# Lecture 04: Self-Attention & Transformer overview

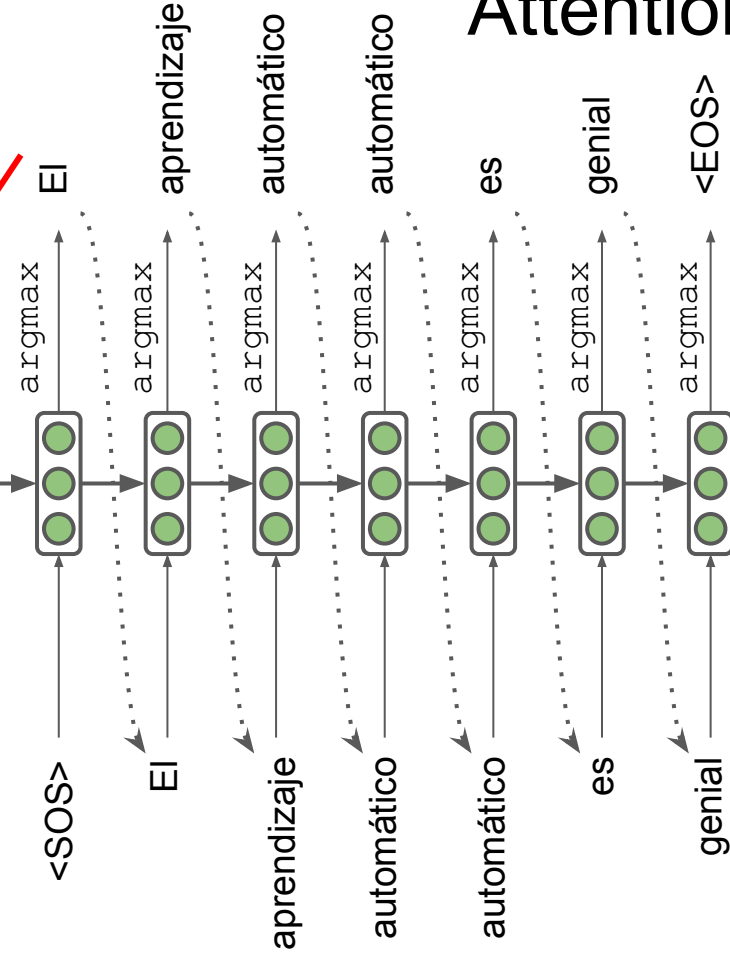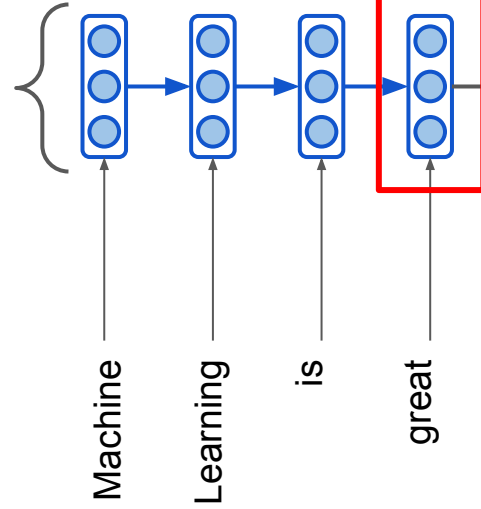**Radoslav Neychev**

Fall 2021, Moscow, Russia

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

Based on: http://web.stanford.edu/class/cs224n/slides/cs224n-2019-lecture08-nmt.pdf
https://jalammar.github.io/illustrated-transformer/

Attention in seq2seq

# Seq2seq with attention

Attention scores

Encoder

# Seq2seq with attention

Attention
distribution

Attention
scores

Encoder

Simply apply softmax to
scores

Seq2seq with attention

Attention output

Weighted sum of all encoder states

Attention distribution

Attention scores

Encoder

Seq2seq with attention

Attention output

Attention distribution

Attention scores

Encoder

Concatenate

Seq2seq with attention

Attention output

Attention distribution

Attention scores

Encoder

y

Seq2seq with attention

Attention output

Attention distribution

Attention scores

Encoder

y

9

Attention
output

Attention
distribution

Attention
scores

Encoder

y

Denote encoder hidden states $\mathbf{h}_1, \ldots, \mathbf{h}_N \in \mathbb{R}^k$

and decoder hidden state at time step t $\quad \mathbf{s}_t \in \mathbb{R}^k$

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

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

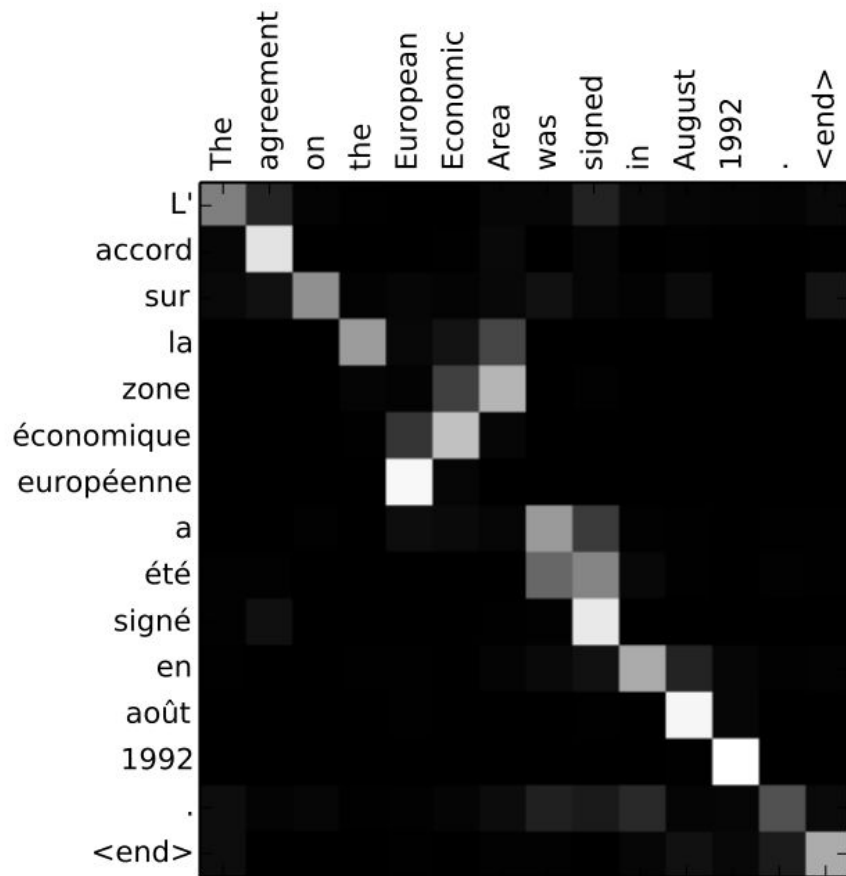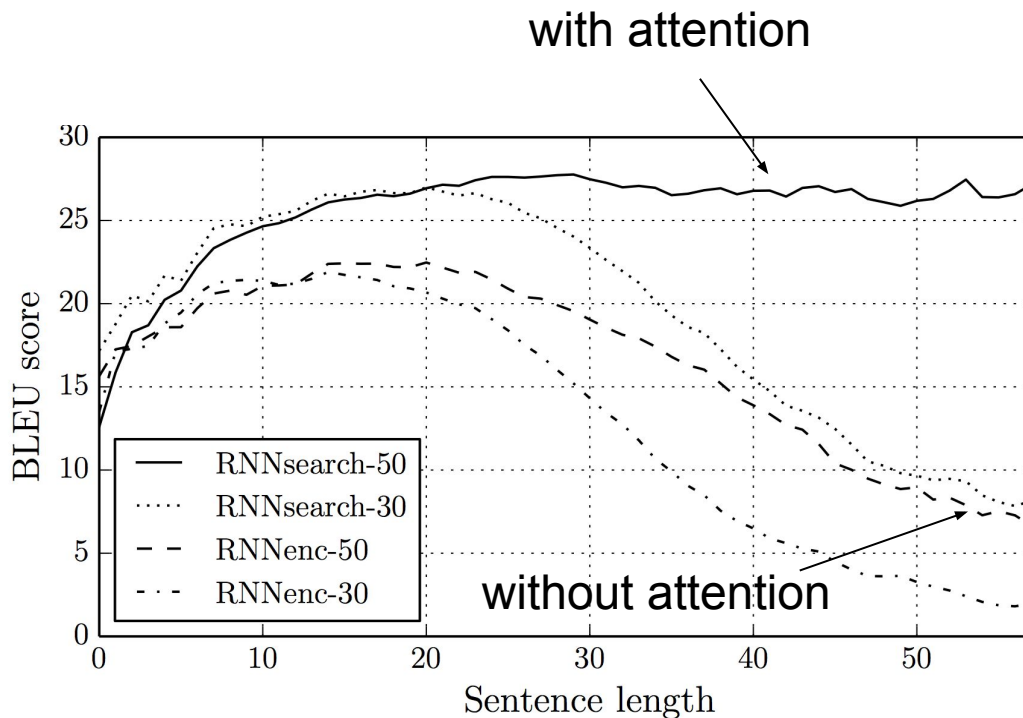Then the attention vector is a linear combination of encoder states

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

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

# Attention advantages

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

with attention

without attention



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

# The Transformer

# The Transformer



INPUT

Je suis étudiant

THE TRANSFORMER

OUTPUT

I am a student

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

$x_1$  $x_2$  $x_3$

Je  suis  étudiant

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

18

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

Can be parallelized

r₁ → Feed Forward Neural Network

r₂ → Feed Forward Neural Network

z₁ z₂

Self-Attention

x₁ Thinking

x₂ Machines

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

19

# 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*
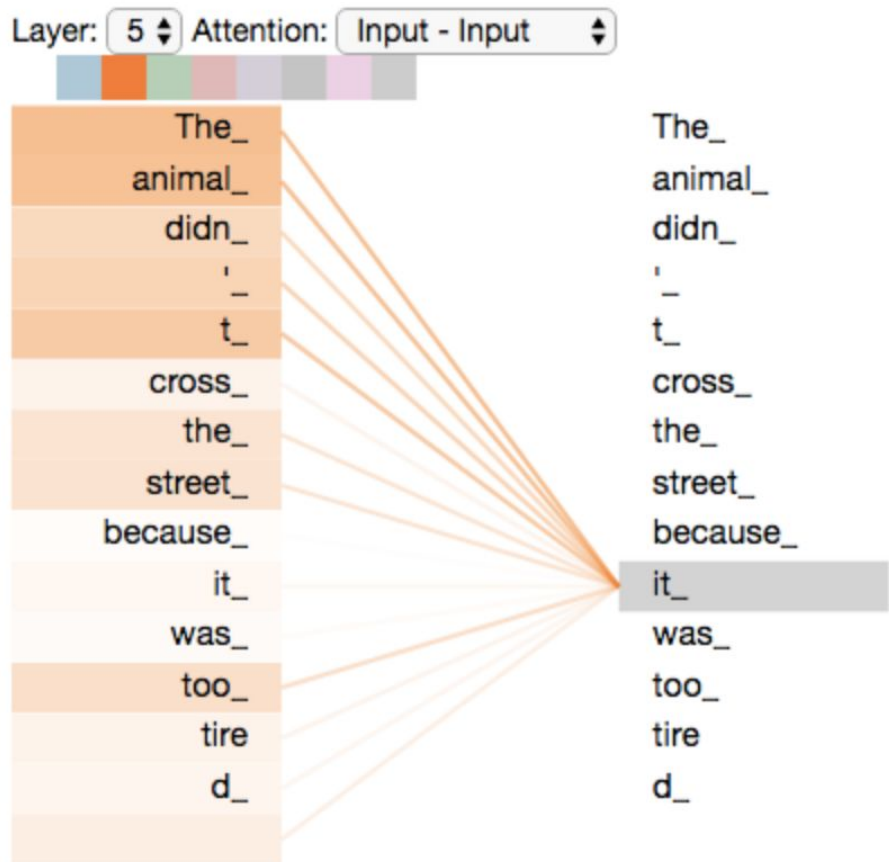- Much faster
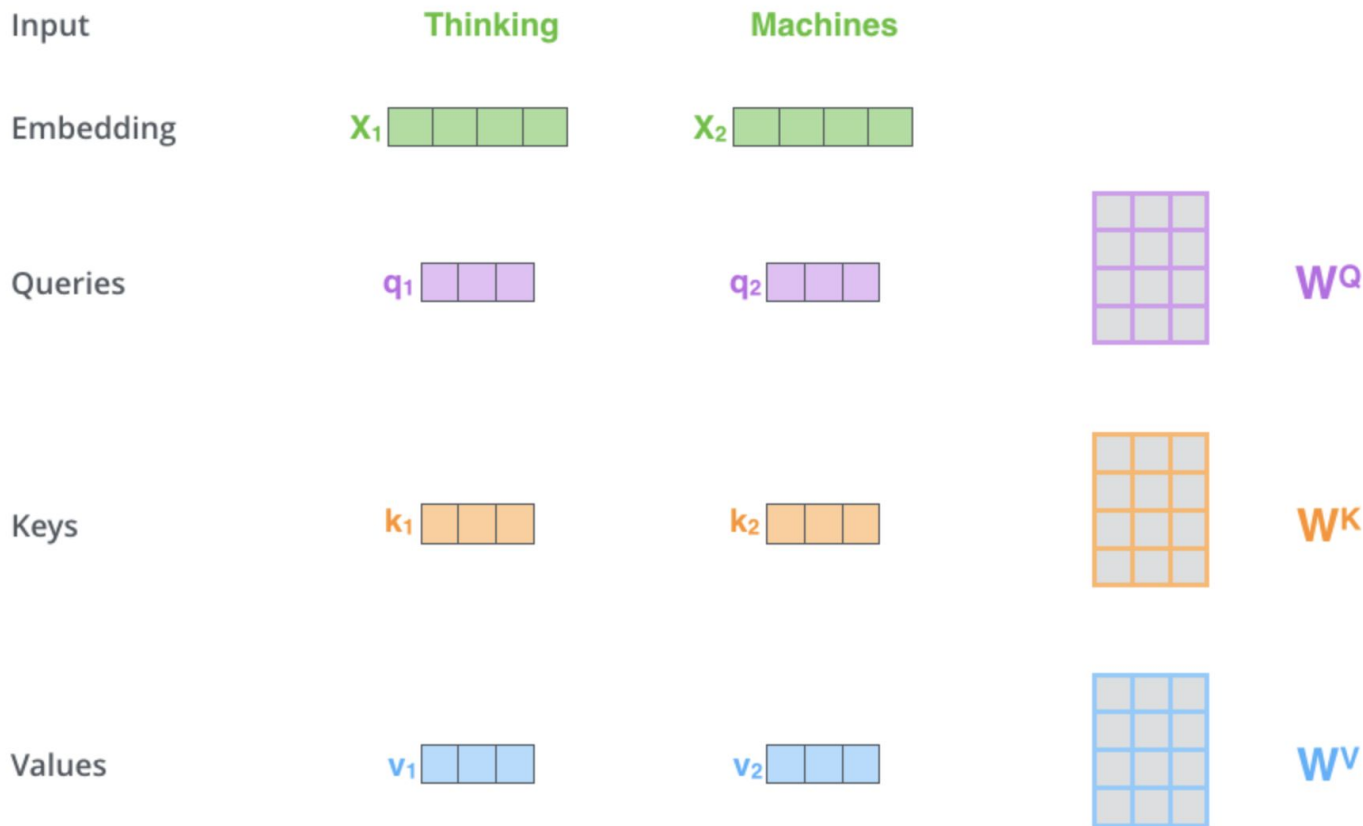- Uses **self-attention** concept

# 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

Layer: [ 5 ⬍ ] Attention: [ Input - Input ⬍ ]

| Left | Right |
|------|-------|
| The_ | The_ |
| animal_ | animal_ |
| didn_ | didn_ |
| '_ | '_ |
| t_ | t_ |
| cross_ | cross_ |
| the_ | the_ |
| street_ | street_ |
| because_ | because_ |
| it_ | it_ |
| was_ | was_ |
| too_ | too_ |
| tire | tire |
| d_ | d_ |

# Self-Attention: detailed explanation

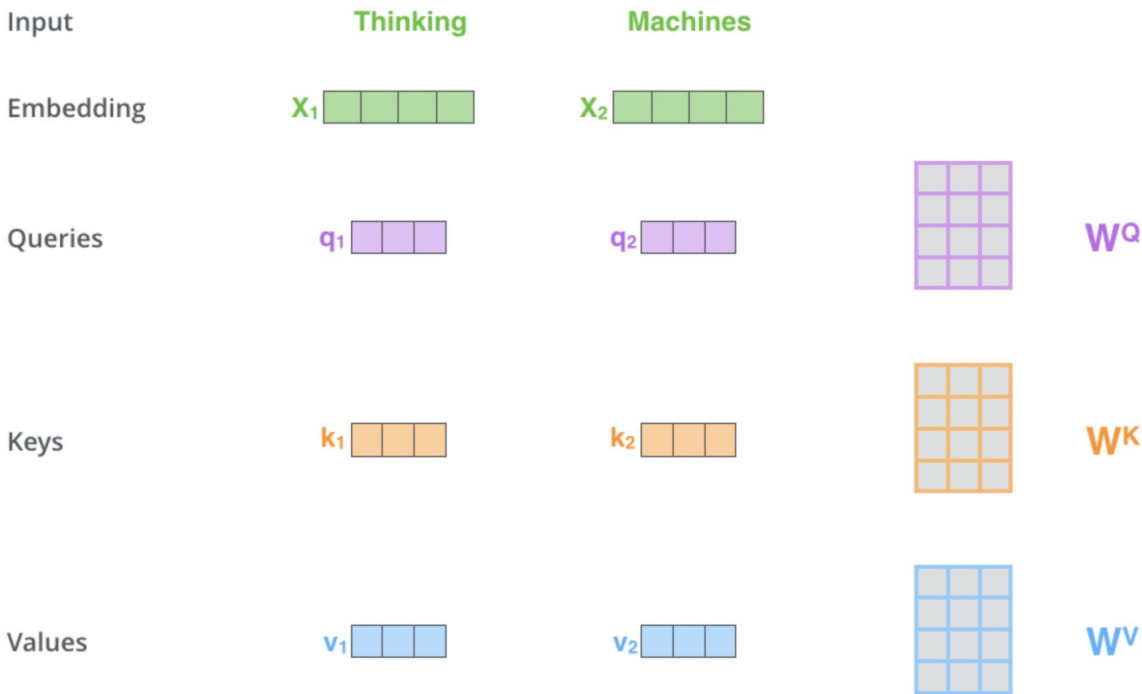Image source: https://jalammar.github.io/illustrated-transformer/

# 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

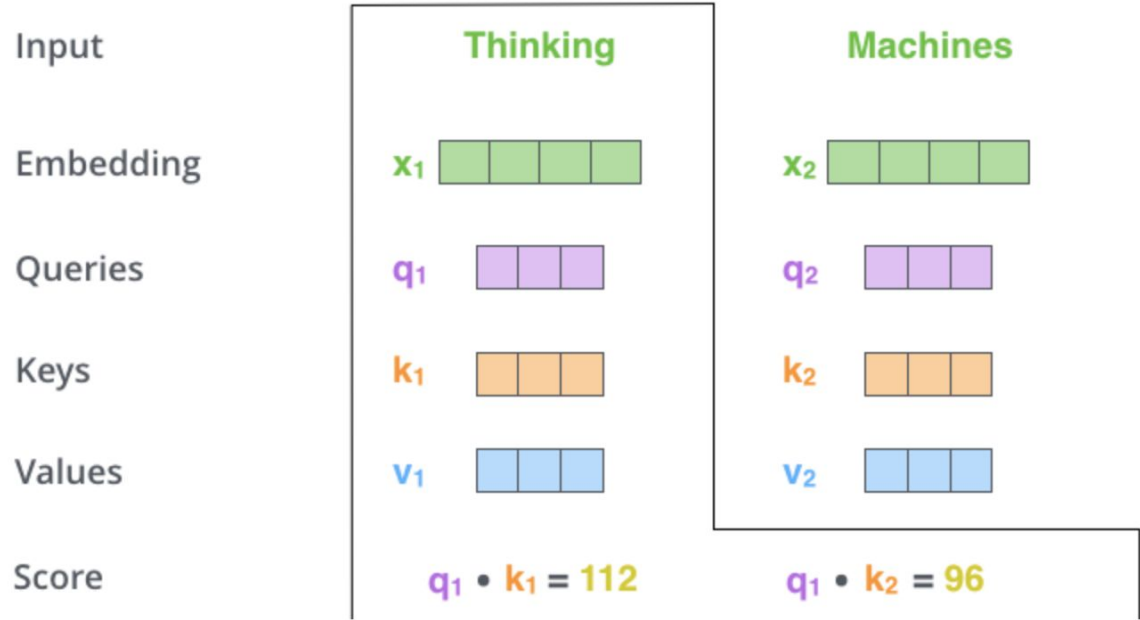What are the **query**, **key**, **value** vectors?

They're abstractions that are useful for

calculating and thinking about attention.

# Self-Attention: detailed explanation

## STEP 2:

calculate a score

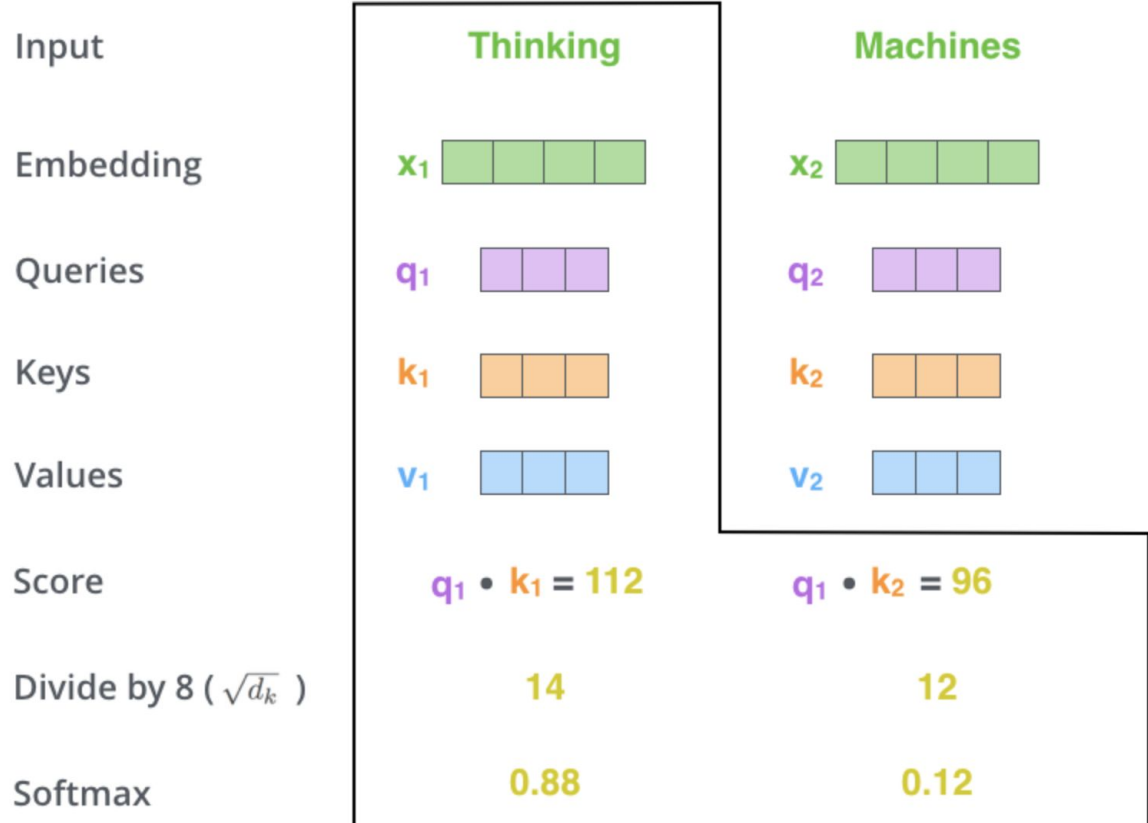(score each word of the input sentence against the current word)

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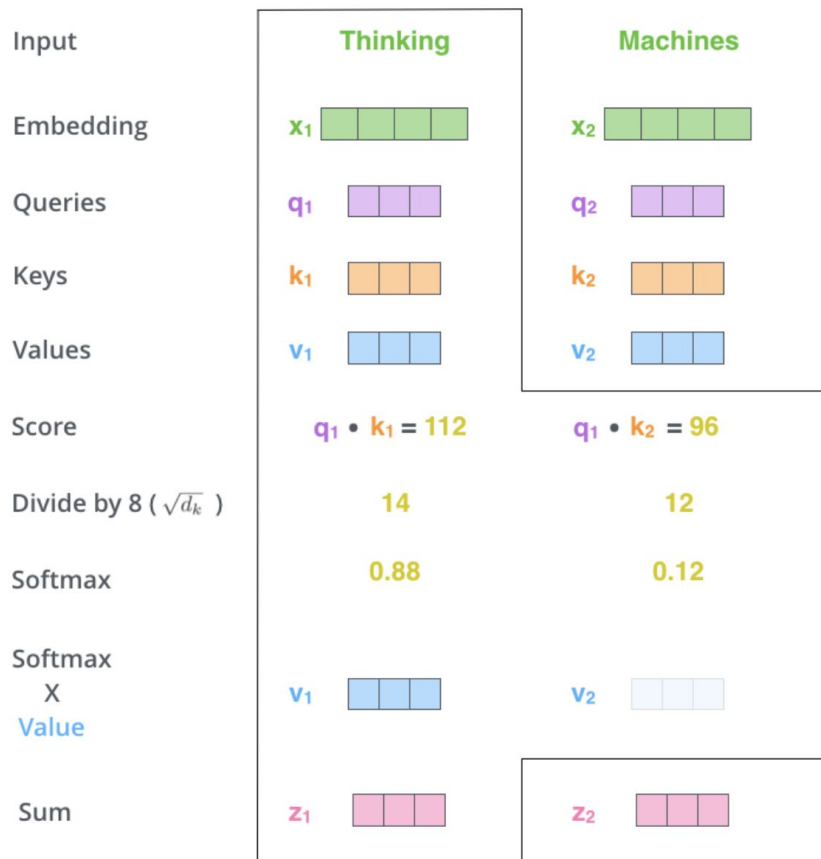
| | Thinking | Machines |
|---|---|---|
| Input | | |
| Embedding | $x_1$ | $x_2$ |
| Queries | $q_1$ | $q_2$ |
| Keys | $k_1$ | $k_2$ |
| Values | $v_1$ | $v_2$ |
| Score | $q_1 \cdot k_1 = 112$ | $q_1 \cdot k_2 = 96$ |
| Divide by 8 ($\sqrt{d_k}$) | 14 | 12 |
| Softmax | 0.88 | 0.12 |

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



| Input | Thinking | Machines |
|---|---|---|
| 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 | $z_1$ | $z_2$ |

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

# Self-Attention

| | Thinking | Machines |
|---|---|---|
| Input | | |
| Embedding | $x_1$ | $x_2$ |
| Queries | $q_1$ | $q_2$ |
| Keys | $k_1$ | $k_2$ |
| Values | $v_1$ | $v_2$ |

**STEP 1:** create Query, Key, Value

| | | |
|---|---|---|
| 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 |

**STEP 2:** calculate scores

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

**STEP 4:** softmax

| | | |
|---|---|---|
| Softmax X Value | $v_1$ | $v_2$ |
| Sum | $z_1$ | $z_2$ |

**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**)

# Self-Attention: Matrix Calculation



$$\text{softmax}\left(\frac{Q \times K^\mathsf{T}}{\sqrt{d_k}}\right) V$$

$$Z =$$

# Multi-Head Attention



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

# 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
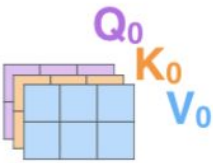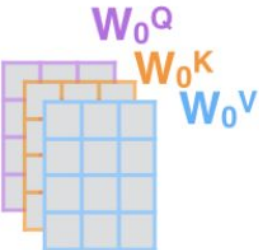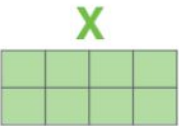
Z

=

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
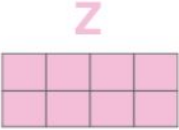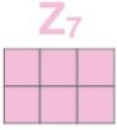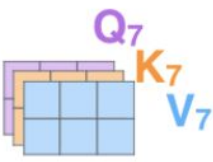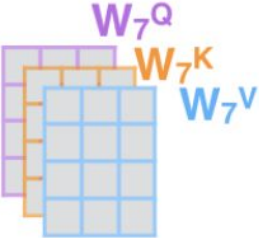
Thinking Machines

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

# Attention head: Who

# Attention head: Did What?

# Attention head: To Whom?

# Attention vs. Multi-Head Attention

**<u>Attention:</u>** a weighted average

**<u>Multi-Head Attention:</u>**

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.

- 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

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

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

# 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

# Positional Encoding

```
 0 :   0 0 0 0      8 :   1 0 0 0
 1 :   0 0 0 1      9 :   1 0 0 1
 2 :   0 0 1 0     10 :   1 0 1 0
 3 :   0 0 1 1     11 :   1 0 1 1
 4 :   0 1 0 0     12 :   1 1 0 0
 5 :   0 1 0 1     13 :   1 1 0 1
 6 :   0 1 1 0     14 :   1 1 1 0
 7 :   0 1 1 1     15 :   1 1 1 1
```
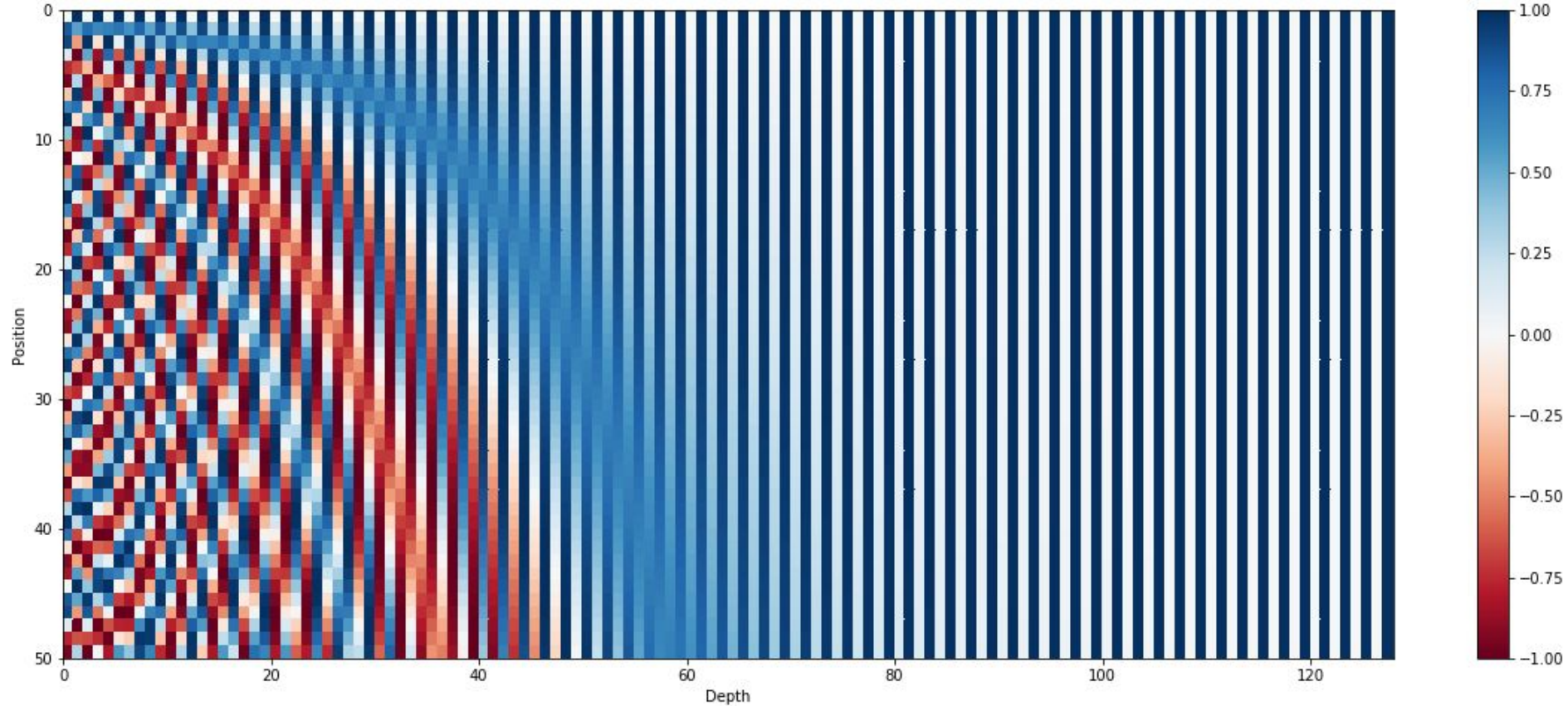
# Positional Encoding



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.*

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

Image source: https://kazemnejad.com/blog/transformer_architecture_positional_encoding/
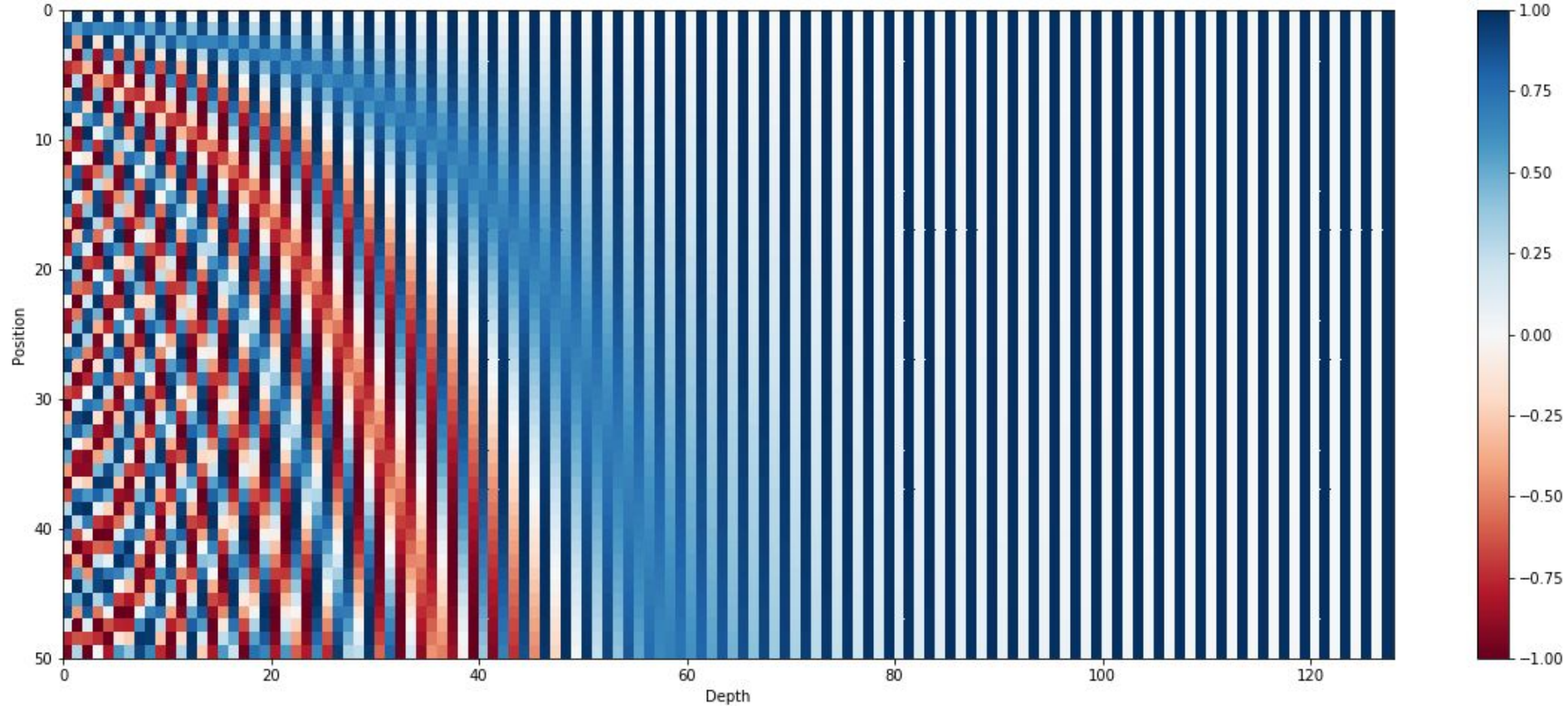
# Positional Encoding: why sin and cos?

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

$$\begin{bmatrix} u_1 & v_1 \\ u_2 & v_2 \end{bmatrix} \begin{bmatrix} \sin(\omega_k t) \\ \cos(\omega_k t) \end{bmatrix} = \begin{bmatrix} \sin(\omega_k t) \cos(\omega_k \phi) + \cos(\omega_k t) \sin(\omega_k \phi) \\ \cos(\omega_k t) \cos(\omega_k \phi) - \sin(\omega_k t) \sin(\omega_k \phi) \end{bmatrix}$$
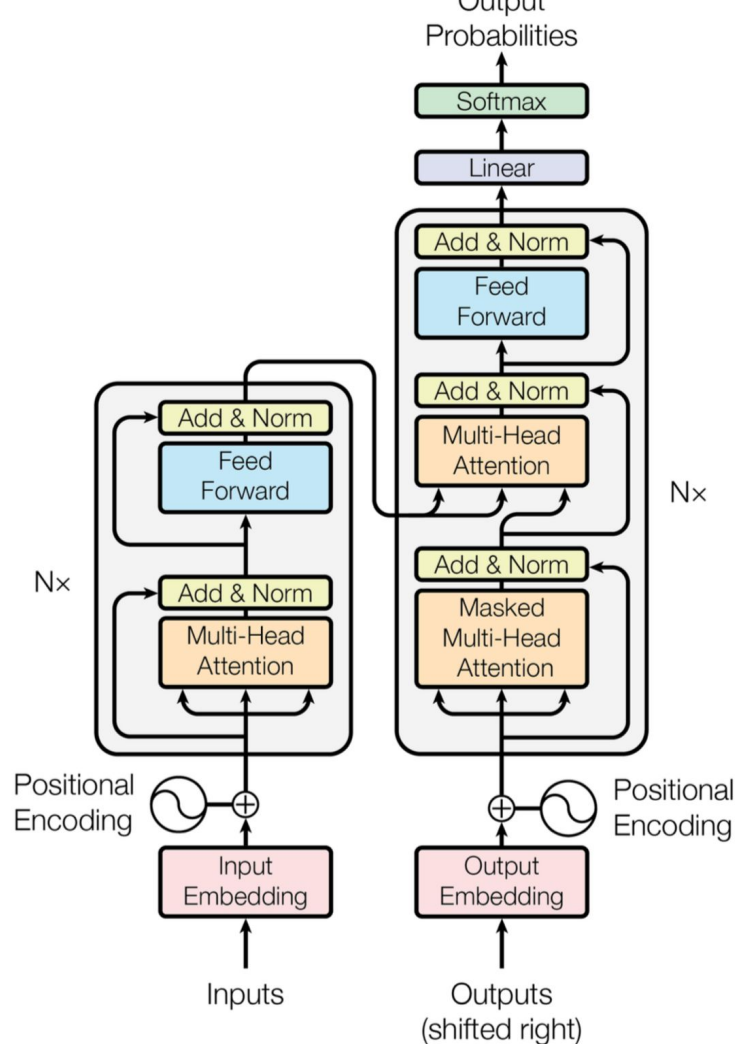
$$M_{\phi,k} = \begin{bmatrix} \cos(\omega_k \phi) & \sin(\omega_k \phi) \\ -\sin(\omega_k \phi) & \cos(\omega_k \phi) \end{bmatrix}$$
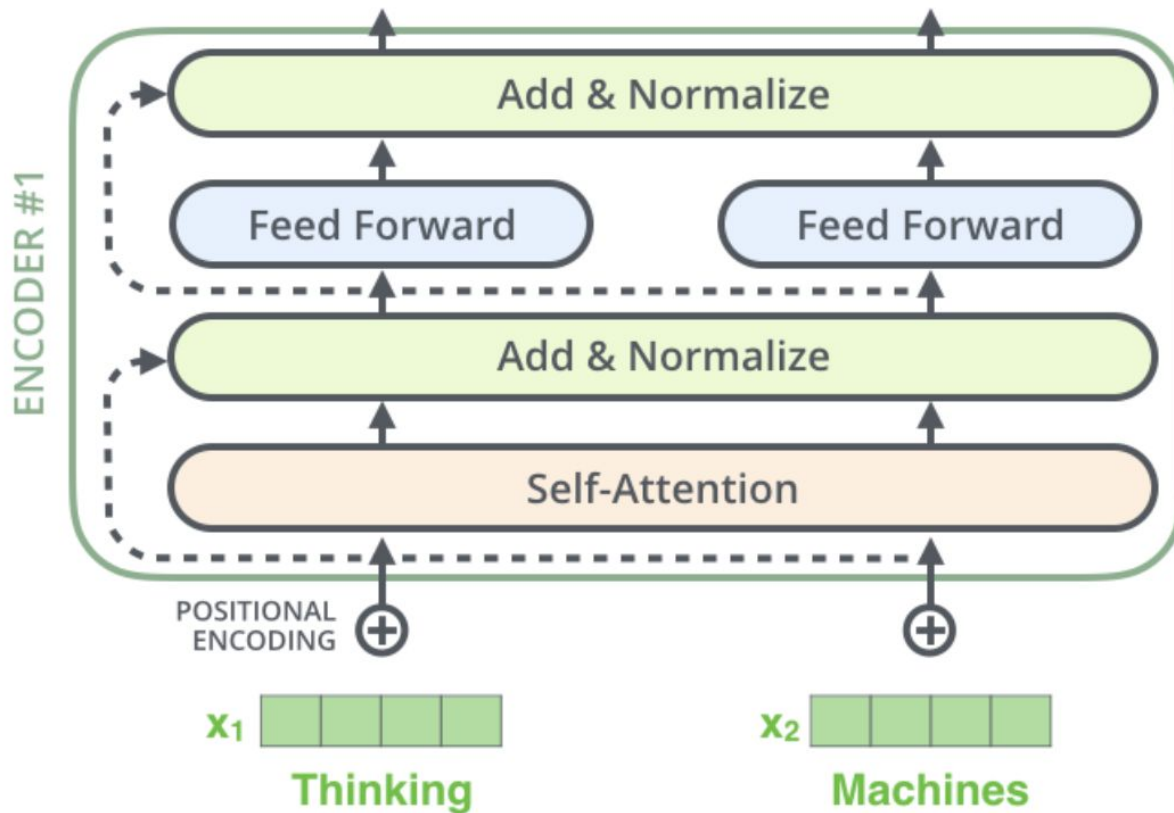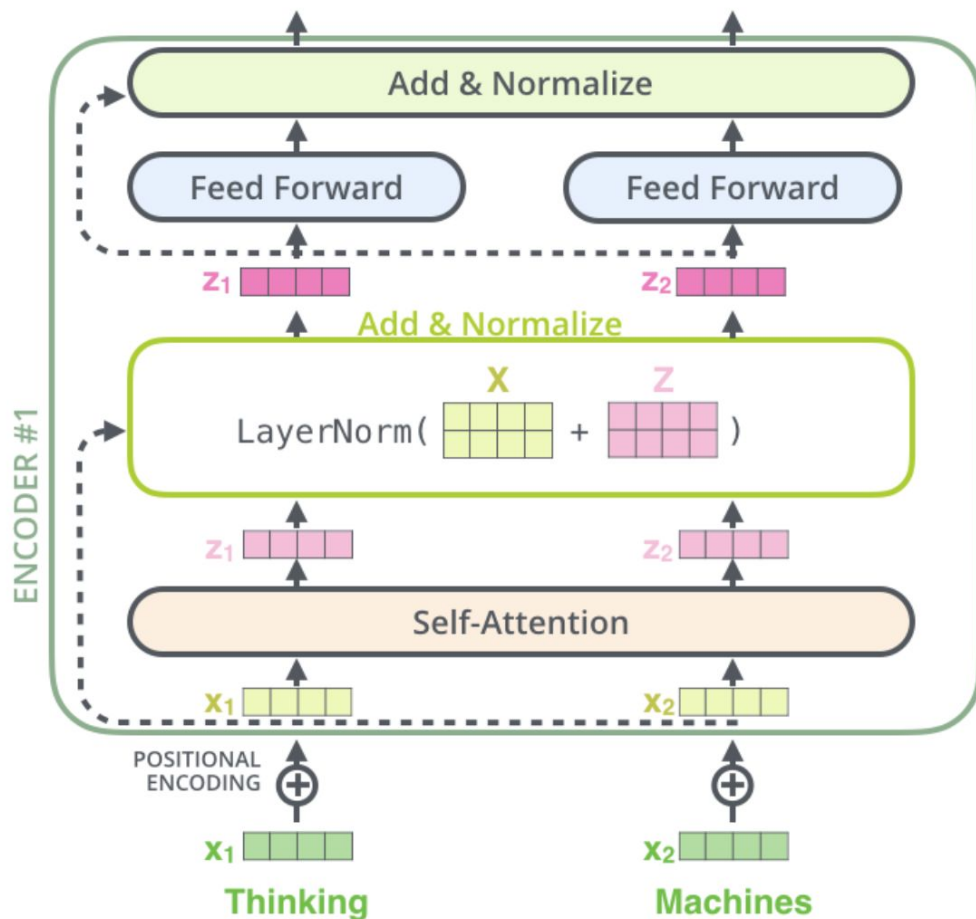
# Positional Encoding



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

# Layer Normalization

# The Transformer: recap

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

# Layer Normalization

# Layer Normalization

Like BatchNorm

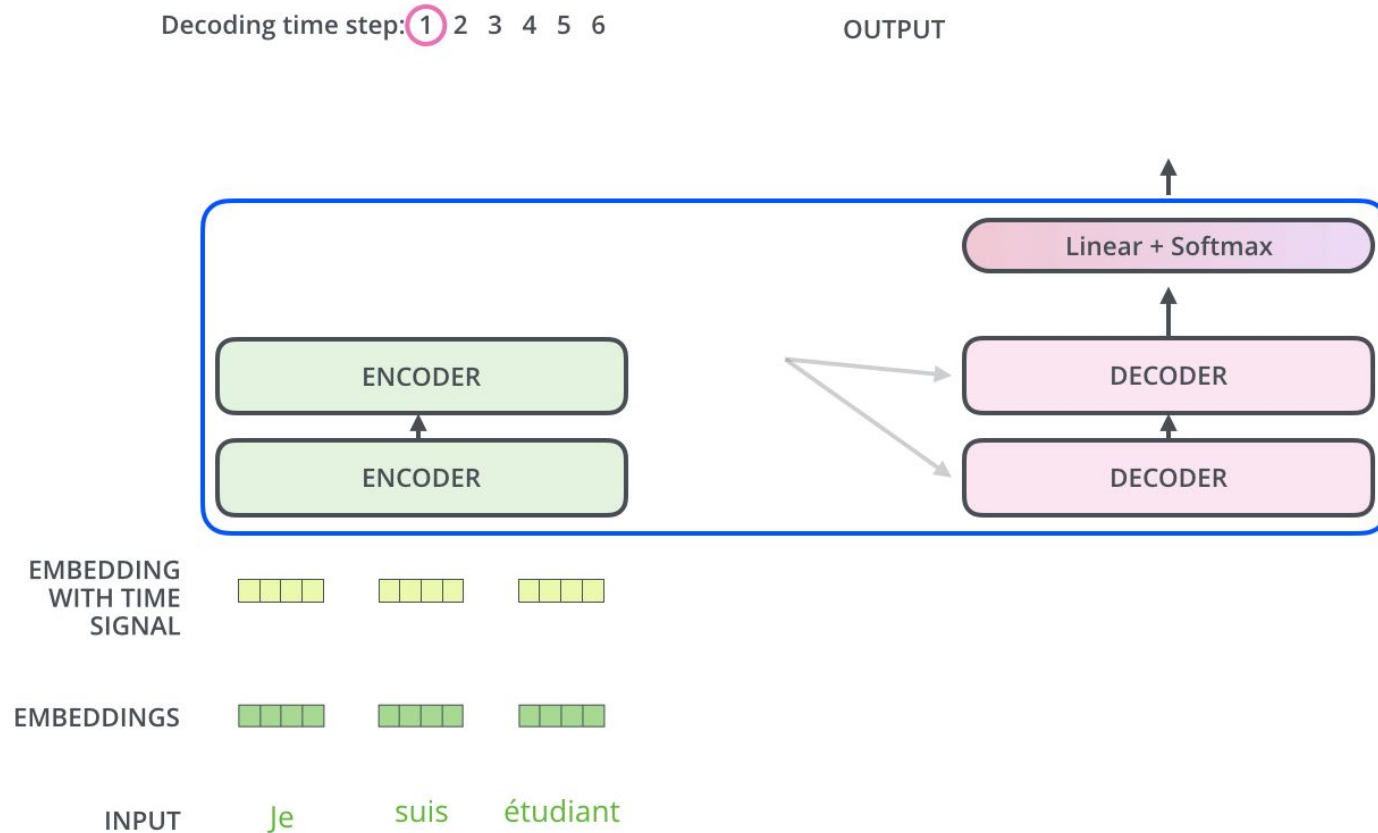but normalize along all features representing latent vector



More info:
Layer Normalization

# Layer Normalization



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

# The Decoder

# The Decoder Side



Decoding time step: (1) 2 3 4 5 6     OUTPUT

Linear + Softmax

ENCODER     DECODER

ENCODER     DECODER

EMBEDDING WITH TIME SIGNAL

EMBEDDINGS

INPUT    Je    suis    étudiant

# The Decoder Side



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

# The Decoder Side

Decoding time step: (1) 2 3 4 5 6      OUTPUT

Linear + Softmax

ENCODER

ENCODER

DECODER

EMBEDDING WITH TIME SIGNAL

EMBEDDINGS

INPUT    Je    suis    étudiant

MOVIECLIPS.com

The masked decoder input

# The Decoder Side



Image source: https://jalammar.github.io/illustrated-transformer/
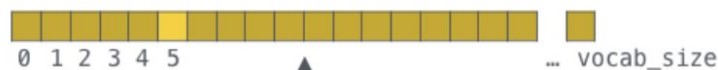
# Final Linear and Softmax Layer

Which word in our vocabulary is associated with this index?

am

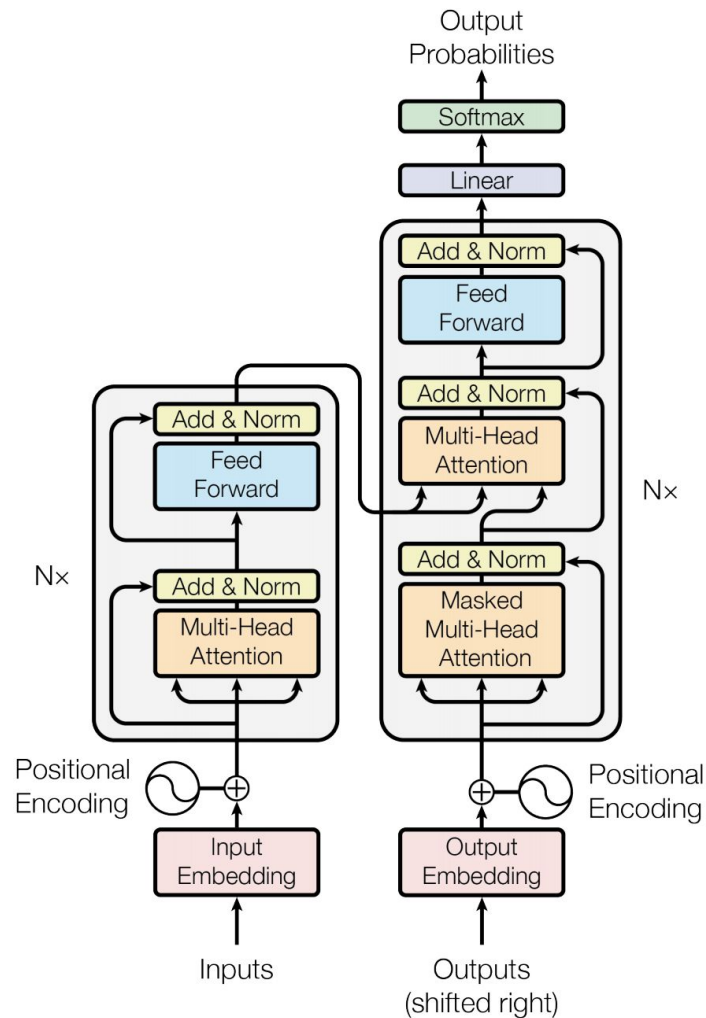Get the index of the cell with the highest value (`argmax`)

5

log_probs

0 1 2 3 4 5 ... vocab_size

Softmax

logits

0 1 2 3 4 5 ... vocab_size

Linear

Decoder stack output

# The Transformer



Output Probabilities

Softmax

Linear

Add & Norm

Feed Forward

Add & Norm

Multi-Head Attention

Add & Norm

Feed Forward

Add & Norm

Masked Multi-Head Attention

Add & Norm

Multi-Head Attention

N×

N×

Positional Encoding

Positional Encoding

Input Embedding

Output Embedding

Inputs

Outputs (shifted right)

- 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