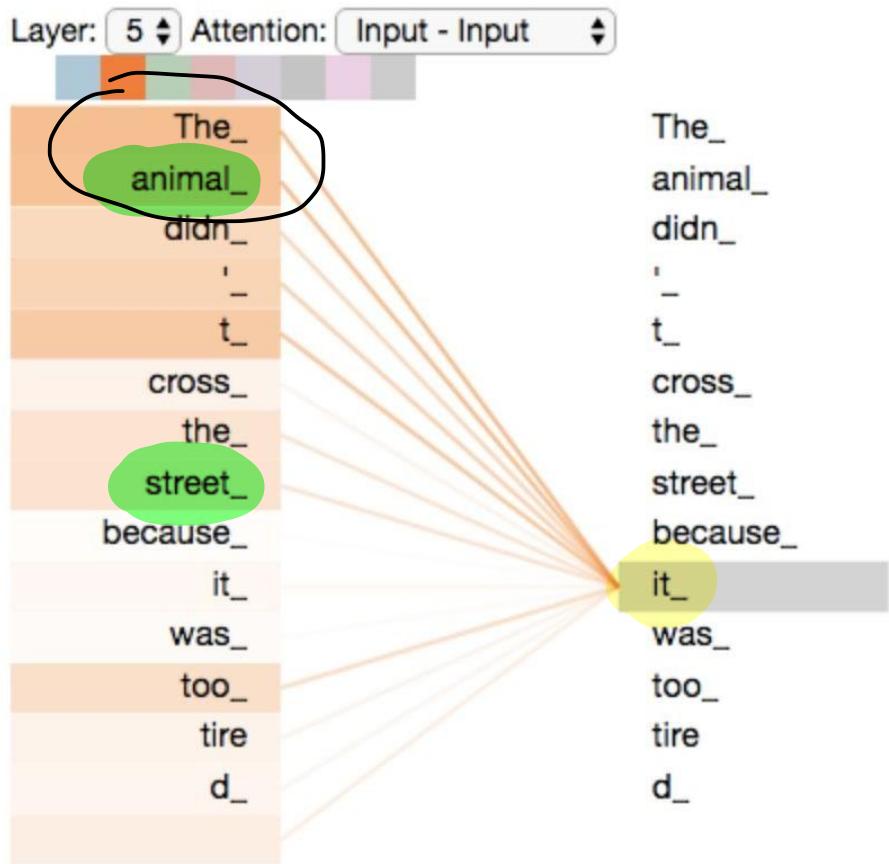


Self-Attention at a High Level

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

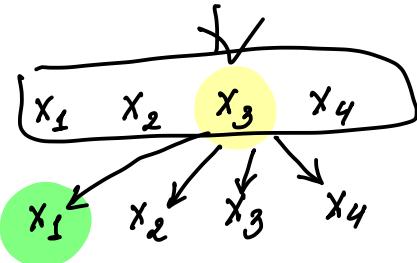
- What does “it” in this sentence refer to?
- We want self-attention to associate “it” with “animal”
- Self-attention is the method the Transformer uses to bake the “understanding” of other relevant words into the one we’re currently processing

Self-Attention at a High Level



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

BPE tokenization



$$x_3 \rightarrow z_3$$

$$z_3 = \sum_{i=1}^4 a_i v_i, \quad \sum_{i=1}^4 a_i = 1$$

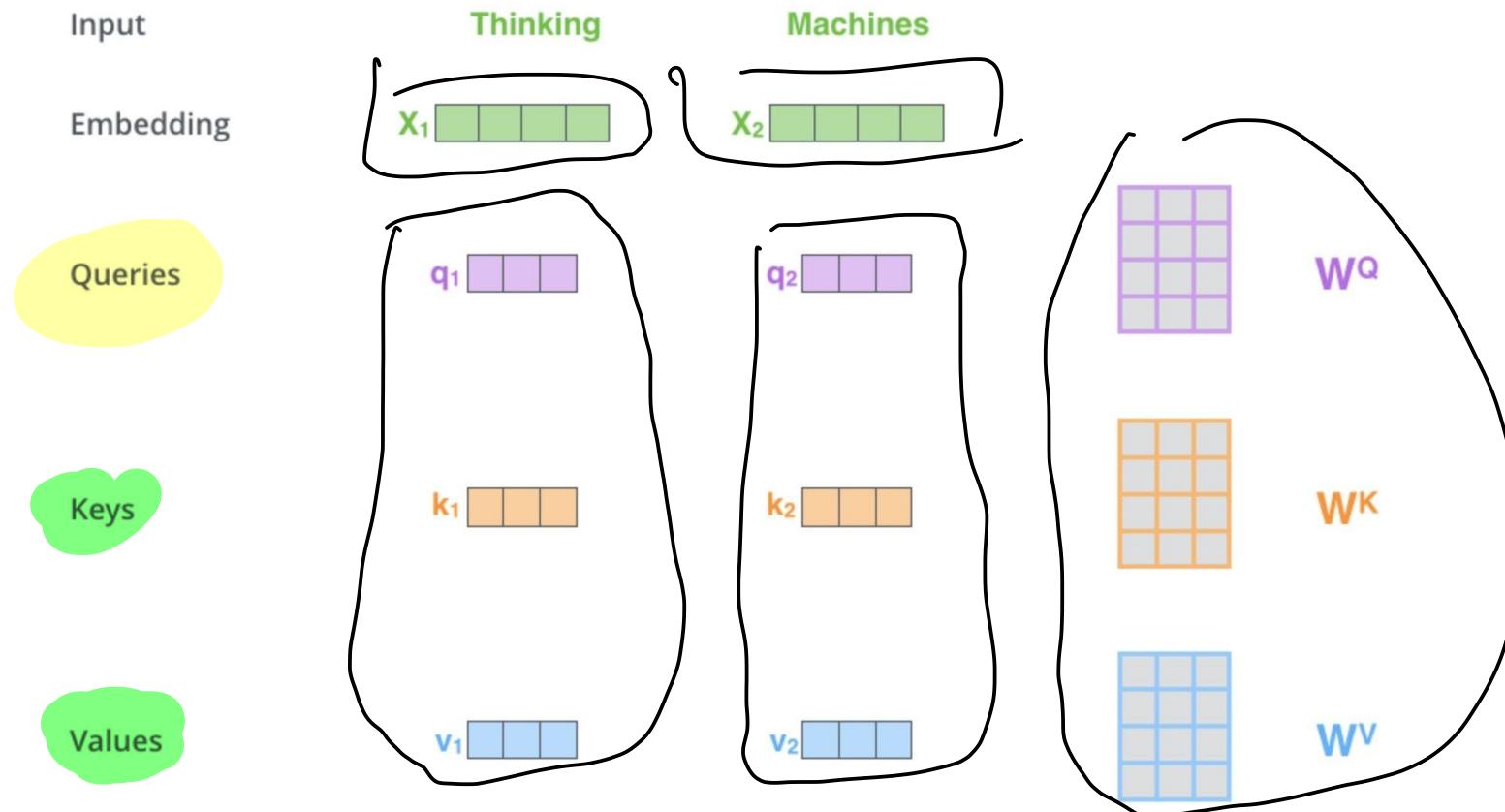
x :
 - value
 - query
 - key
 record

$$\begin{aligned} & \left\{ \begin{array}{c} a_i \\ \vdots \\ i=1 \end{array} \right\}^4 \xrightarrow{\text{softmax}} \left\{ \begin{array}{c} a_i \\ \vdots \\ i=1 \end{array} \right\}^4 \\ & \hat{a}_i - ? \quad \hat{a}_i = f(x_3, \vec{x}_i) \\ & \hat{a}_i = \langle q_3, k_i \rangle \end{aligned}$$

(\vec{x}_i) - emb

$q_i, k_i, v_i \times W.$

Self-Attention: detailed explanation

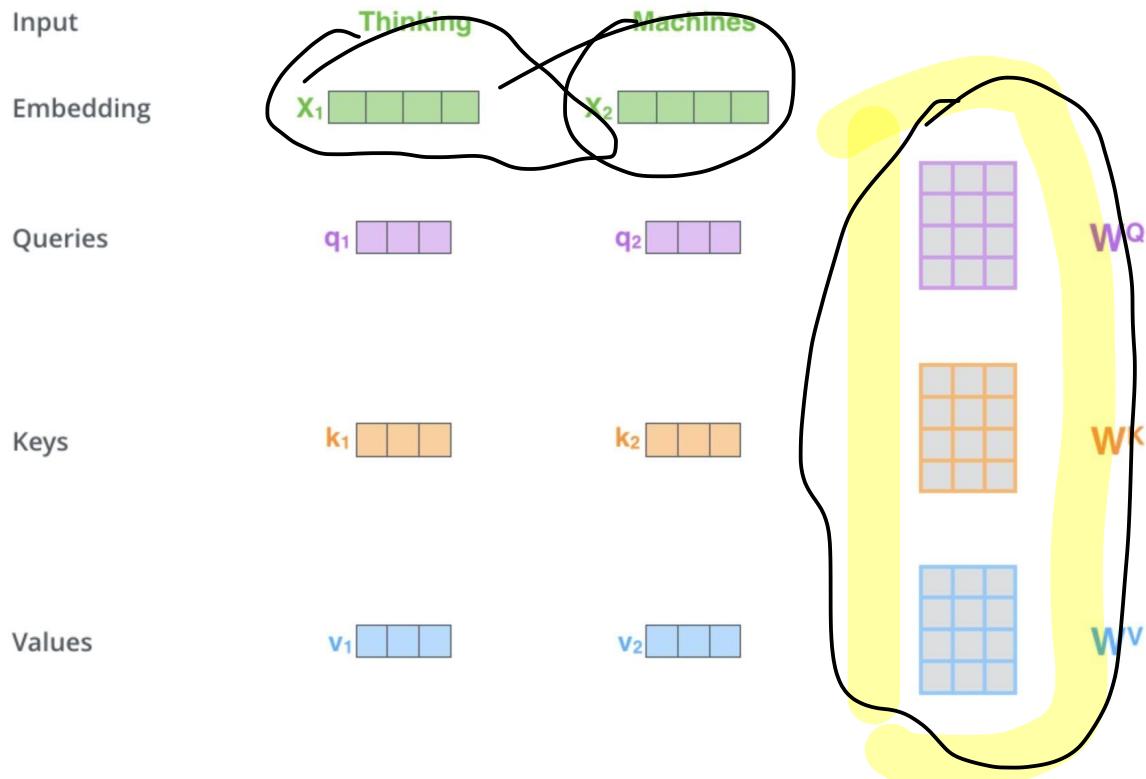


Self-Attention: detailed explanation

STEP 1:

create 3 vectors
(query, key, value)

from each of the encoder's
input vectors



$$\begin{aligned}x - (1, n) \\ \underline{W_q, W_k, W_v - (n, m)}\end{aligned}$$

$$\left\{ \begin{array}{l} q = x \cdot W_q \\ k = x \cdot W_k \\ v = x \cdot W_v \end{array} \right. \Rightarrow q, k, v - (1, m)$$

Self-Attention: detailed explanation

What are the **query**, **key**, **value** vectors?

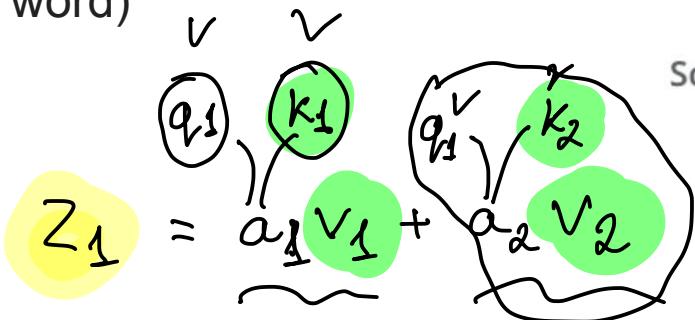
They're abstractions that are useful for calculating and thinking about attention.

Self-Attention: detailed explanation

STEP 2:

calculate a score

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



Input

Embedding

Queries

Keys

Values

Score

Thinking

x_1

q_1

k_1

v_1

Machines

x_2

q_2

k_2

v_2

$$q_1 \cdot k_1 = 112$$

$$q_1 \cdot k_2 = 96$$

$$x_1 \rightarrow z_1$$

Self-Attention: detailed explanation

STEP 3:

divide the scores by 8

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

STEP 4:

$$\sqrt{64} = 8$$

softmax

Input

Embedding

Queries

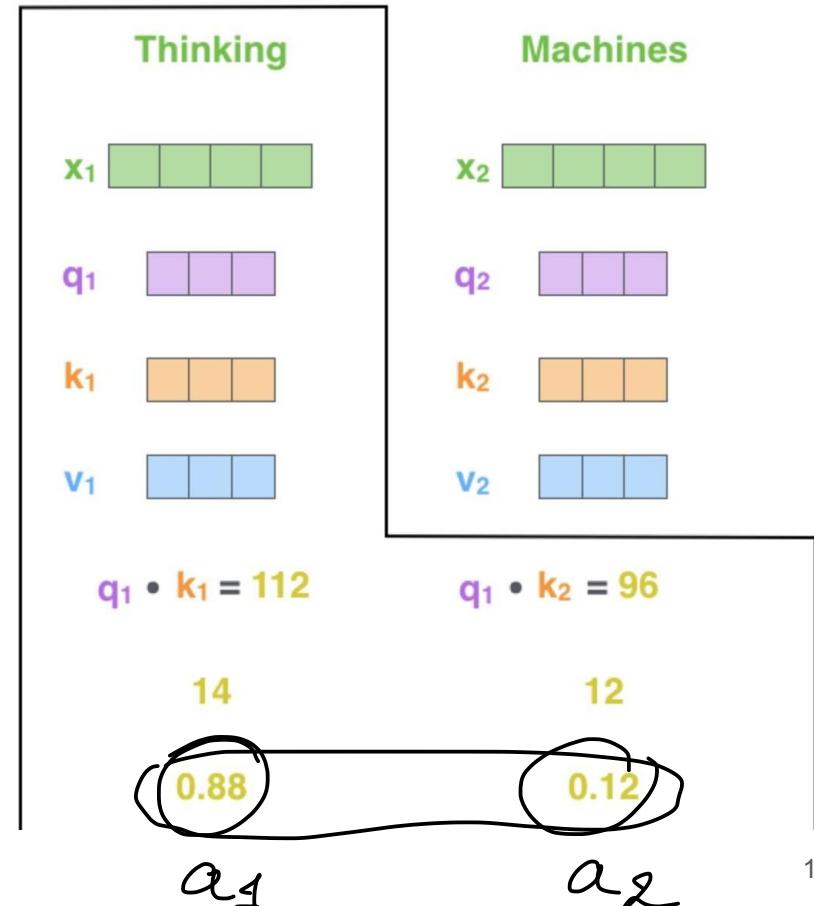
Keys

Values

Score

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

Softmax



$$\langle q, k \rangle = \sum_{i=1}^m q_i k_i = \sum_{i=1}^m (\underbrace{q_i}_{\xi_i}) \cdot \underbrace{\left(\xi_i \sim N(0, 1) \right)}_{k_i}$$

$$f = \langle q, k \rangle = \sum_{i=1}^m \xi_i \quad | \quad E_f = E \sum \xi_i = \sum_{i=1}^m E \xi_i = m E \xi_i = 0$$

$$D_f = D \sum_{i=1}^m \xi_i = \sum_{i=1}^m D \xi_i = \boxed{m}$$

$$\frac{D_f}{a} = \frac{D\xi}{a^2}$$

$$\frac{D\xi}{\sqrt{m}} = \frac{D\xi}{m} = \underline{\xi} \Rightarrow \underbrace{\frac{\xi}{\sqrt{m}} \sim N(0, 1)}_{\xi_i}$$

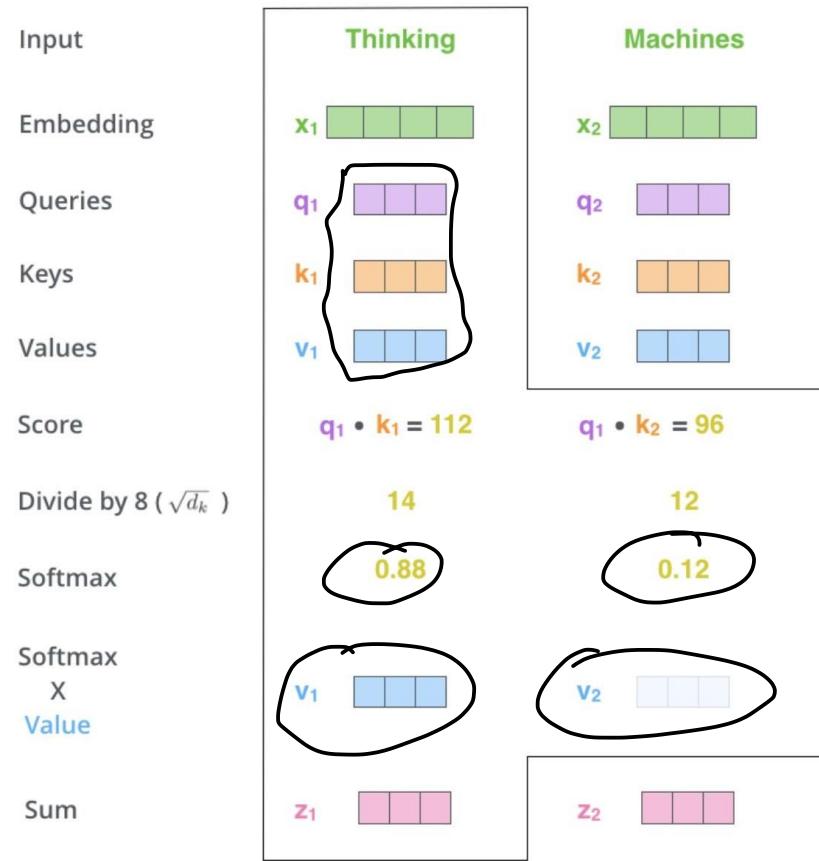
Self-Attention: detailed explanation

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 ($\sqrt{d_k}$)

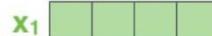
Softmax

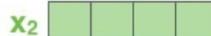
Softmax
X
Value

Sum

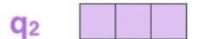
Thinking

Machines

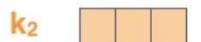
x_1 

x_2 

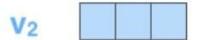
q_1 

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 $\sqrt{d_k}$

STEP 3: divide by

STEP 4: softmax

STEP 5: multiply each value vector by the softmax score

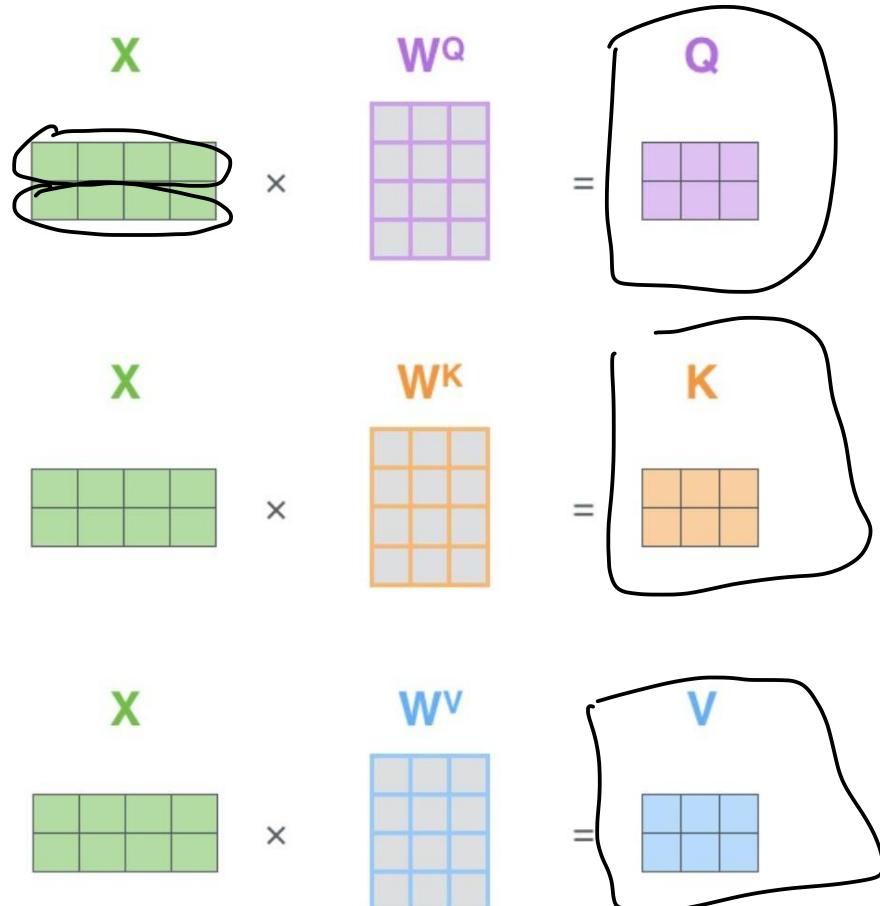
STEP 6: sum up the weighted value vectors

$\otimes W_q$

Pack embeddings into
matrix X

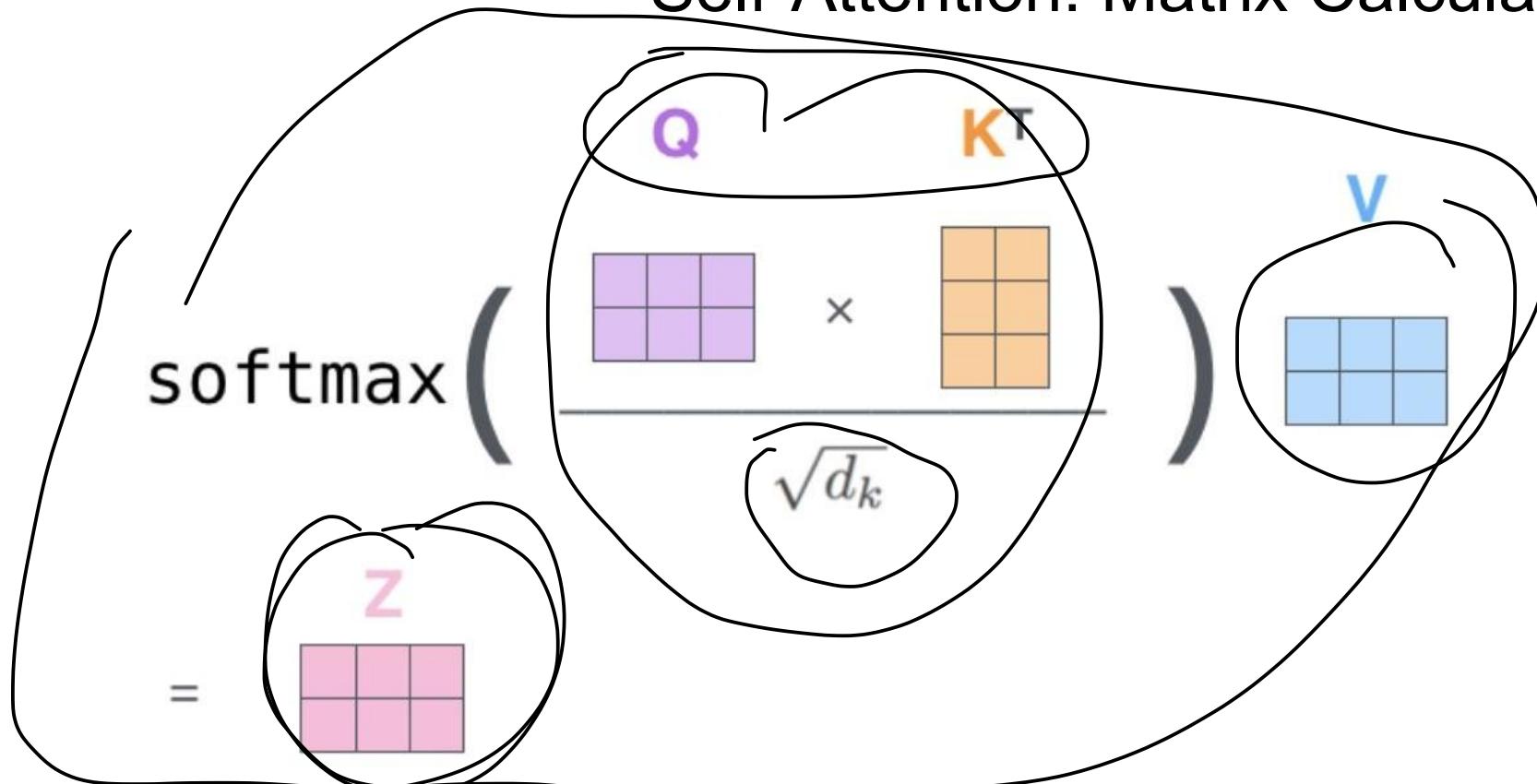
Multiply X by weight
matrices we've trained
(W_k , W_q , W_v)

Self-Attention: Matrix Calculation



$$x \cdot W_q = q$$
$$\left(\begin{matrix} x_1 \\ x_2 \\ x_3 \\ \vdots \\ x_n \end{matrix} \right) \left(\begin{matrix} & & & \\ & \text{...} & & \\ & & \text{...} & \\ & & & m \end{matrix} \right) = \left(\begin{matrix} q_1 \\ q_2 \\ q_3 \\ \vdots \\ q_n \end{matrix} \right)$$
$$(L, n)$$

Self-Attention: Matrix Calculation



$$X - (e, n) \Rightarrow Q, K, V - (e, m)$$

$$q_1 \quad q_2 \quad q_3 \quad \cdot \quad \begin{pmatrix} \parallel & \parallel & \dots & \parallel \\ k_1 & k_2 & \dots & k_L \end{pmatrix} = i \quad \hat{a}$$

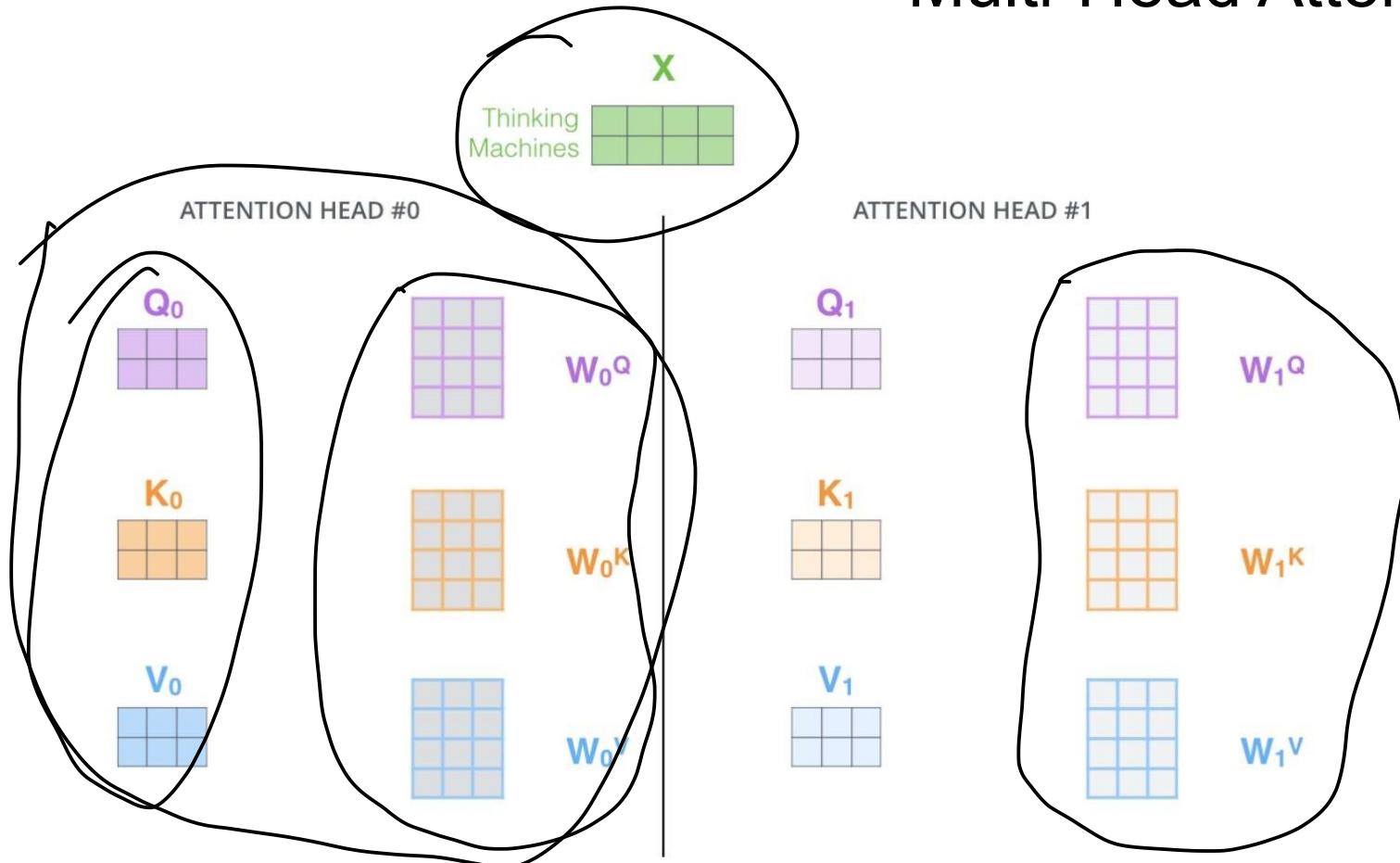
$$\hat{A} - (e, l) \quad z_i = \sum a_i v_i$$

$$A = \text{softmax}(\hat{A}) = \frac{\exp(\hat{A}^T)}{\sqrt{m}}$$

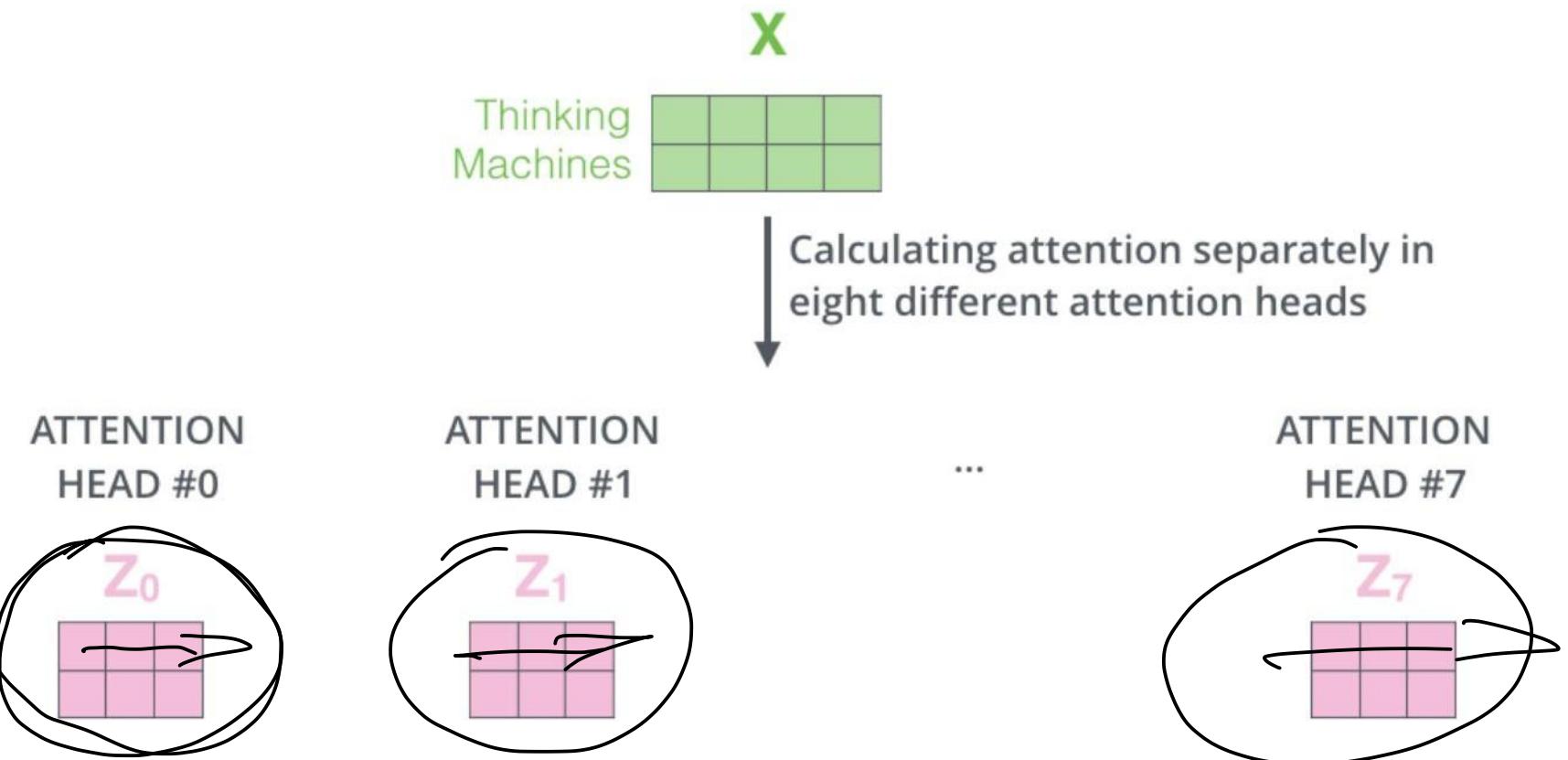
$$A \cdot V = Z \quad ; \quad \begin{pmatrix} a_1 & \dots & a_l \\ \vdots & \ddots & \vdots \\ a_m \end{pmatrix} \cdot \begin{pmatrix} v_1 \\ v_2 \\ v_3 \end{pmatrix} = i \quad \begin{pmatrix} z_1 \\ z_2 \\ z_3 \end{pmatrix}$$

$$z_1 = \sum a_i v_i.$$

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

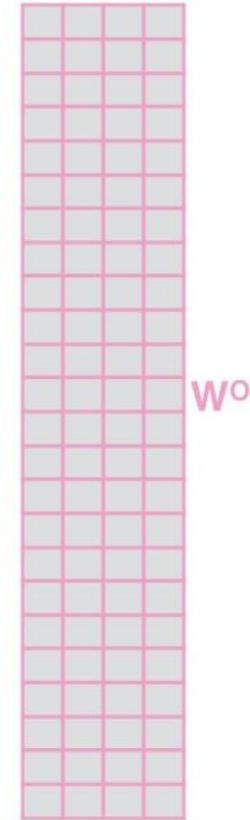
x

$$(\ell, \delta m) \quad \delta m \rightarrow m \Rightarrow 2 - (\ell, m)$$

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

$$= \begin{matrix} Z \\ \text{---} \\ \text{---} \end{matrix} \rightarrow (\delta m, m)$$

$\sum_i \rightarrow \sum_i^1$



1) This is our input sentence*

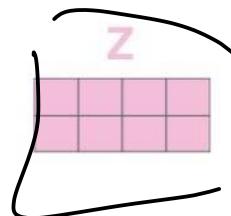
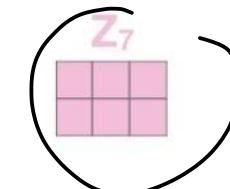
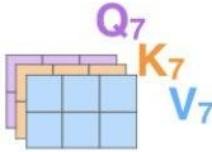
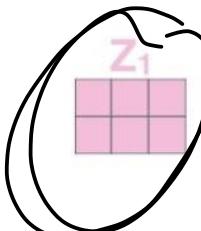
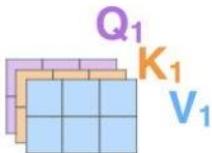
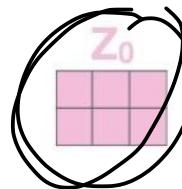
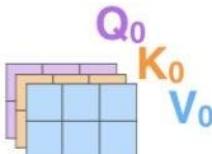
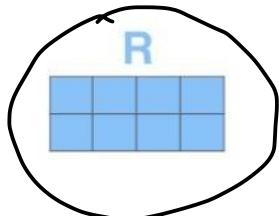
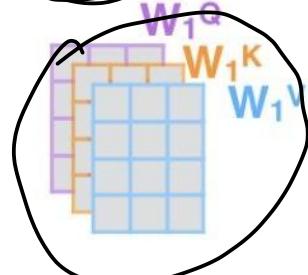
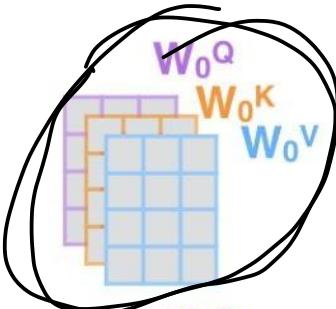
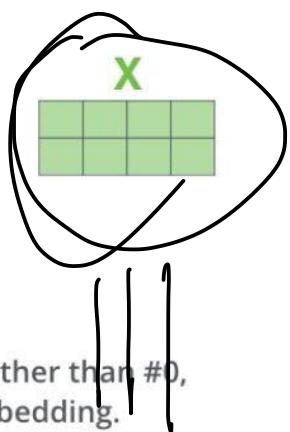
2) We embed each word*

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

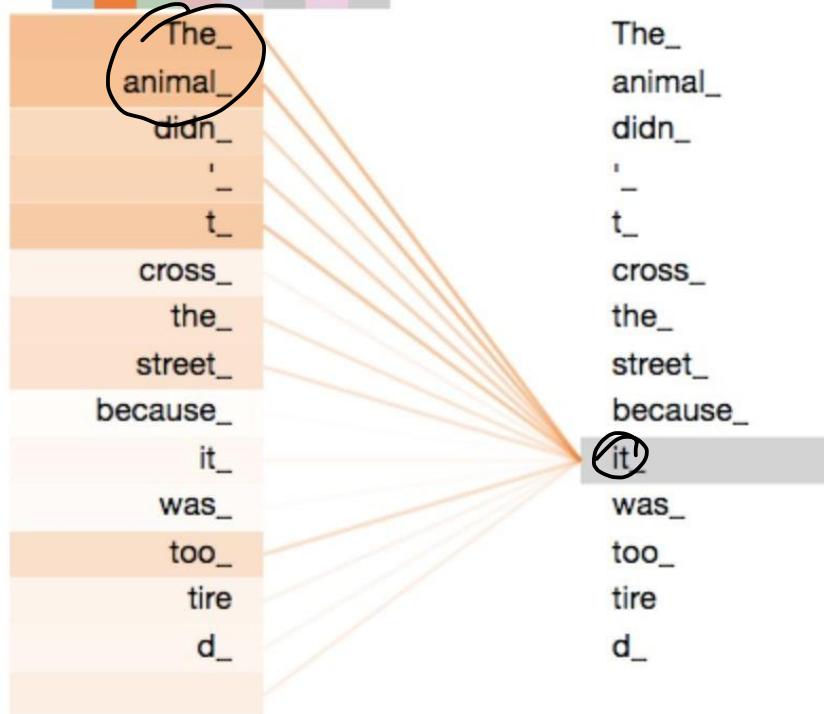
Thinking
Machines



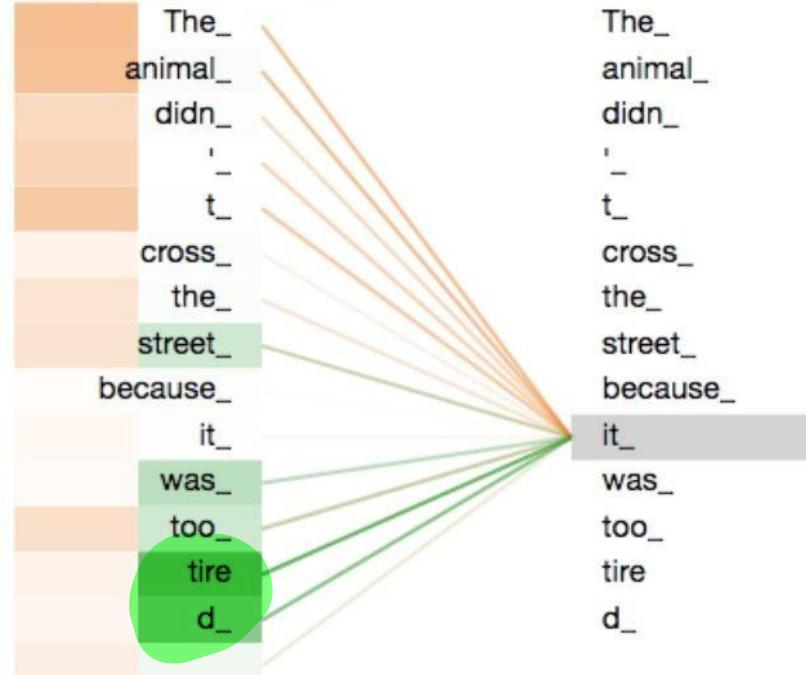
* In all encoders other than #0, we don't need embedding. We start directly with the output of the encoder right below this one

Multi-Head Attention

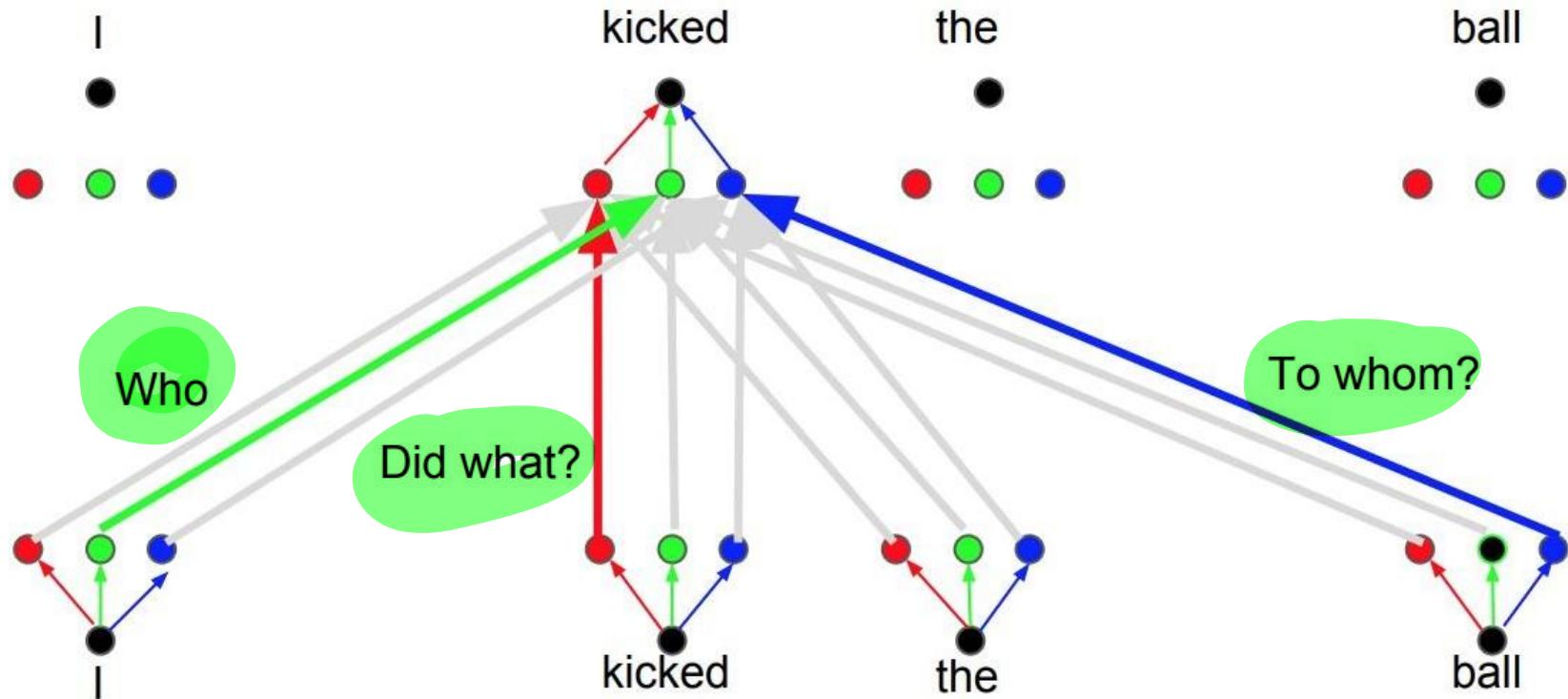
Layer: 5 ⬆️ Attention: Input - Input ⬆️



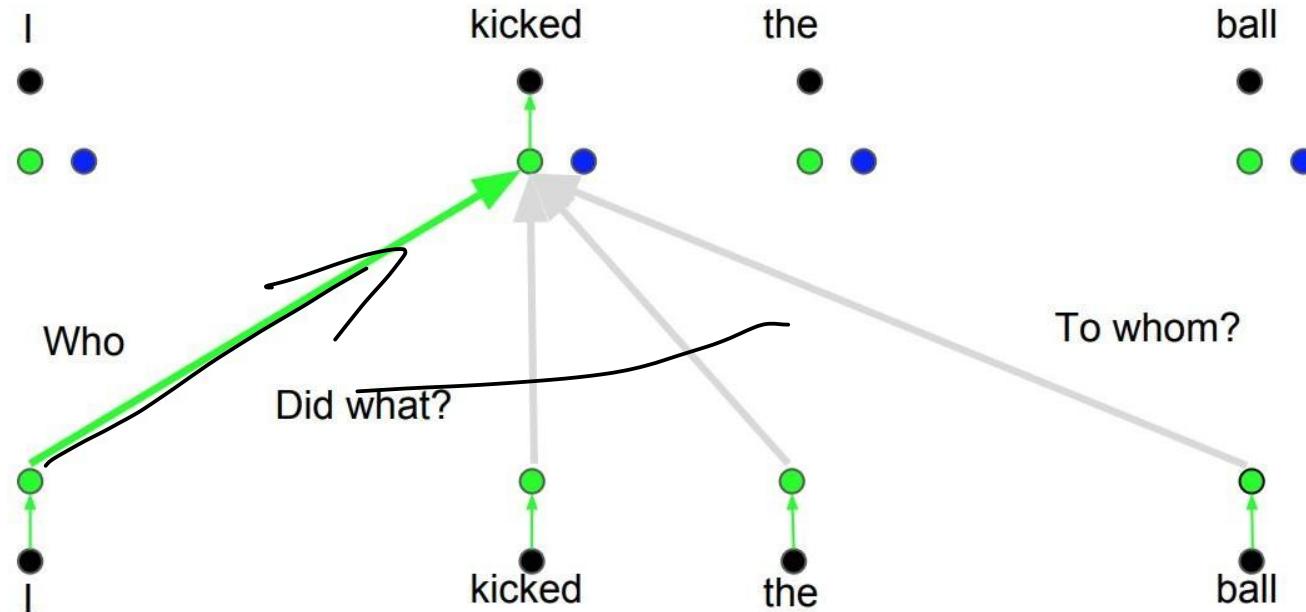
Layer: 5 ⬆️ Attention: Input - Input ⬆️



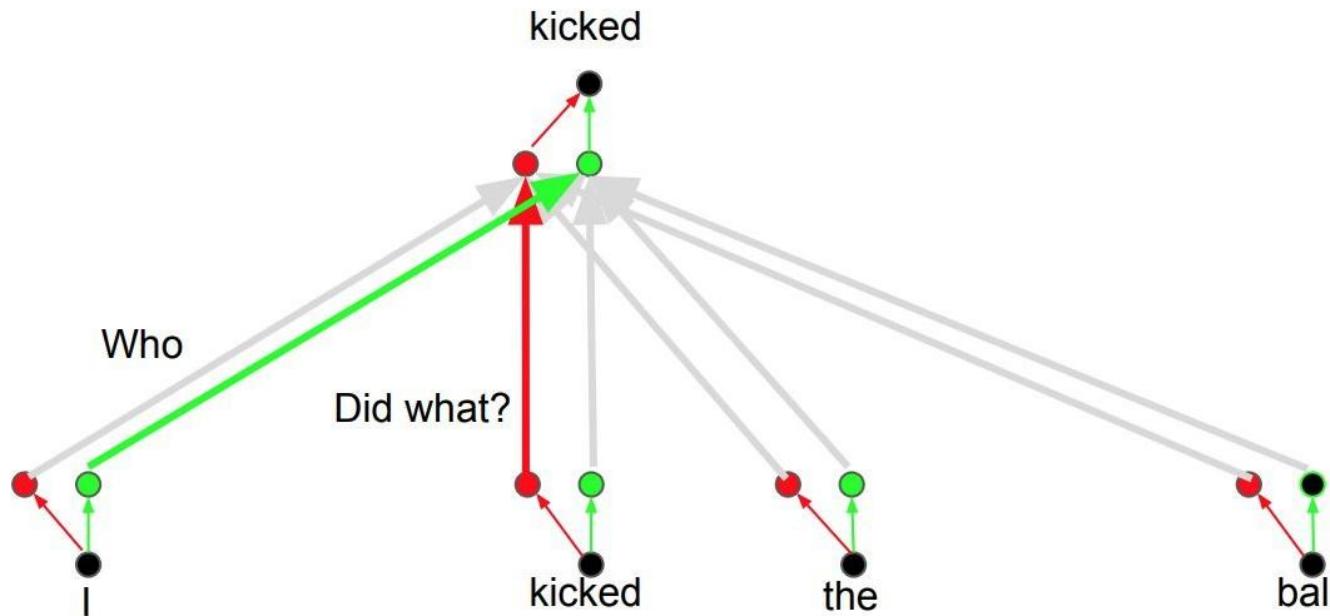
Why Multi-Head Attention?



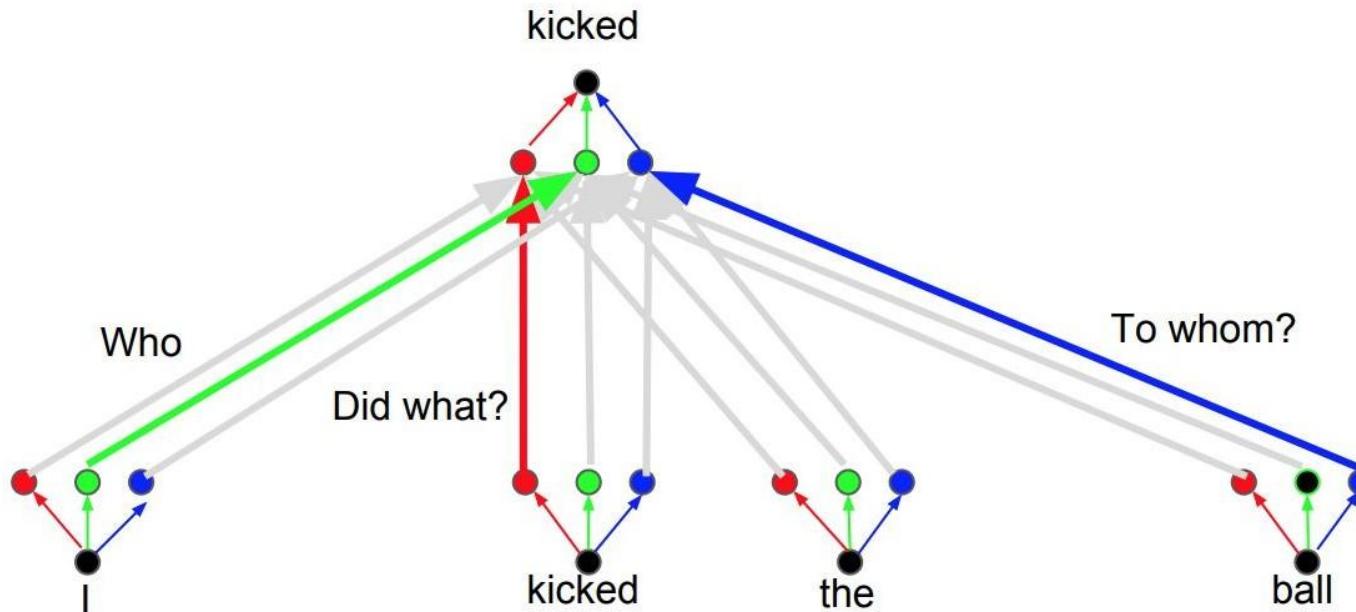
Attention head: Who



Attention head: Did What?

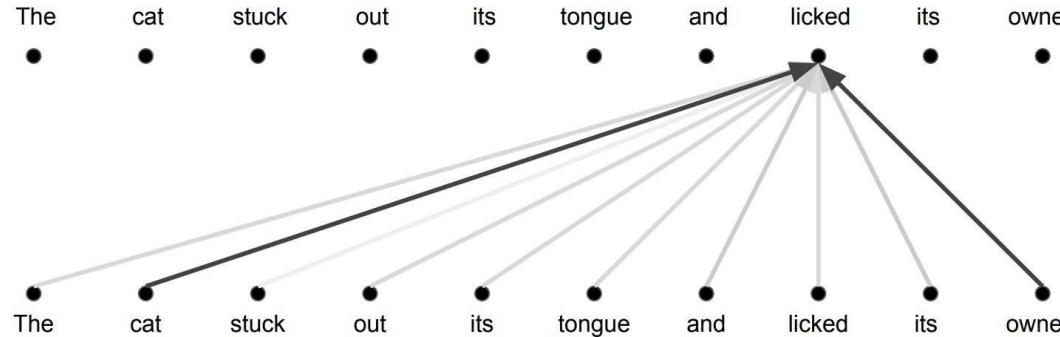


Attention head: To Whom?

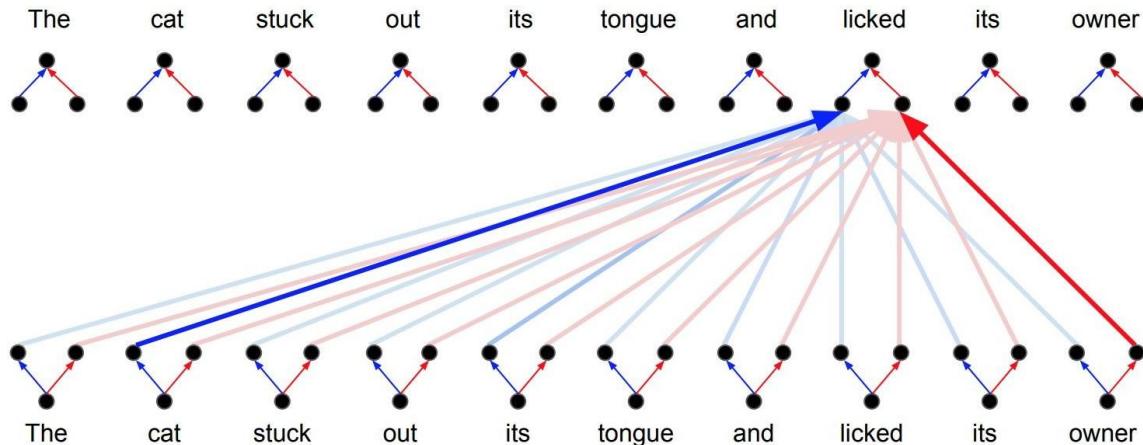


Attention vs. Multi-Head Attention

Attention: a weighted average



Multi-Head Attention: parallel attention layers with different linear transformations on input and output.



Performance: WMT 2014 BLEU

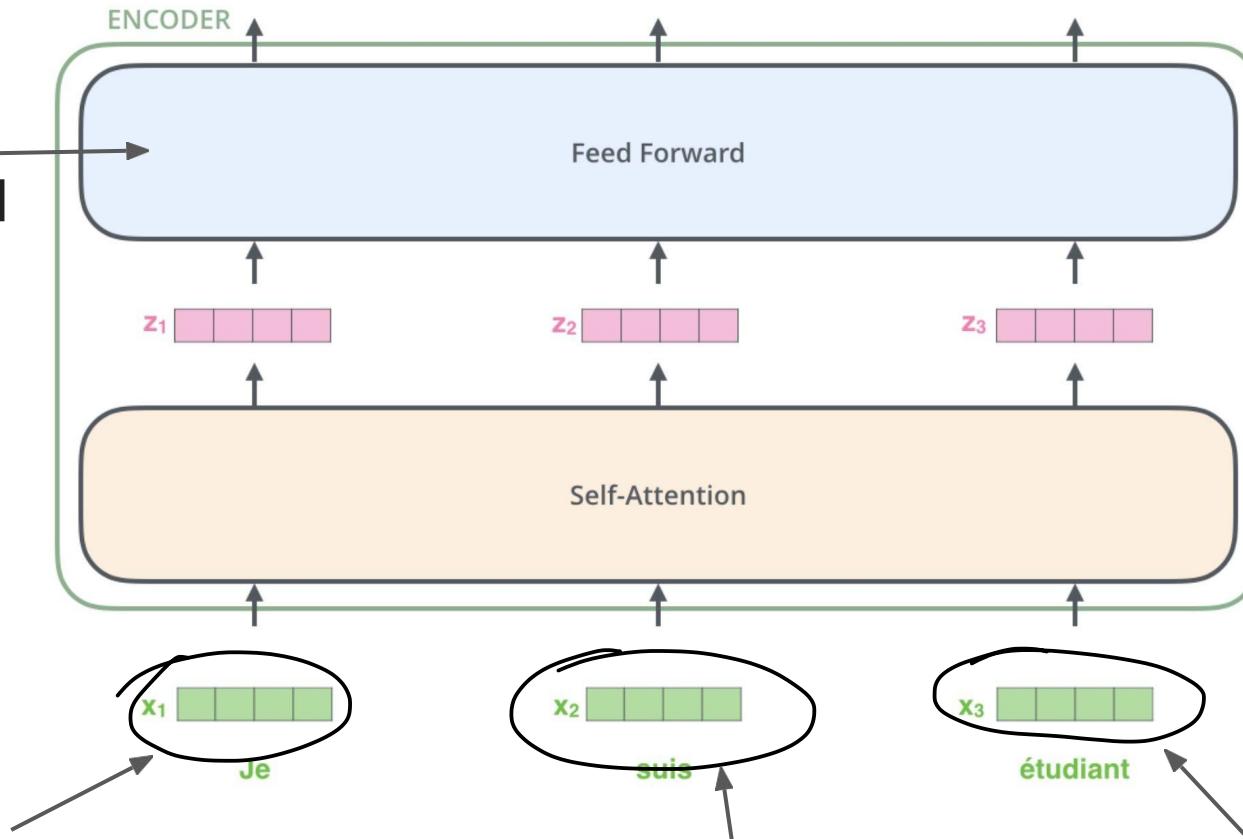
	EN-DE	EN-FR
GNMT (orig)	24.6	39.9
ConvSeq2Seq	25.2	40.5
Transformer*	28.4	41.8

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

Positional Encoding

The Encoder Side

Can be parallelized



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

Positional encoding requirements

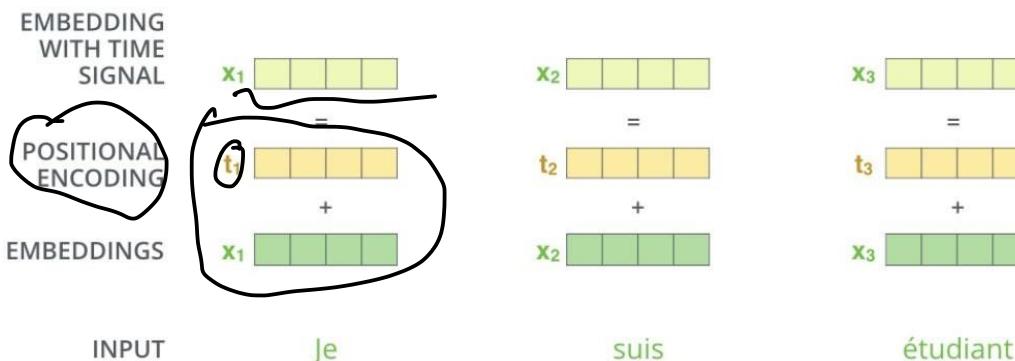
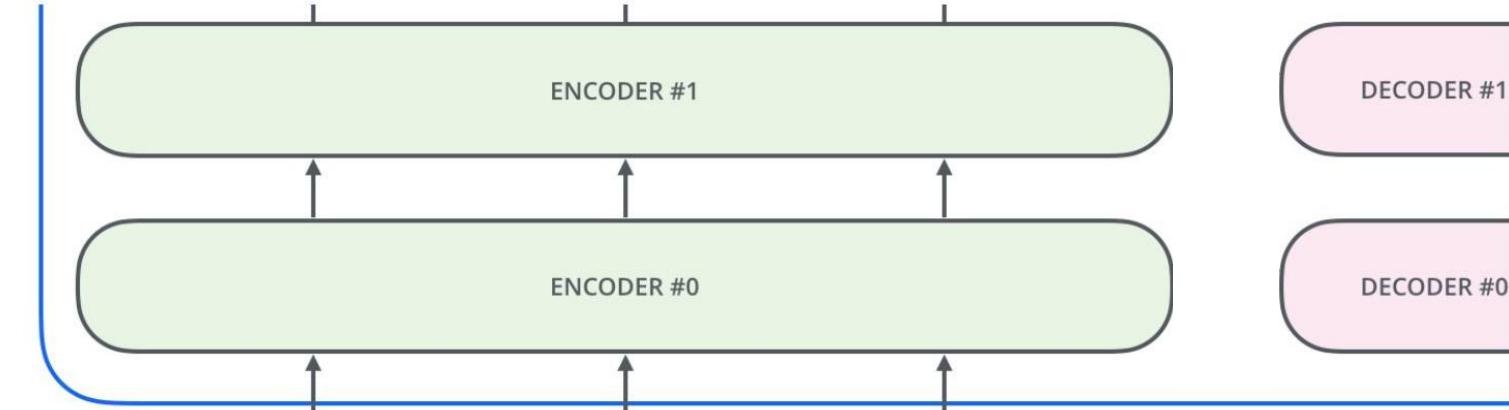
- ✓ • ✓ Positional encoding should be unique for every position in the sequence
- ? • Distance between two same positions should be preserved with sequences of different length
- ✓ • ✓ The positional encoding should be deterministic
- ✓ • *It would be great if it would work with long sequences (longer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

1, 4

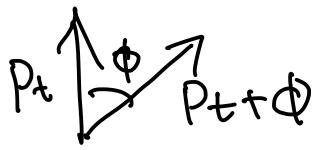
9, 7

5, 8

Positional Encoding



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



$$\vec{p}_t^{(i)} = f(t)^{(i)}$$

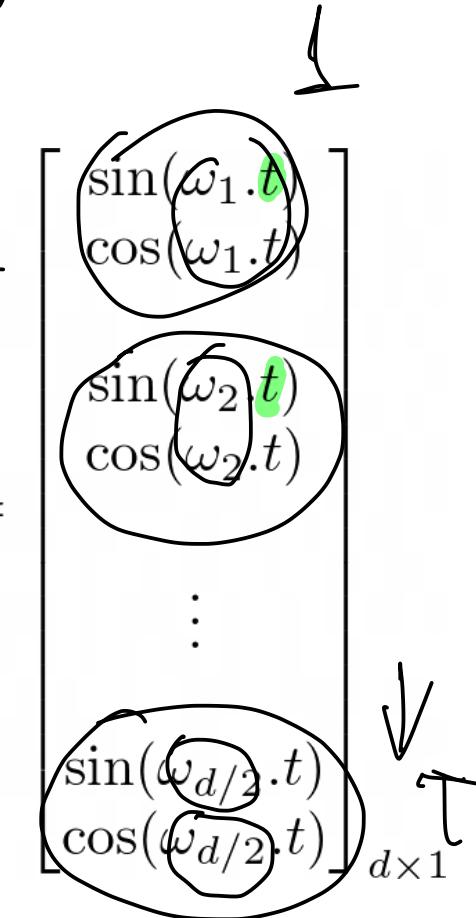
Positional Encoding: why sin and cos?

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

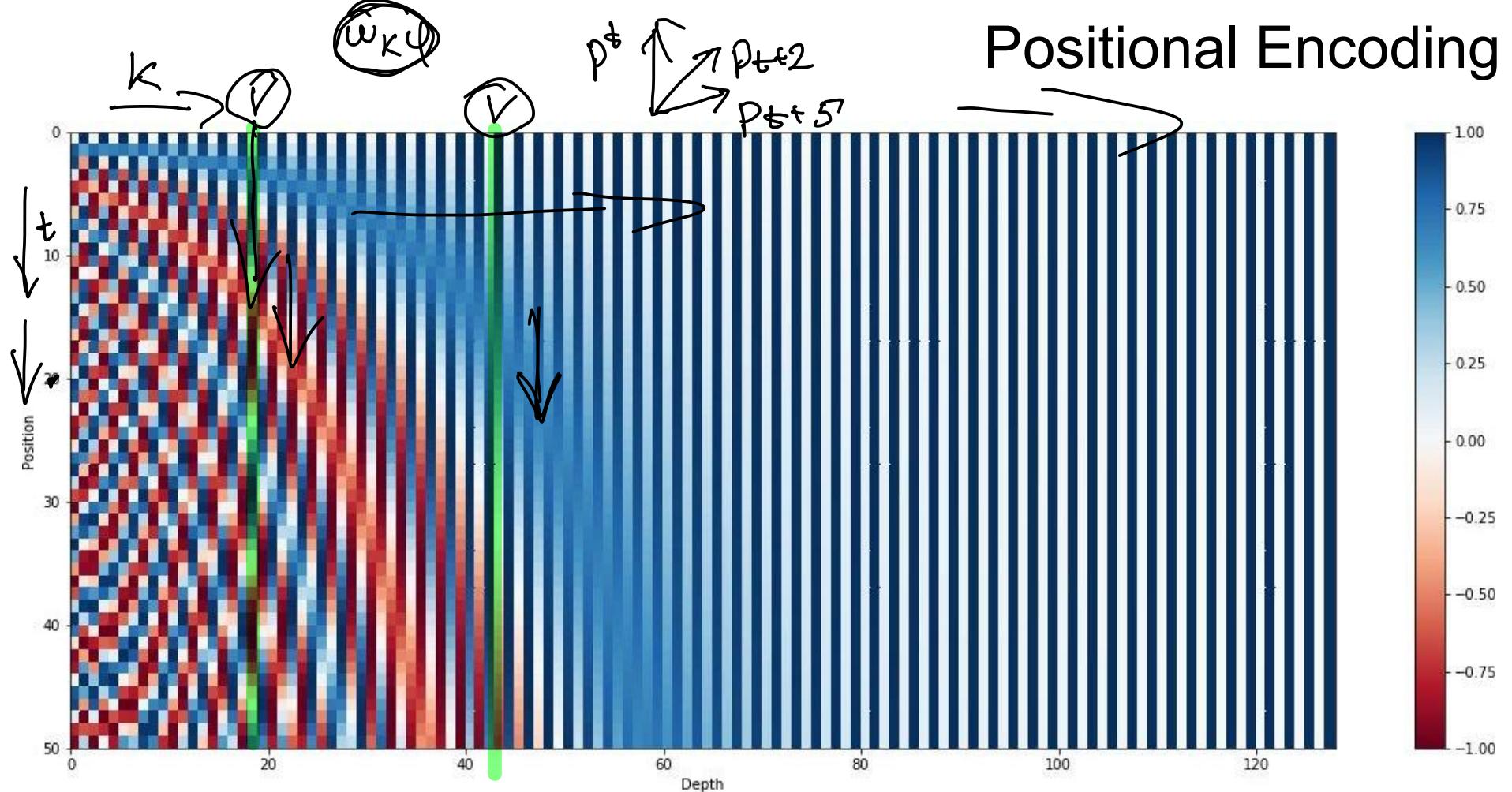
$$\omega_k = \frac{1}{10000^{2k/d}}$$

t stays for position in the original sequence
 k is the index of the element in the positional vector

$$\vec{p}_t = \sqrt{\frac{d}{2}}$$



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 , PE_{pos+k} can be represented as a linear function of PE_{pos} .

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

$P_t \rightarrow P_{t+\phi}$

$P_t \rightarrow P_{t+\phi}$

Positional Encoding: why sin and

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

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

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

P_t

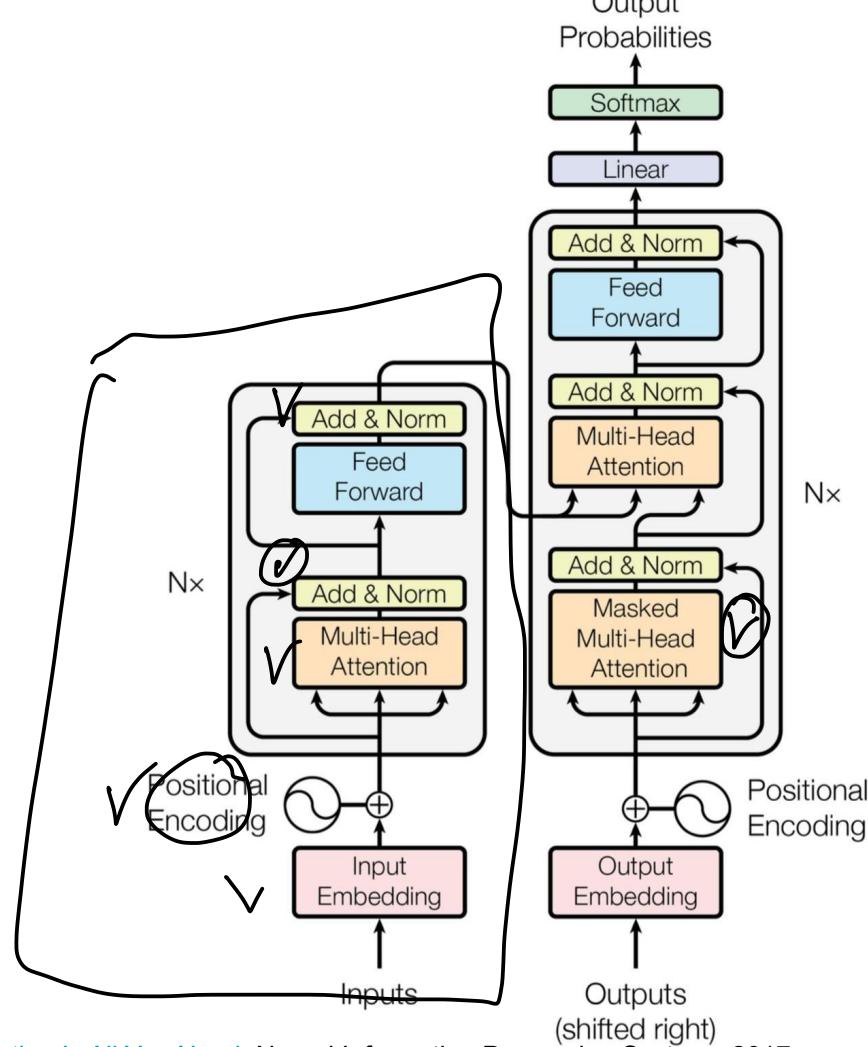
Layer Normalization

$$q = 3$$

1, 4

5, 8

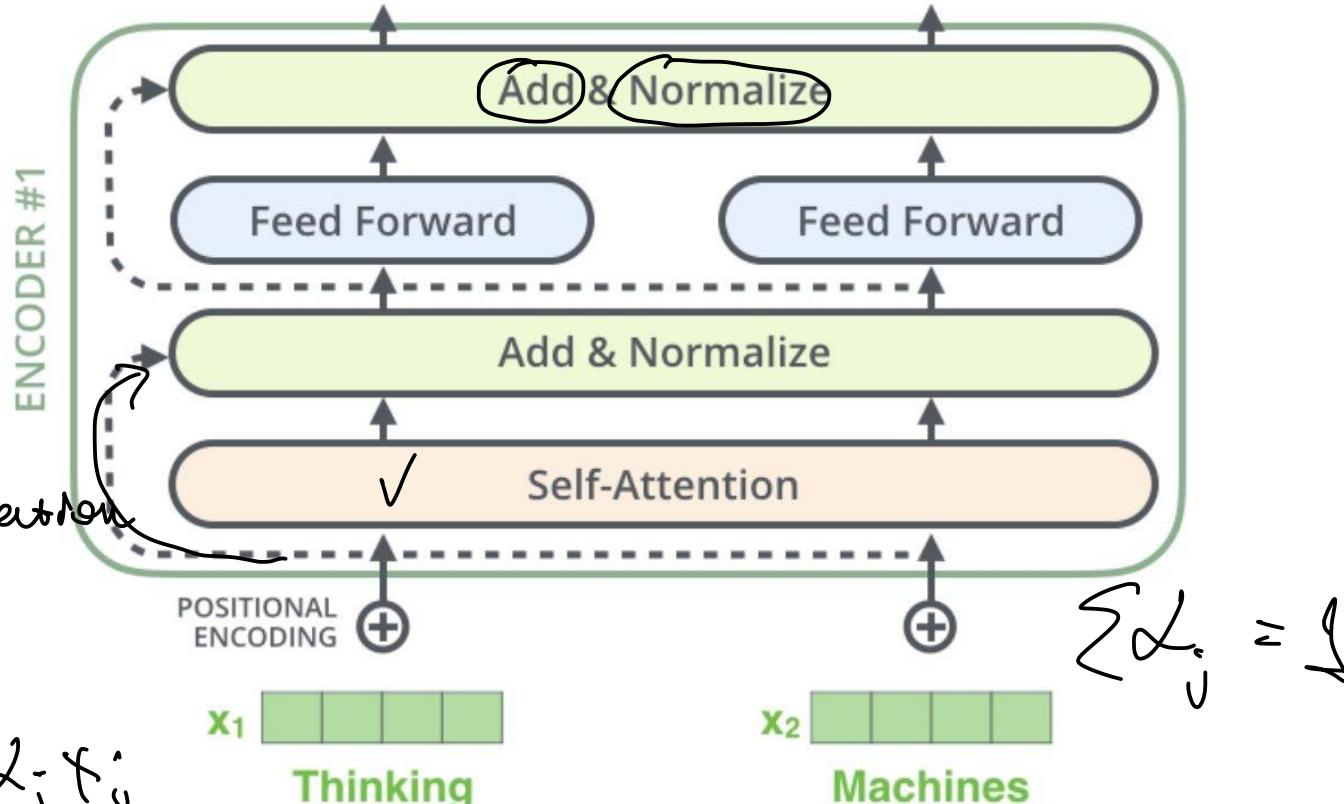
The Transformer: recap



Layer Normalization

$$z = \gamma z + \beta$$

$$z \sim N(0, 1)$$

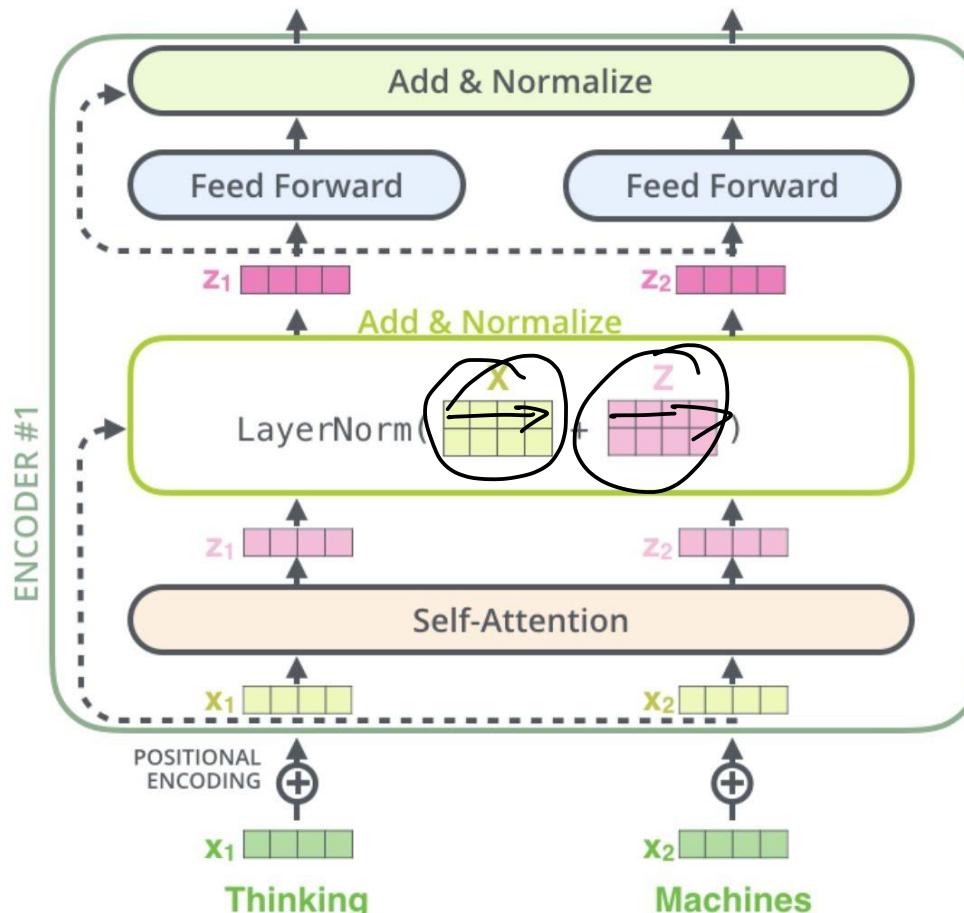


Layer Normalization

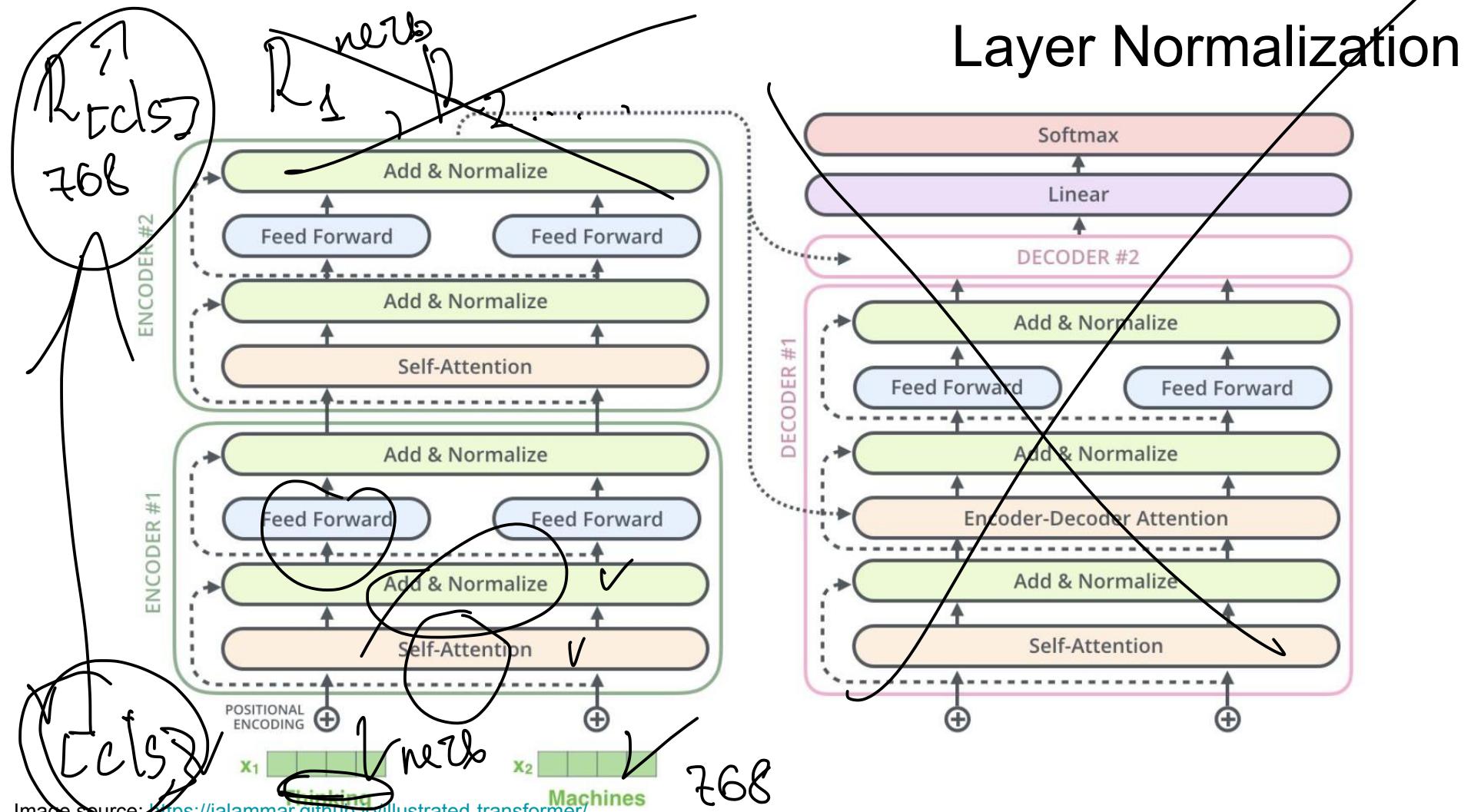
Like BatchNorm

but normalize along
all features
representing latent
vector

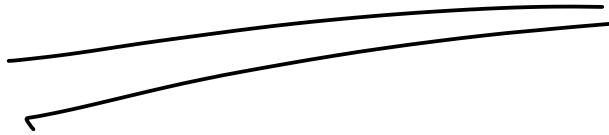
More info:
[Layer Normalization](#)



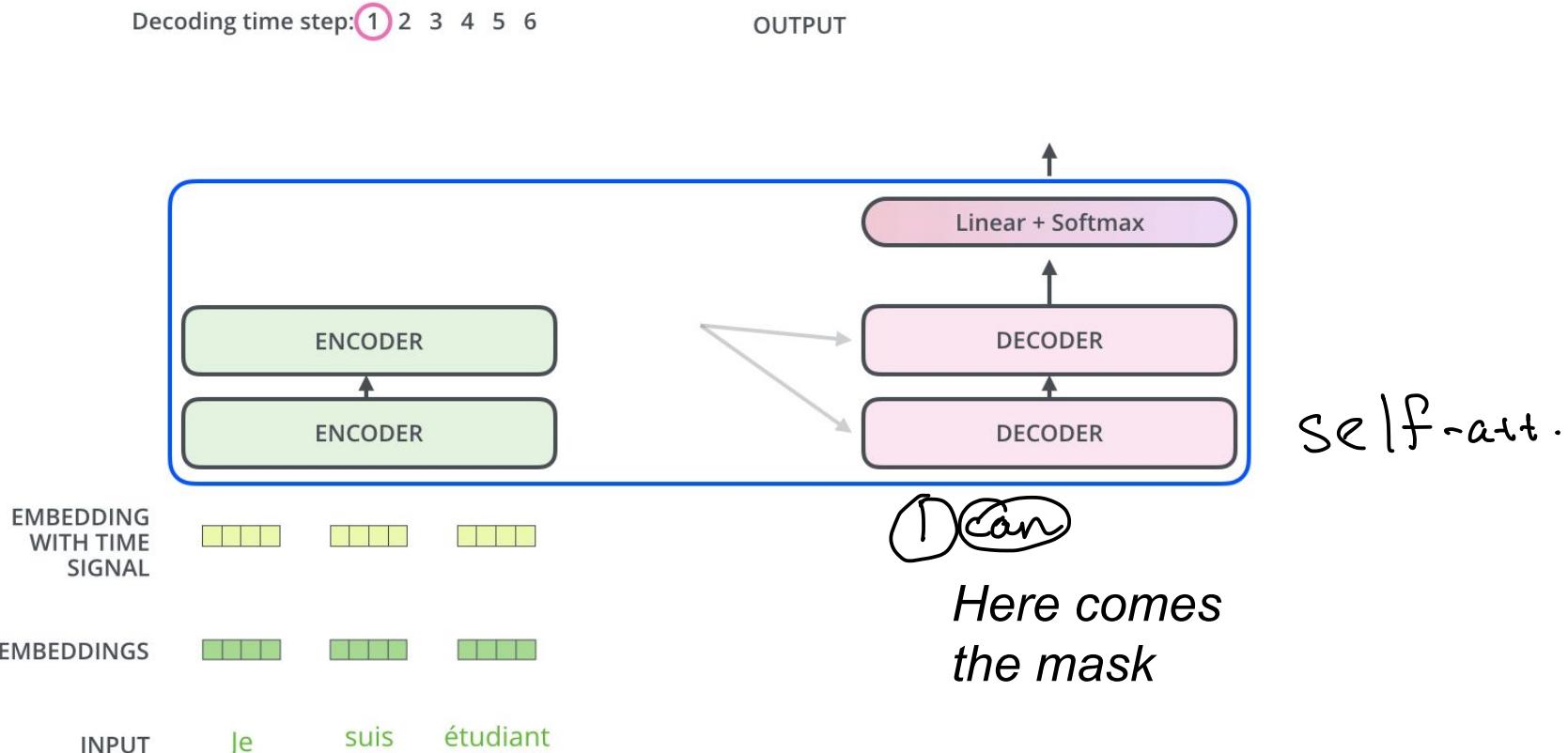
Layer Normalization



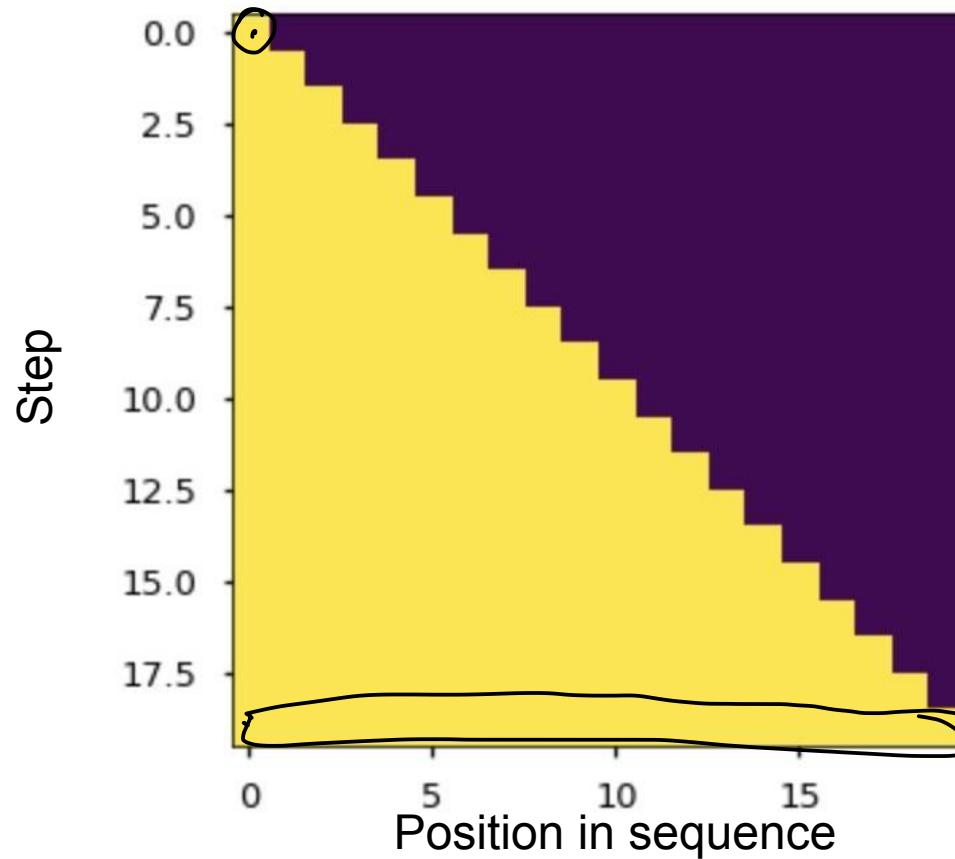
The Decoder



The Decoder Side



The masked decoder input



$(\ell, \text{emb-dim})$

Decoding time step: 1 2 3 4 5 6

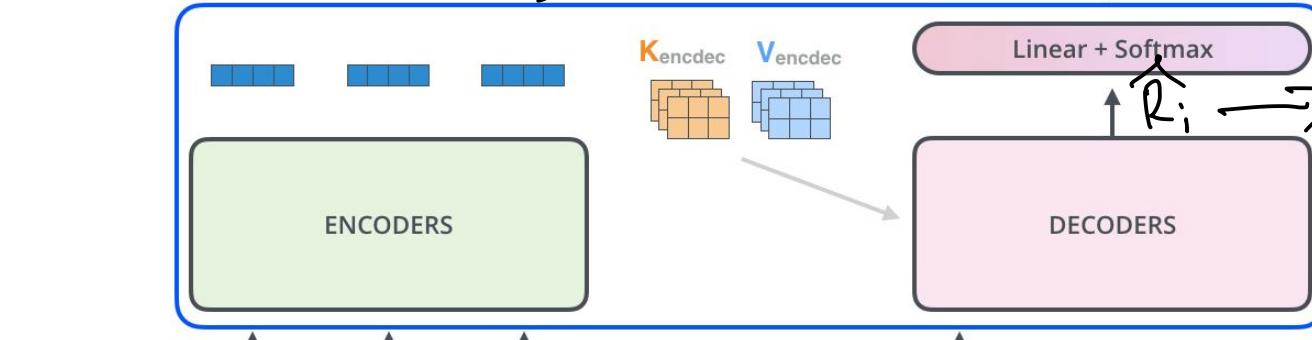
W_K, W_V, W_Q

OUTPUT

The Decoder Side

$R_1 R_2 R_3$

$$K_{i,V} = W_{K,V} \cdot [R_1 \dots R_3]$$



EMBEDDING
WITH TIME
SIGNAL

EMBEDDINGS

INPUT

Je

suis

étudiant

PREVIOUS
OUTPUTS

Final Linear and Softmax Layer

Which word in our vocabulary
is associated with this index?

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

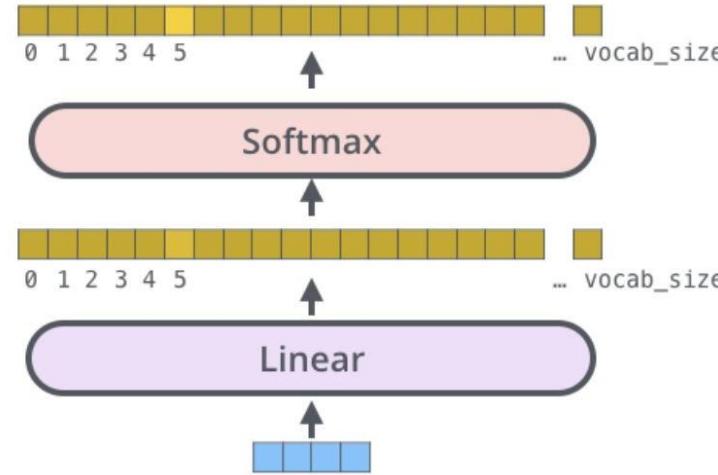
`log_probs`

`logits`

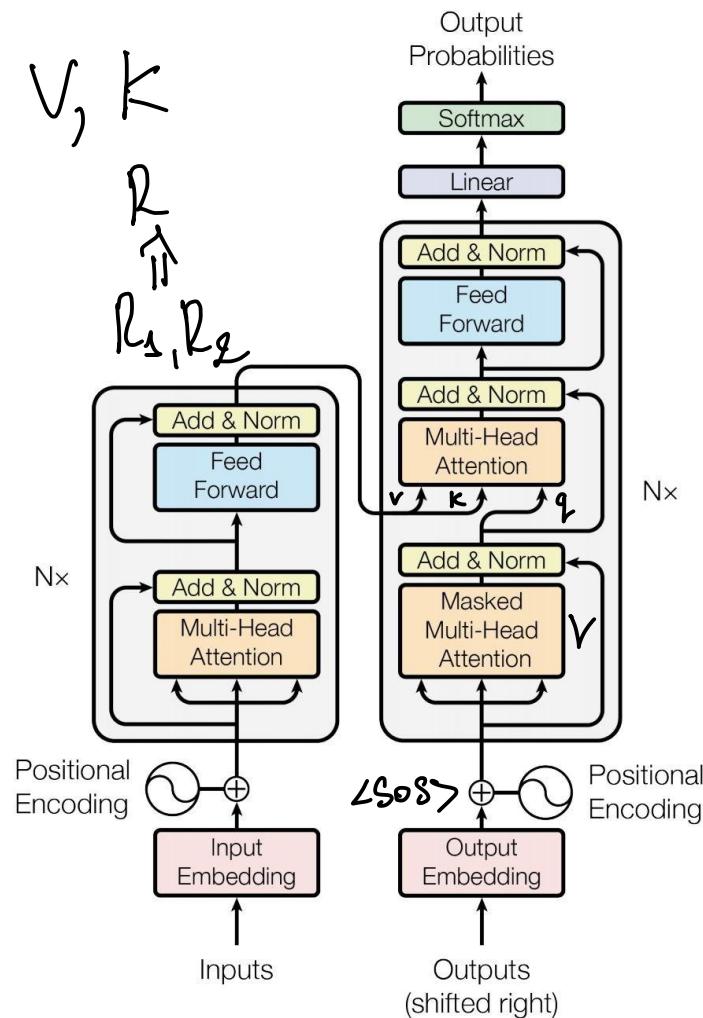
Decoder stack output

am

5



The Transformer



$$\hat{a}_1 = 1$$

$$z_1 = a_1 v_1$$

Training The Transformer

$$\beta = 0.99$$

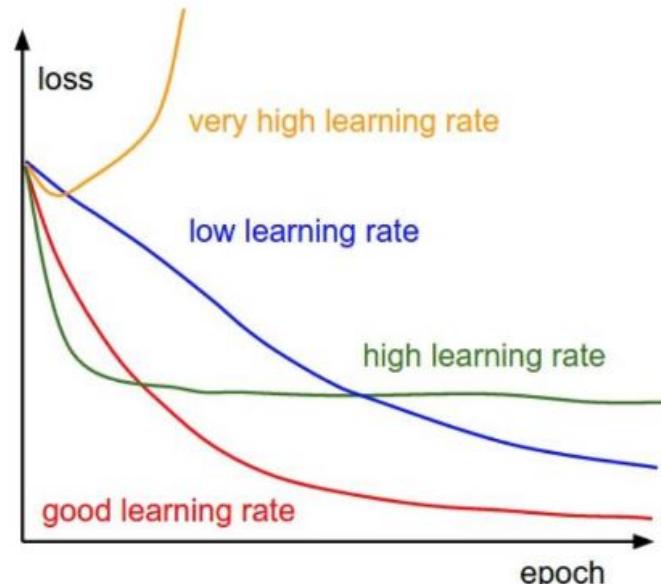
- Consider Adam optimizer

$$\mu_{t+1} = \beta_1 \cdot \mu_t + (1 - \beta_1) \cdot \nabla_{\theta} L$$

$$v_{t+1} = \beta_2 \cdot v_t + (1 - \beta_2) \cdot \|\nabla_{\theta} L\|^2$$

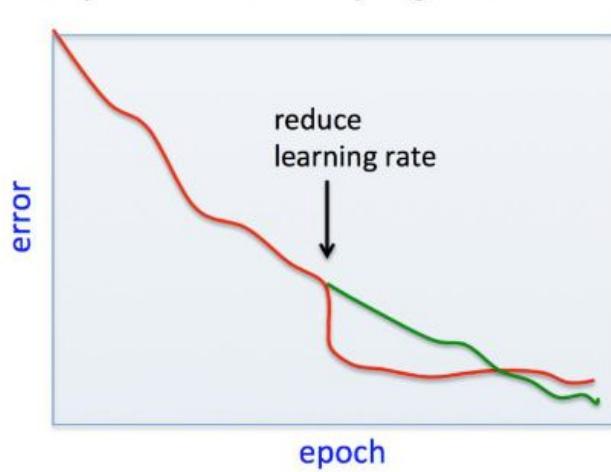
$$\theta_{t+1} = \theta_t - \frac{\alpha}{\sqrt{v_t} + \epsilon} \mu_t$$

- The choice of α is crucial!



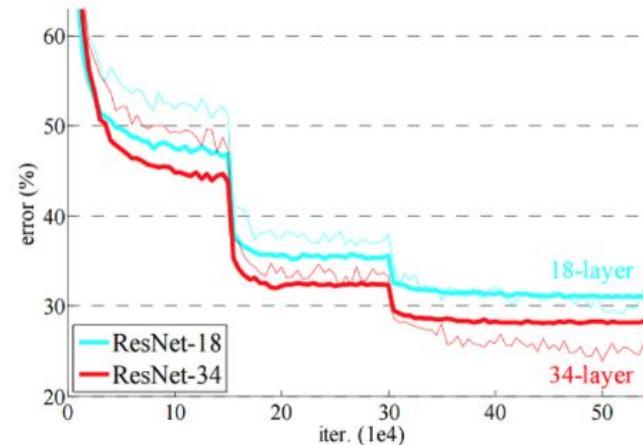
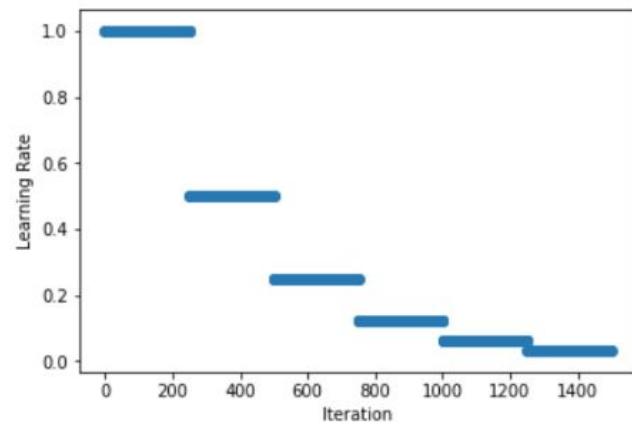
Training The Transformer

- Traditional approach: decrease learning rate in stages
 - every k steps or whenever progress slows down



Training The Transformer

- Traditional approach: decrease learning rate in stages
 - every **k** steps or whenever progress slows down



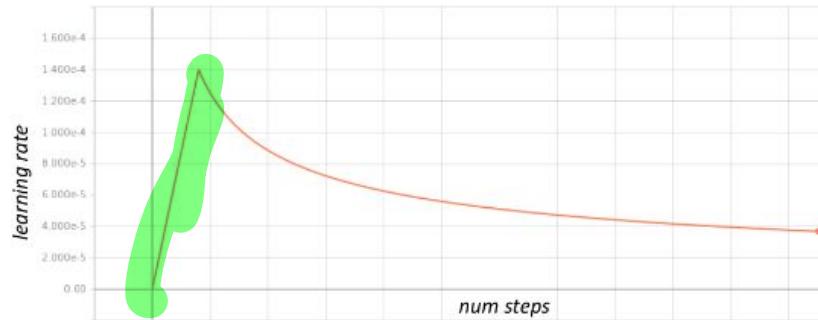
Training The Transformer

- Problem: first k steps of Adam are unstable
 - it needs time to accumulate statistics
- Use warmup time!
keep lr small over first epochs: $\alpha = \alpha_{base} \cdot \min(growth(t), decay(t))$

$$growth(t) = \frac{t}{T_{warmup}}$$

$$decay(t) = \sqrt{\frac{T_{warmup}}{t}}$$

$$t^{-1/2}$$



Training The Transformer

Transformers trainings requires non-trivial efforts

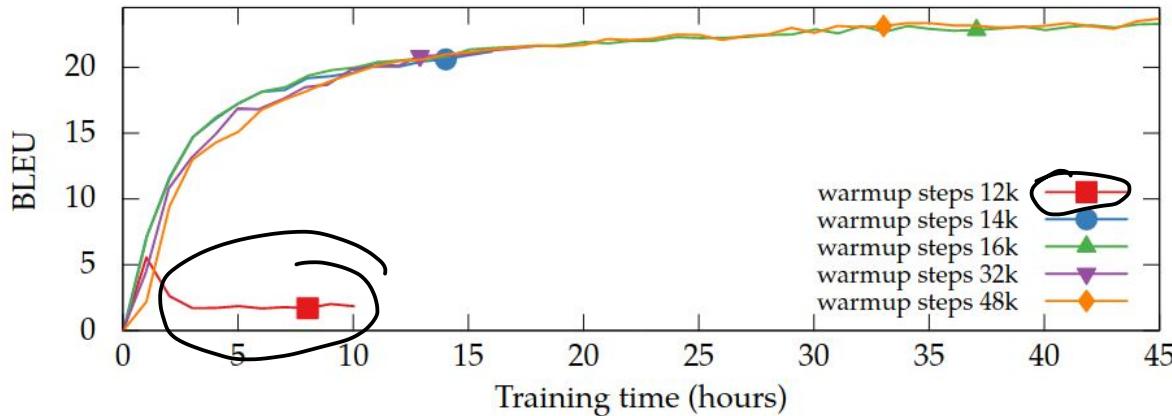


Figure 8: Effect of the warmup steps on a single GPU. All trained on CzEng 1.0 with the default batch size (1500) and learning rate (0.20).