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/girafe-ai/ml-mipt/blob/master/week1_04_Transformer/week04_Transformer.pdf
http://web.stanford.edu/class/cs224n/slides/cs224n-2019-lecture08-nmt.pdf
https://jalammar.github.io/illustrated-transformer/
https://github.com/yandexdataschool/nlp_course

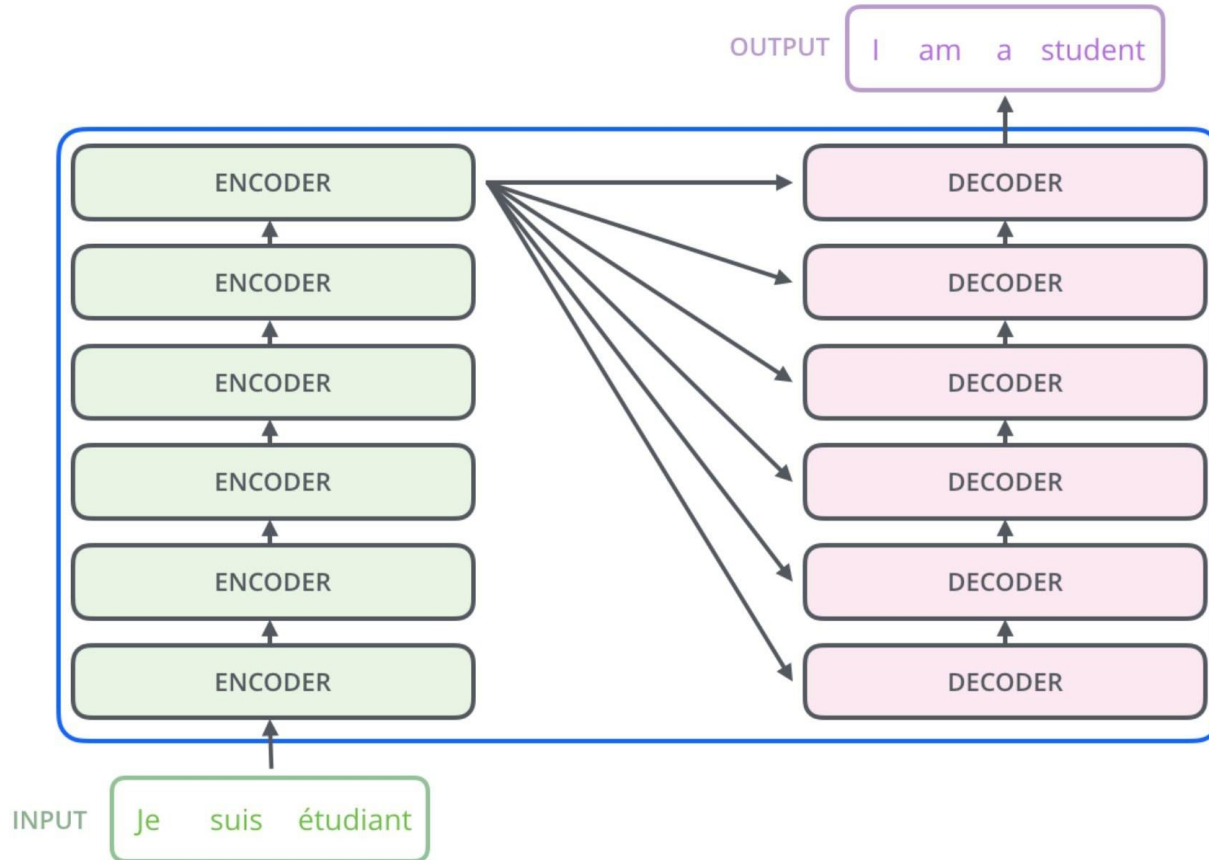# Deep Encoder-Decoder Models (GNMT)



Wu et al. 2016

# The Transformer

Output Probabilities

Softmax

Linear

Add & Norm

Feed Forward

Add & Norm

Multi-Head Attention

Nx

Add & Norm

Feed Forward

Add & Norm

Add & Norm

Multi-Head Attention

Nx

Masked Multi-Head Attention

Positional Encoding

Input Embedding

Inputs

Positional Encoding

Output Embedding

Outputs (shifted right)

Image source: [Attention Is All You Need](), Neural Information Processing Systems 2017

# The Transformer



INPUT

Je suis étudiant

THE TRANSFORMER

OUTPUT

I am a student

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

# The Transformer

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

# The Transformer

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

Can be parallelized

ENCODER

Feed Forward

$z_1$ $z_2$ $z_3$

Self-Attention

$x_1$ Je $x_2$ suis $x_3$ étudiant

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

7

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

# 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

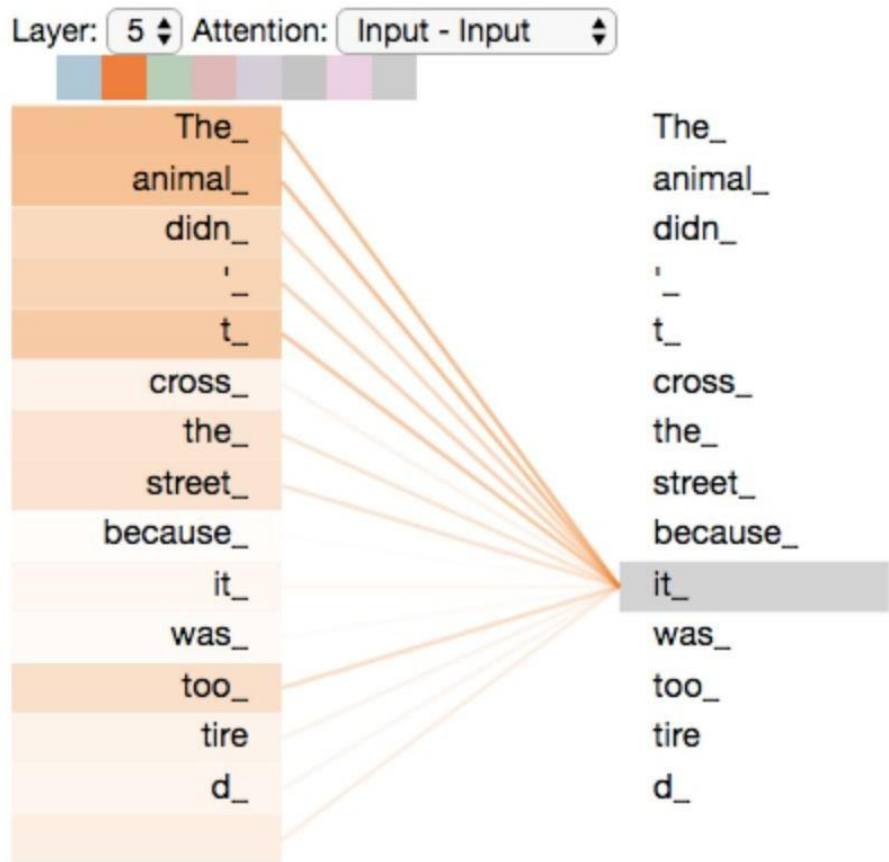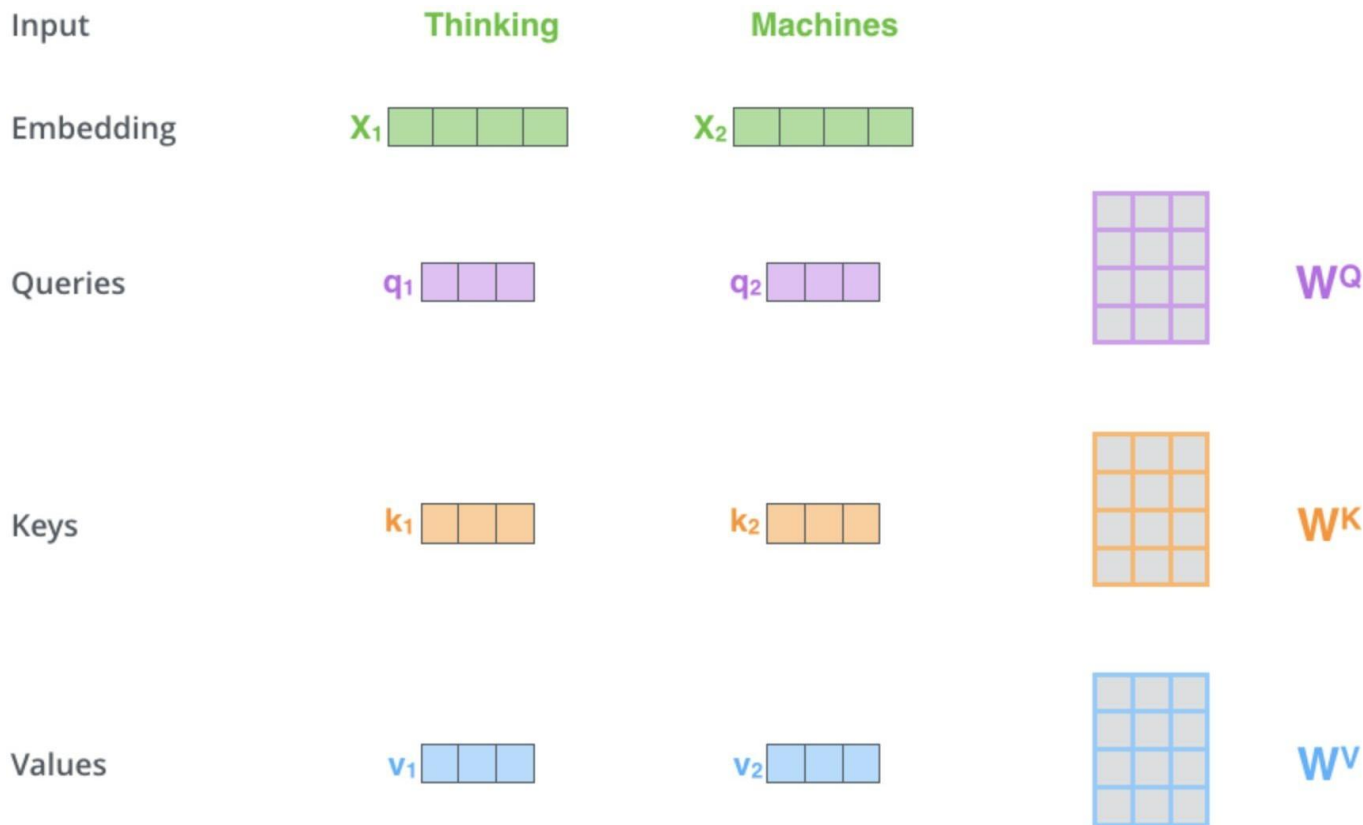- Uses **self-attention** concept

# Self-Attention

"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

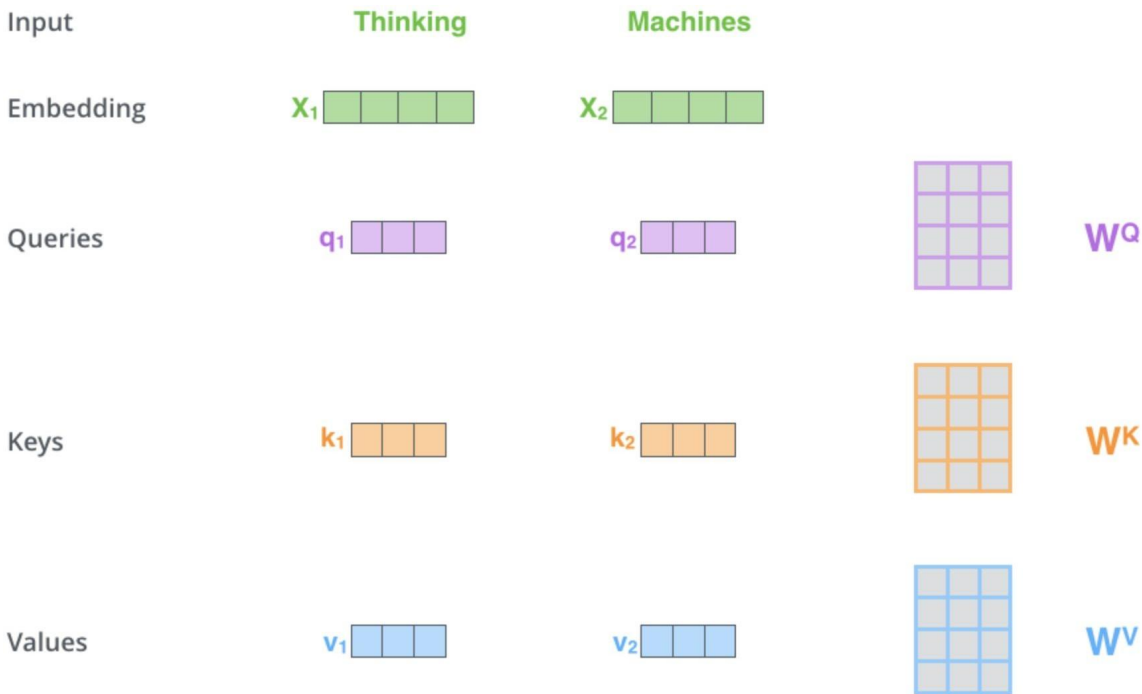# Self-Attention: detailed explanation

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

# Self-Attention: detailed explanation

## STEP 1:

create 3 vectors
(**query**, **key**, **value**)

from each of the encoder's input vectors

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

# Self-Attention: detailed explanation
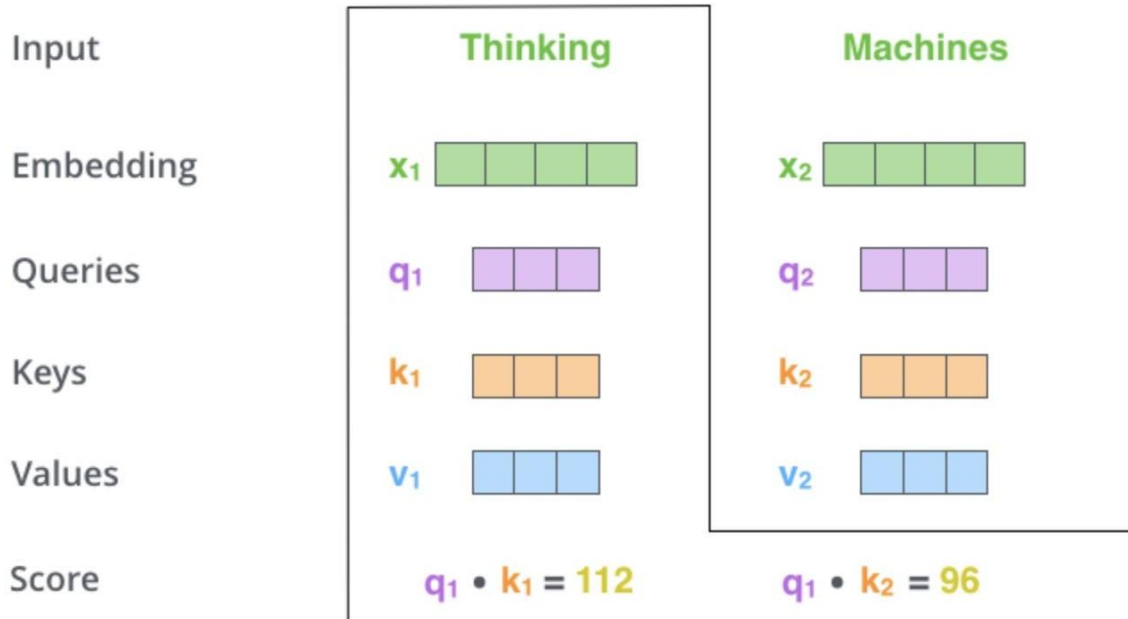
What are the **query**, **key**, **value** vectors?


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)

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

# Self-Attention: detailed explanation

## STEP 3:

divide the scores by 8

(the square root of the dimension of the key vectors)
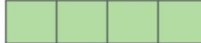
## STEP 4:

softmax



| | Thinking | Machines |
|---|---|---|
| Input | | |
| Embedding | $x_1$ | $x_2$ |
| Queries | $q_1$ | $q_2$ |
| Keys | $k_1$ | $k_2$ |
| Values | $v_1$ | $v_2$ |
| Score | $q_1 \cdot k_1 = 112$ | $q_1 \cdot k_2 = 96$ |
| Divide by 8 ( $\sqrt{d_k}$ ) | 14 | 12 |
| Softmax | 0.88 | 0.12 |

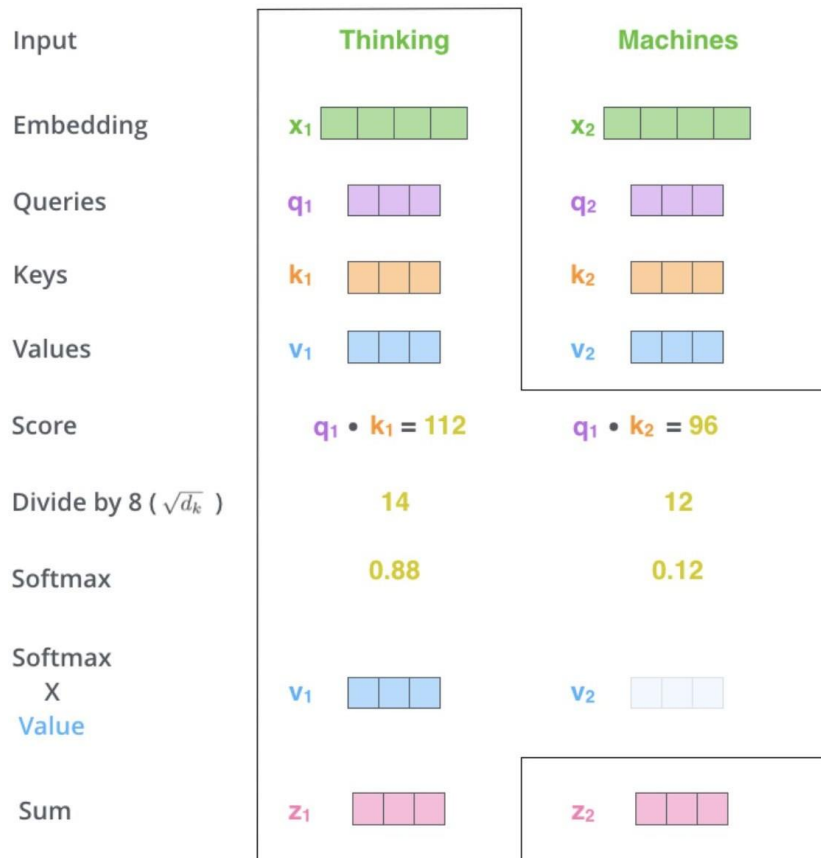Image source: https://jalammar.github.io/illustrated-transformer/

# Self-Attention: detailed explanation

**STEP 5:**

multiply each value
vector by the softmax
score

**STEP 6:**

sum up the weighted
value vectors



| | Thinking | Machines |
|---|---|---|
| Input | | |
| Embedding | $x_1$ | $x_2$ |
| Queries | $q_1$ | $q_2$ |
| Keys | $k_1$ | $k_2$ |
| Values | $v_1$ | $v_2$ |
| Score | $q_1 \cdot k_1 = 112$ | $q_1 \cdot k_2 = 96$ |
| Divide by 8 ($\sqrt{d_k}$) | 14 | 12 |
| Softmax | 0.88 | 0.12 |
| Softmax X Value | $v_1$ | $v_2$ |
| Sum | $z_1$ | $z_2$ |

# Self-Attention



| | Thinking | Machines |
|---|---|---|
| Input | | |
| Embedding | $x_1$ ▢▢▢▢ | $x_2$ ▢▢▢▢ |
| Queries | $q_1$ ▢▢▢ | $q_2$ ▢▢▢ |
| Keys | $k_1$ ▢▢▢ | $k_2$ ▢▢▢ |
| Values | $v_1$ ▢▢▢ | $v_2$ ▢▢▢ |
| Score | $q_1 \cdot k_1 = 112$ | $q_1 \cdot k_2 = 96$ |
| Divide by 8 ( $\sqrt{d_k}$ ) | 14 | 12 |
| Softmax | 0.88 | 0.12 |
| Softmax X Value | $v_1$ ▢▢▢ | $v_2$ ▢▢▢ |
| Sum | $z_1$ ▢▢▢ | $z_2$ ▢▢▢ |

**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

30

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

# Self-Attention: Matrix Calculation

Pack embeddings into matrix **X**

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

# Self-Attention: Matrix Calculation

# Multi-Head Attention

# Multi-Head Attention

# Multi-Head Attention

1) Concatenate all the attention heads

$Z_0$ $Z_1$ $Z_2$ $Z_3$ $Z_4$ $Z_5$ $Z_6$ $Z_7$

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

X

$W^O$

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

$Z$

=

Image source: https://jalammar.github.io/illustrated-transformer/
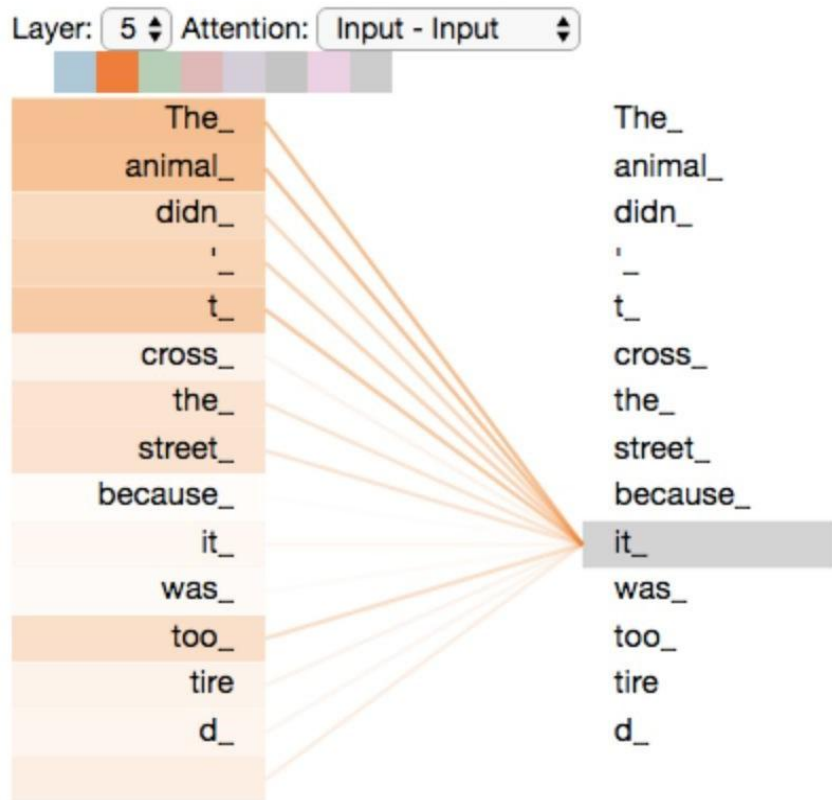
1) This is our input sentence*

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

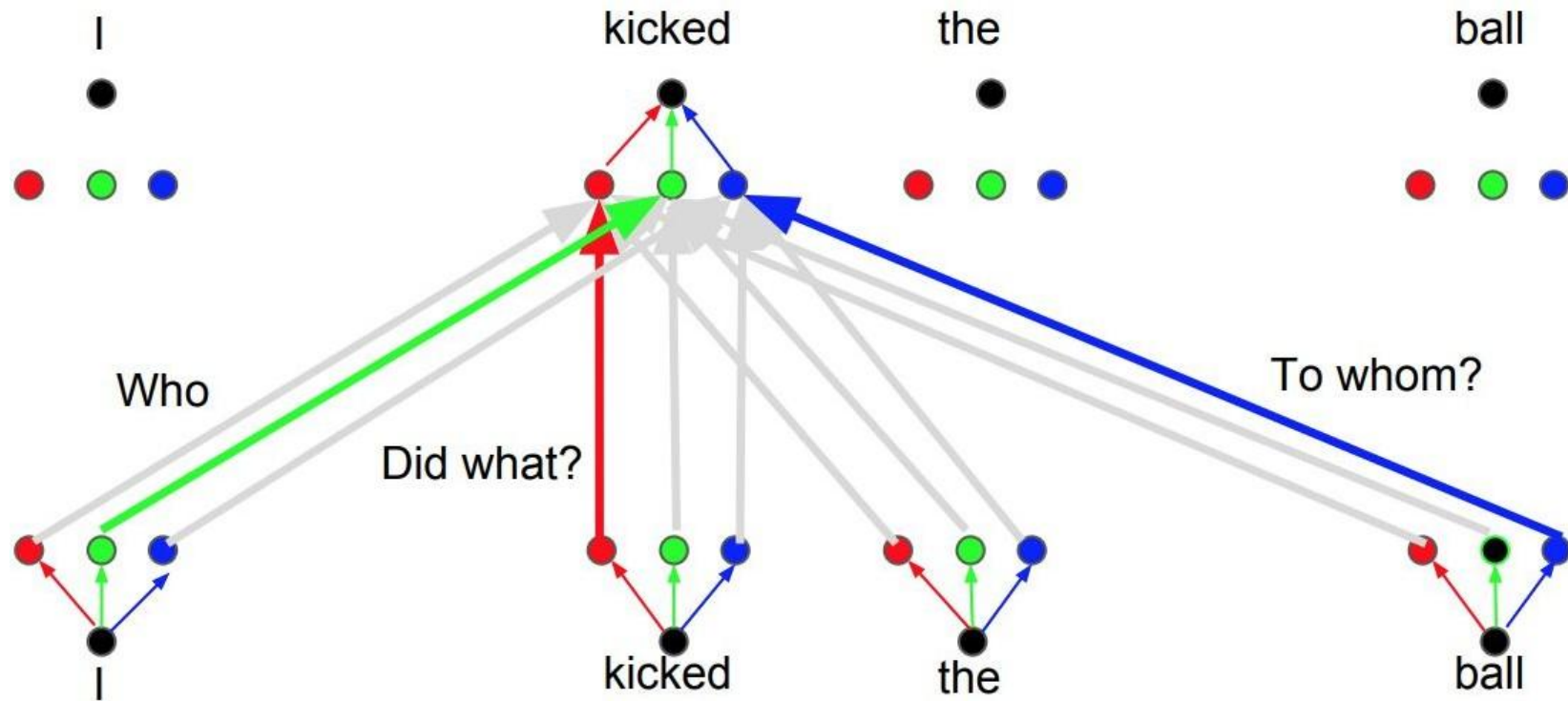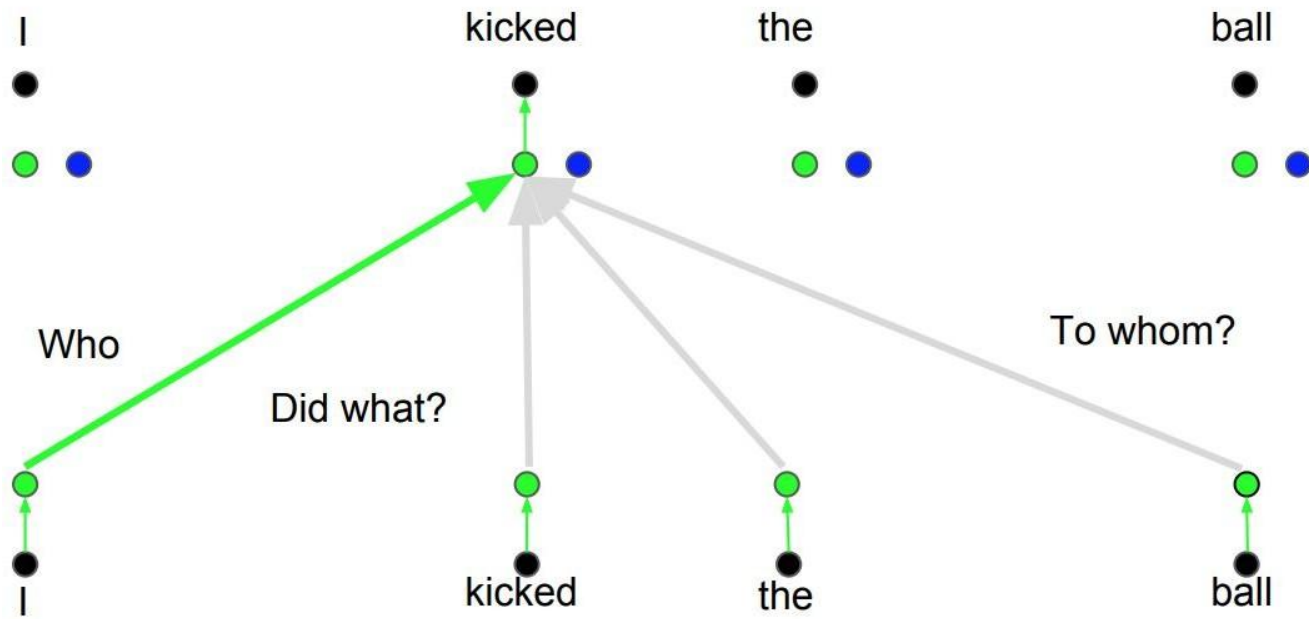Thinking Machines

X

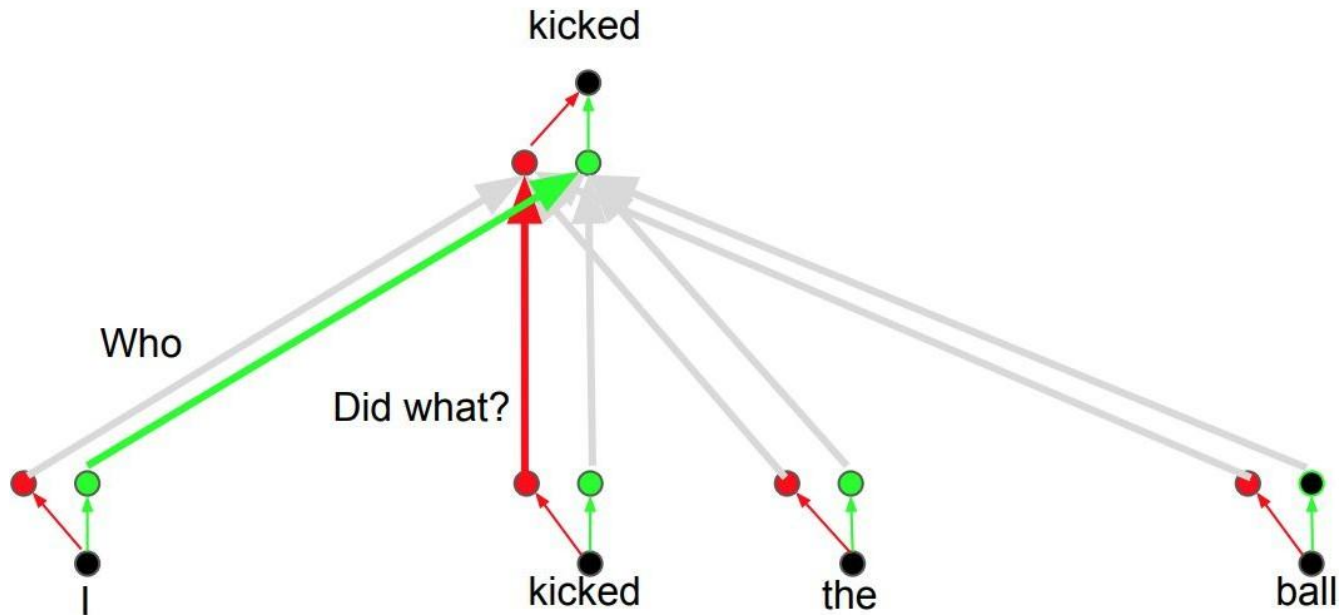* In all encoders other than #0, we don't need embedding. We start directly with the output of the encoder right below this one

R

$W_0^Q$
$W_0^K$
$W_0^V$

$W_1^Q$
$W_1^K$
$W_1^V$

$W_7^Q$
$W_7^K$
$W_7^V$

$Q_0$
$K_0$
$V_0$

$Q_1$
$K_1$
$V_1$

$Q_7$
$K_7$
$V_7$

$Z_0$

$Z_1$

$Z_7$

$W^O$

$Z$

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

# Multi-Head Attention

# Why Multi-Head Attention?



I kicked the ball

Who

Did what?

To whom?

I kicked the ball

# Attention head: Who

# Attention head: Did What?

# Attention head: To Whom?



kicked

Who

Did what?

To whom?

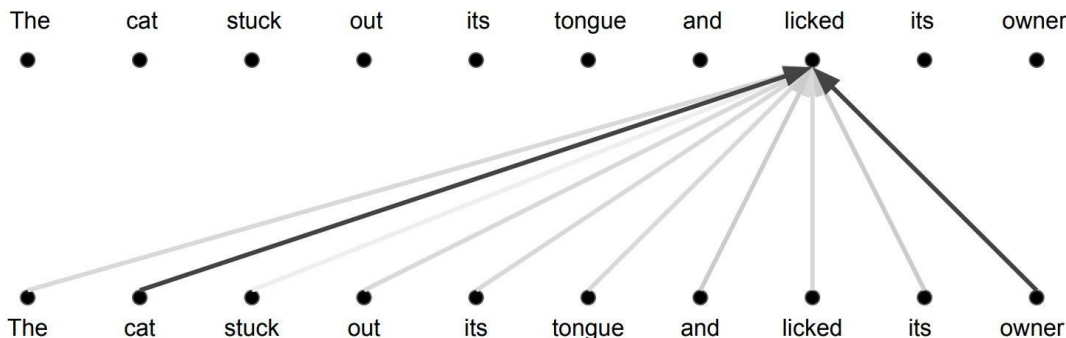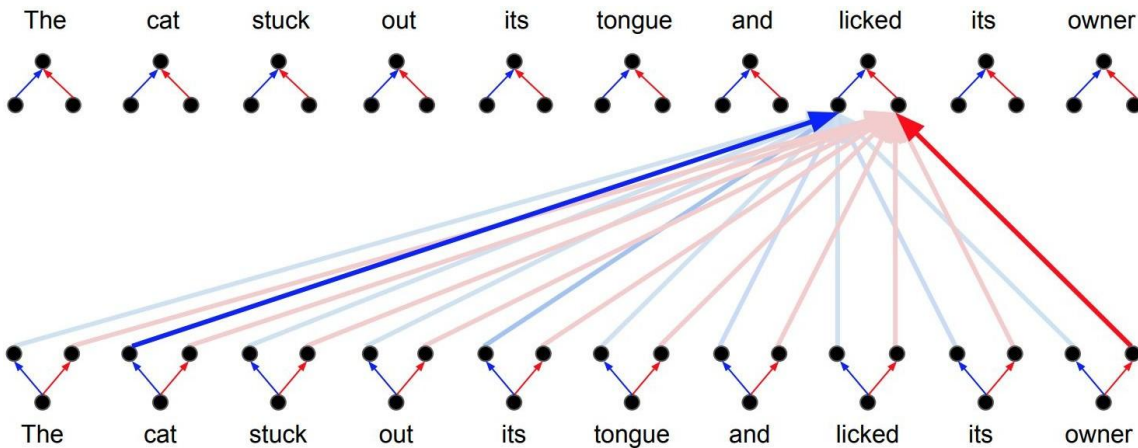I          kicked          the          ball

# Attention vs. Multi-Head Attention

**Attention:** a weighted average

**Multi-Head Attention:**

parallel attention layers with different linear transformations on input and output.
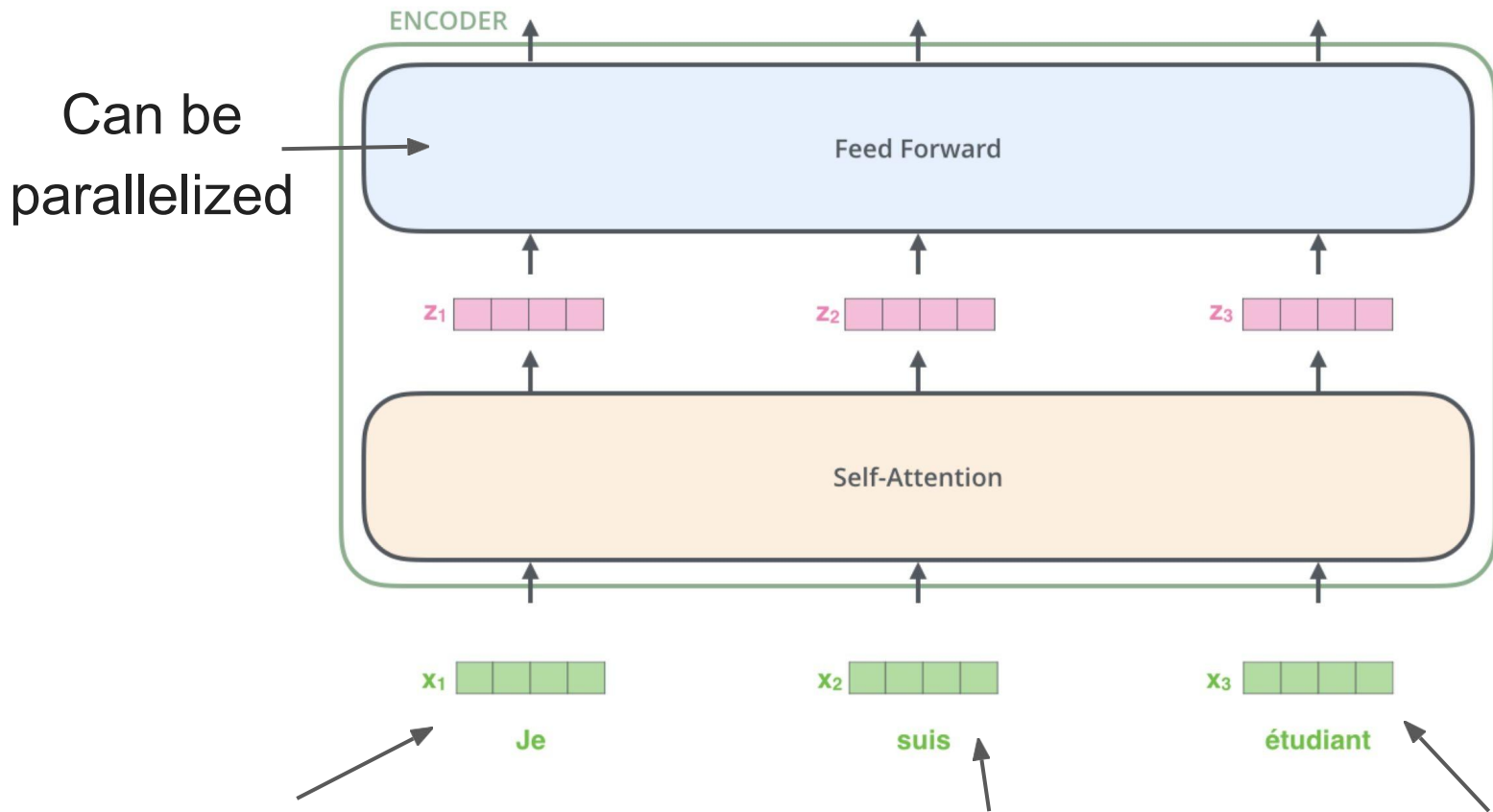
# Performance: WMT 2014 BLEU

|  | 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

ENCODER

Feed Forward

$z_1$ $z_2$ $z_3$

Self-Attention

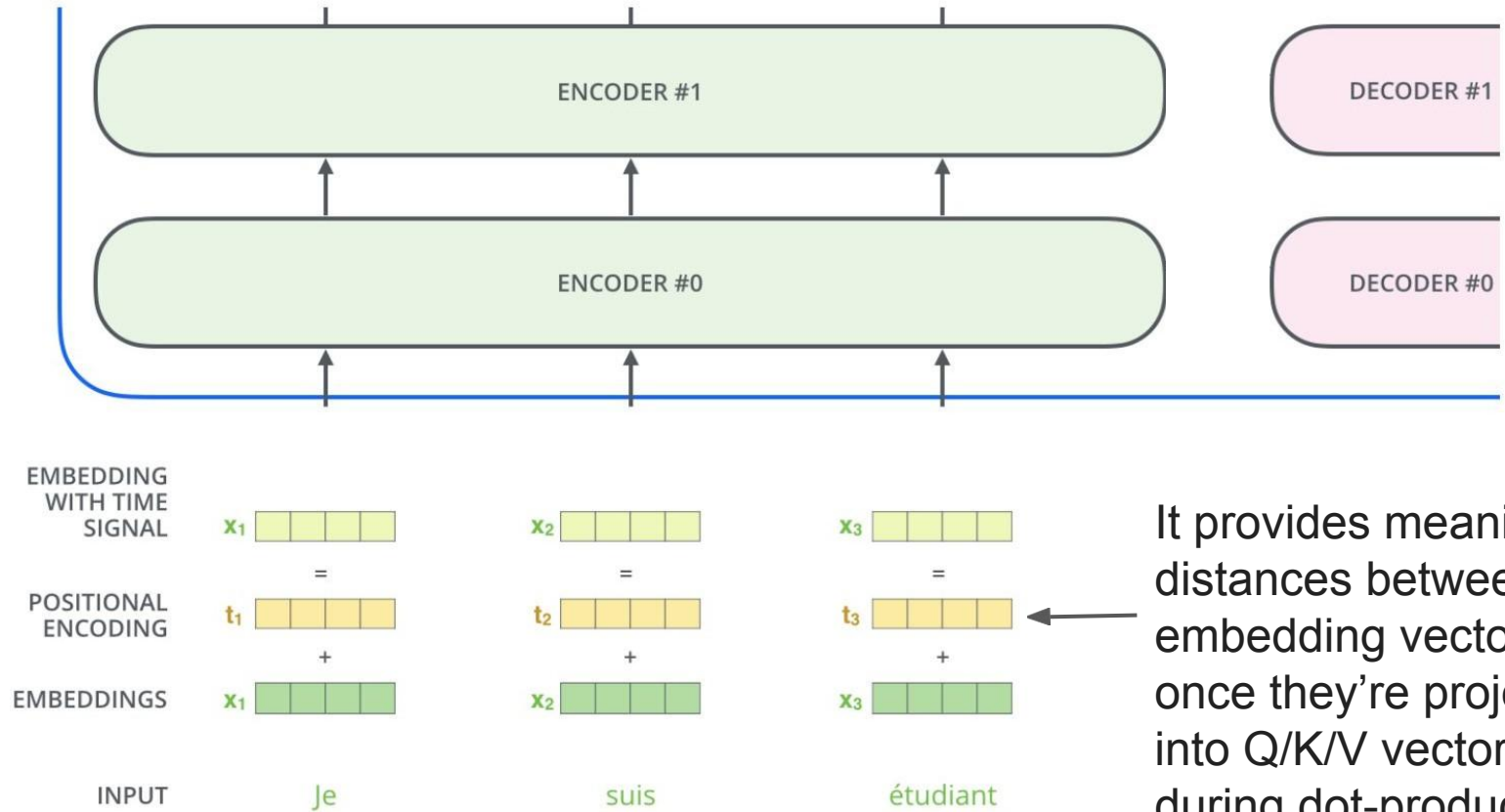$x_1$ Je  $x_2$ suis  $x_3$ étudiant

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



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

# Positional Encoding: why sin and cos?

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

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

$$\vec{p_t} = \begin{bmatrix} \sin(\omega_1.t) \\ \cos(\omega_1.t) \\ \\ \sin(\omega_2.t) \\ \cos(\omega_2.t) \\ \\ \vdots \\ \\ \sin(\omega_{d/2}.t) \\ \cos(\omega_{d/2}.t) \end{bmatrix}_{d \times 1}$$

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

Image source: https://kazemnejad.com/blog/transformer_architecture_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.*

$$M \begin{bmatrix} \sin(\omega_k t) \\ \cos(\omega_k t) \end{bmatrix} = \begin{bmatrix} \sin(\omega_k (t + \phi)) \\ \cos(\omega_k (t + \phi)) \end{bmatrix}$$
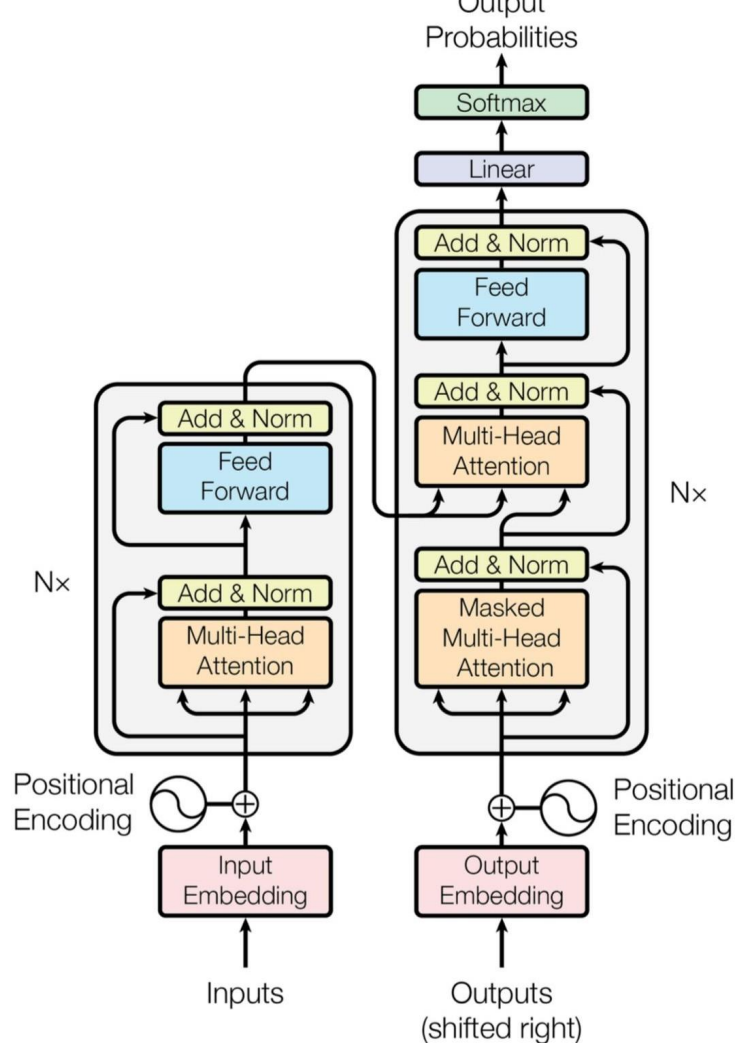
# Positional Encoding: why sin and cos?

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

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

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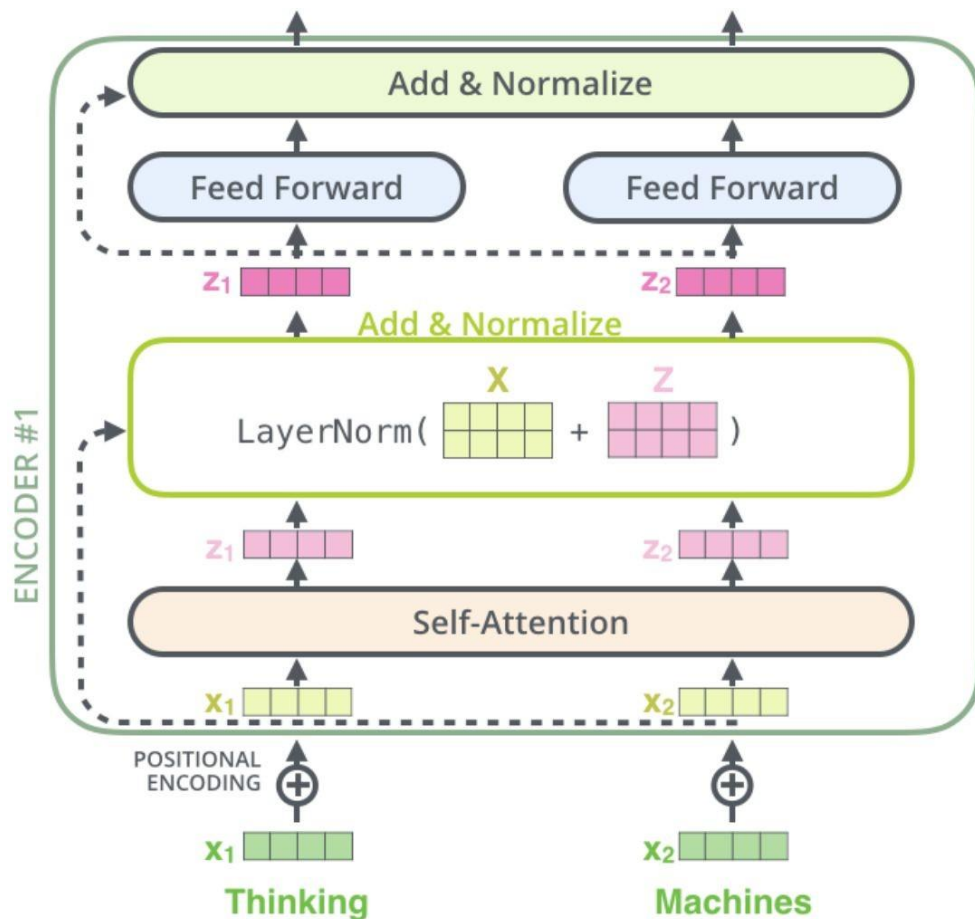
# Layer Normalization

# The Transformer: recap

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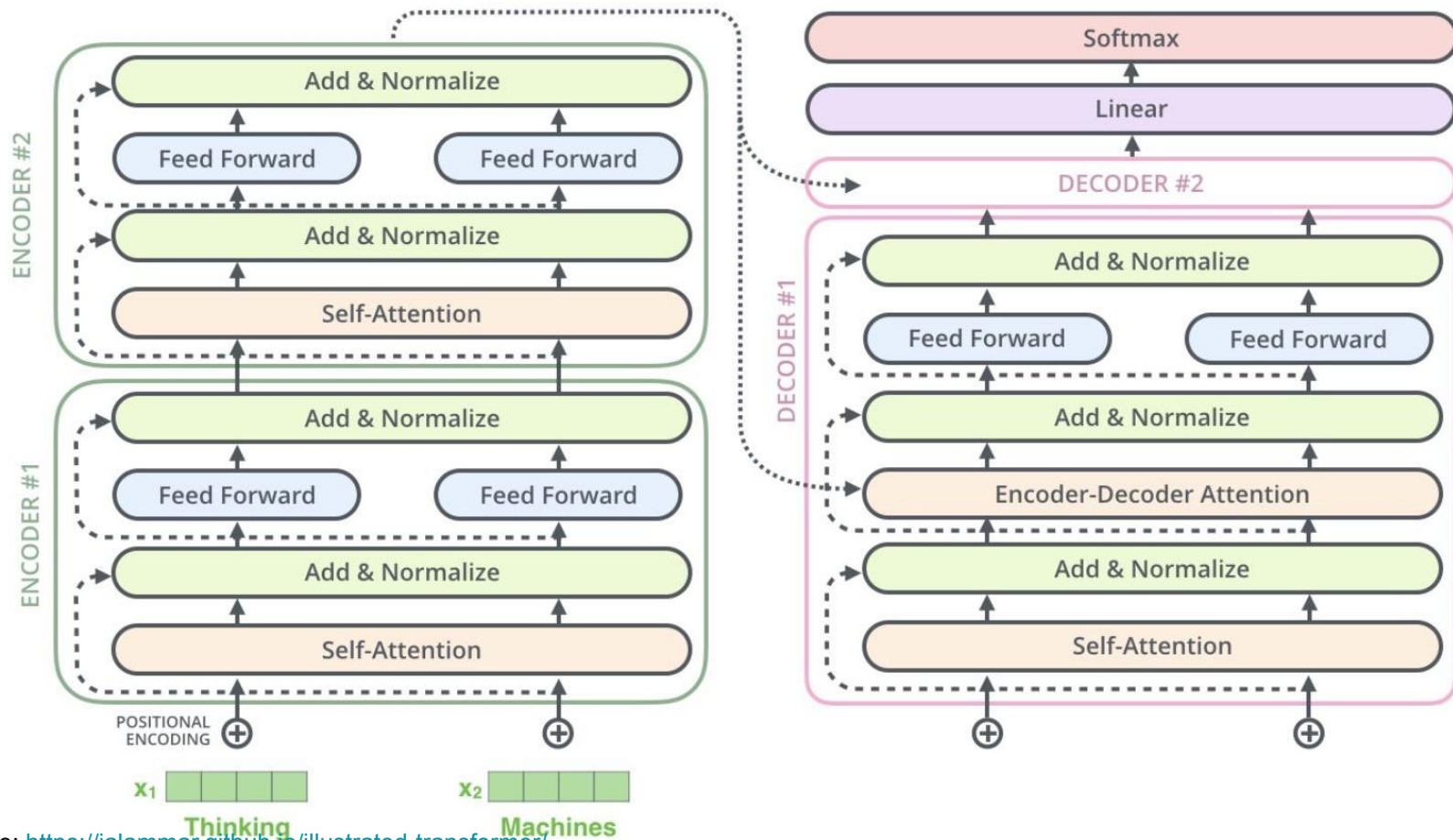
# Layer Normalization

# Layer Normalization

Like BatchNorm

but normalize along all features representing latent vector

More info:
Layer Normalization
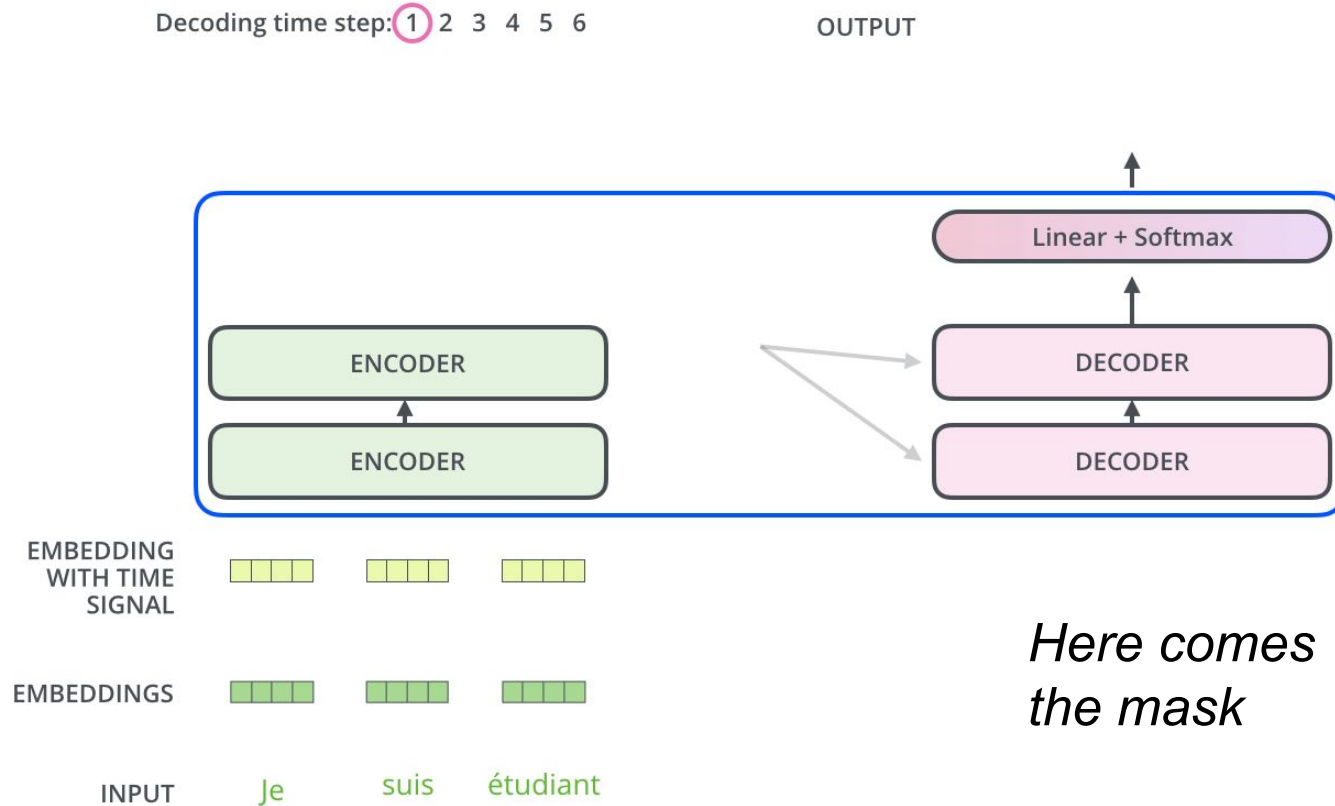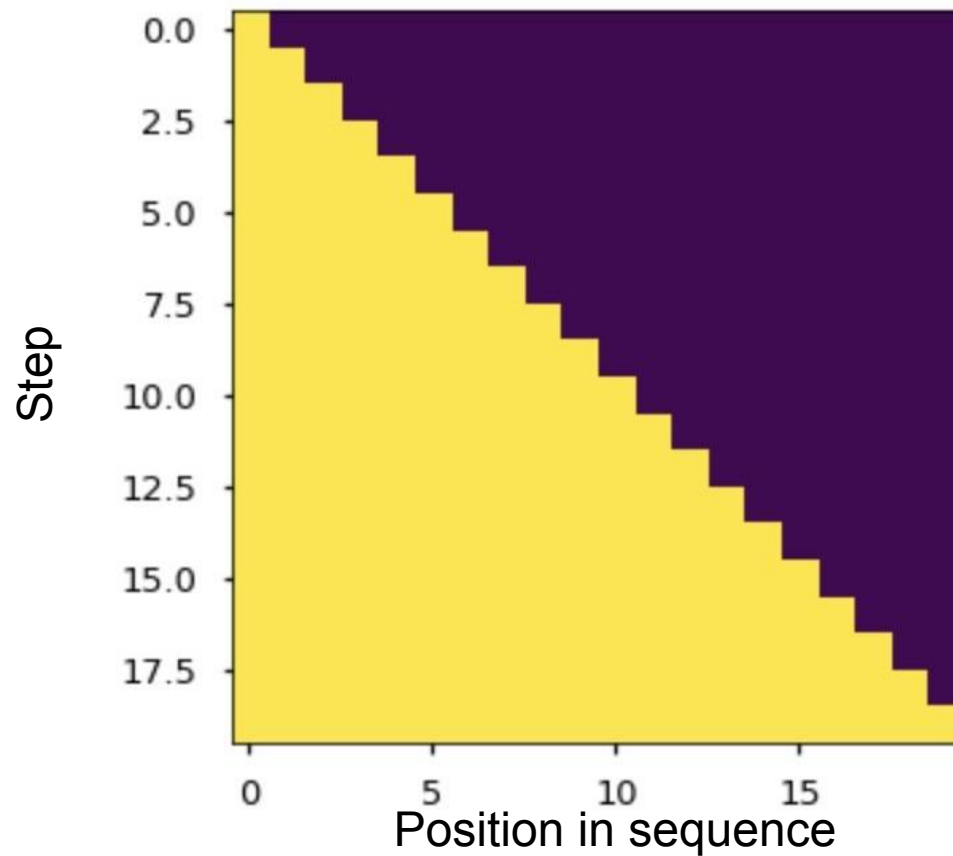
# Layer Normalization

# The Decoder

# The Decoder Side
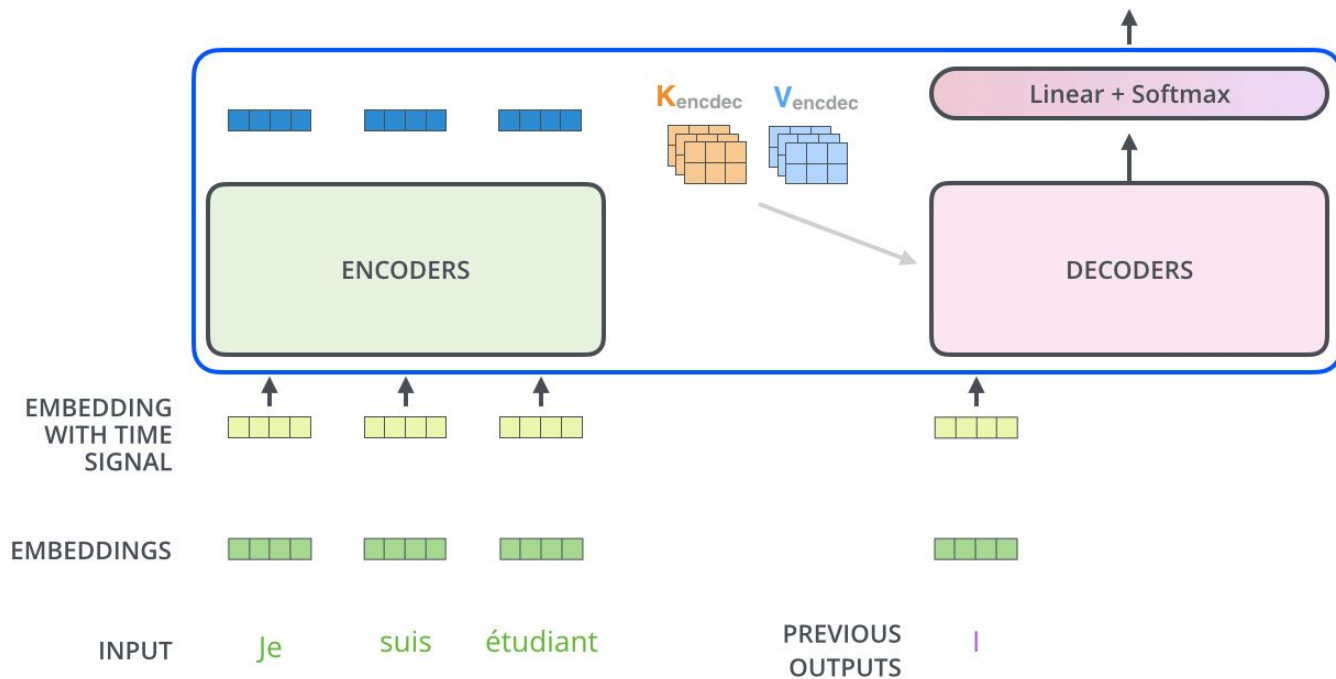


Image source: https://jalammar.github.io/illustrated-transformer/

# The masked decoder input

# The Decoder Side



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

# Final Linear and Softmax Layer



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

# The Transformer



Output Probabilities

Softmax

Linear

Add & Norm

Feed Forward

Add & Norm

Multi-Head Attention

Add & Norm

Masked Multi-Head Attention

Add & Norm

Feed Forward

Add & Norm

Multi-Head Attention

N×

N×

Positional Encoding

Positional Encoding

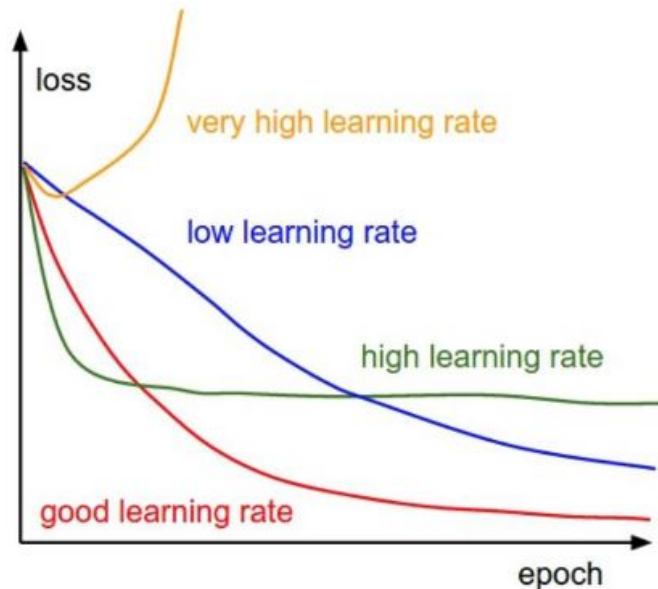Input Embedding

Output Embedding

Inputs

Outputs (shifted right)

# Training The Transformer

➤ Consider Adam optimizer

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

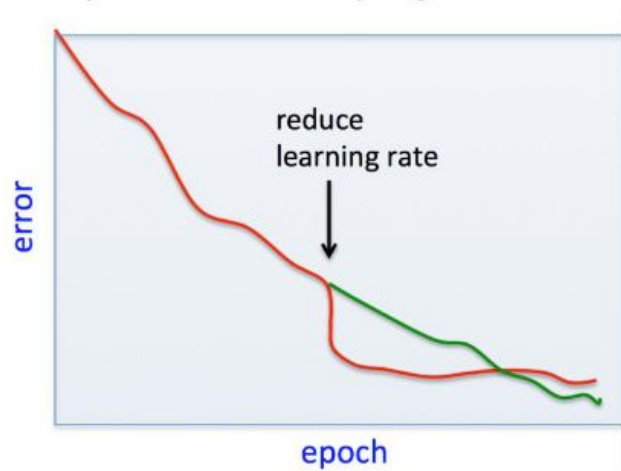$$v_{t+1} = \beta_2 \cdot v_t + (1 - \beta_2) \cdot || \nabla_\theta L ||^2$$

$$\theta_{t+1} = \theta_t - \frac{\alpha}{\sqrt{v_t} + \epsilon} \mu_t$$
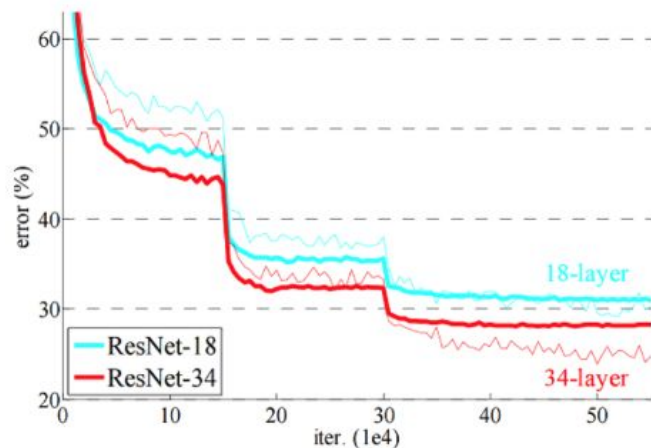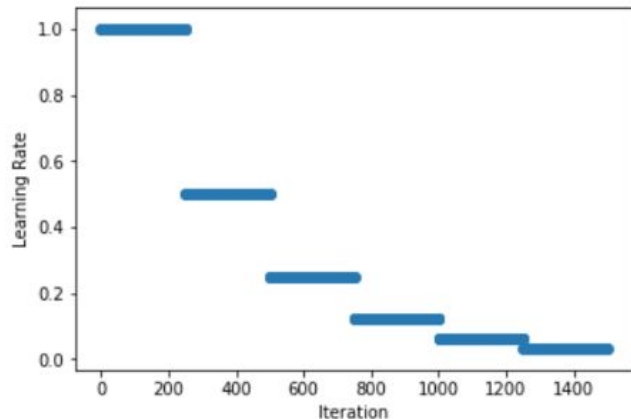
➤ **The choice of α is crucial!**

➤ 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
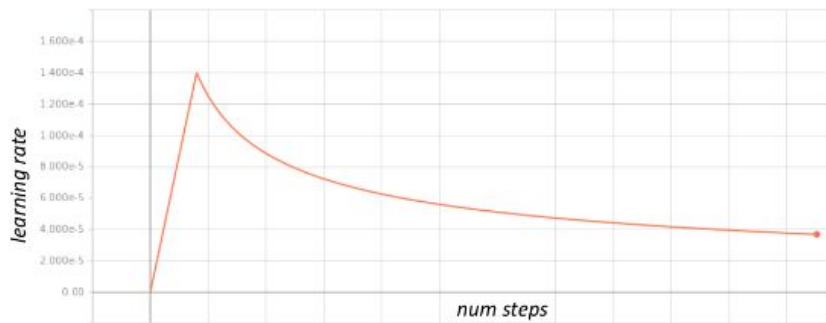   ○ every **k** steps or whenever progress slows down

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

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

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

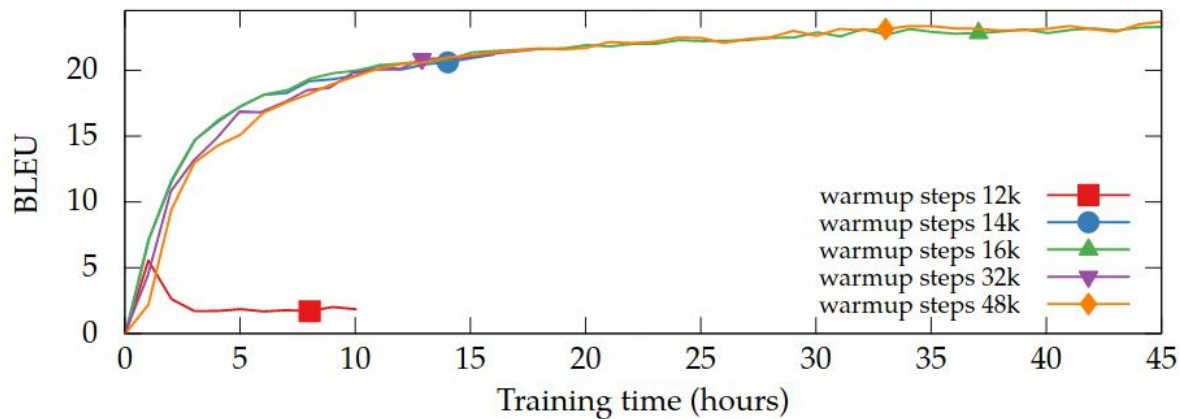Transformers trainings requires non-trivial efforts



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).