# Lecture 4: Self-attention & Transformer
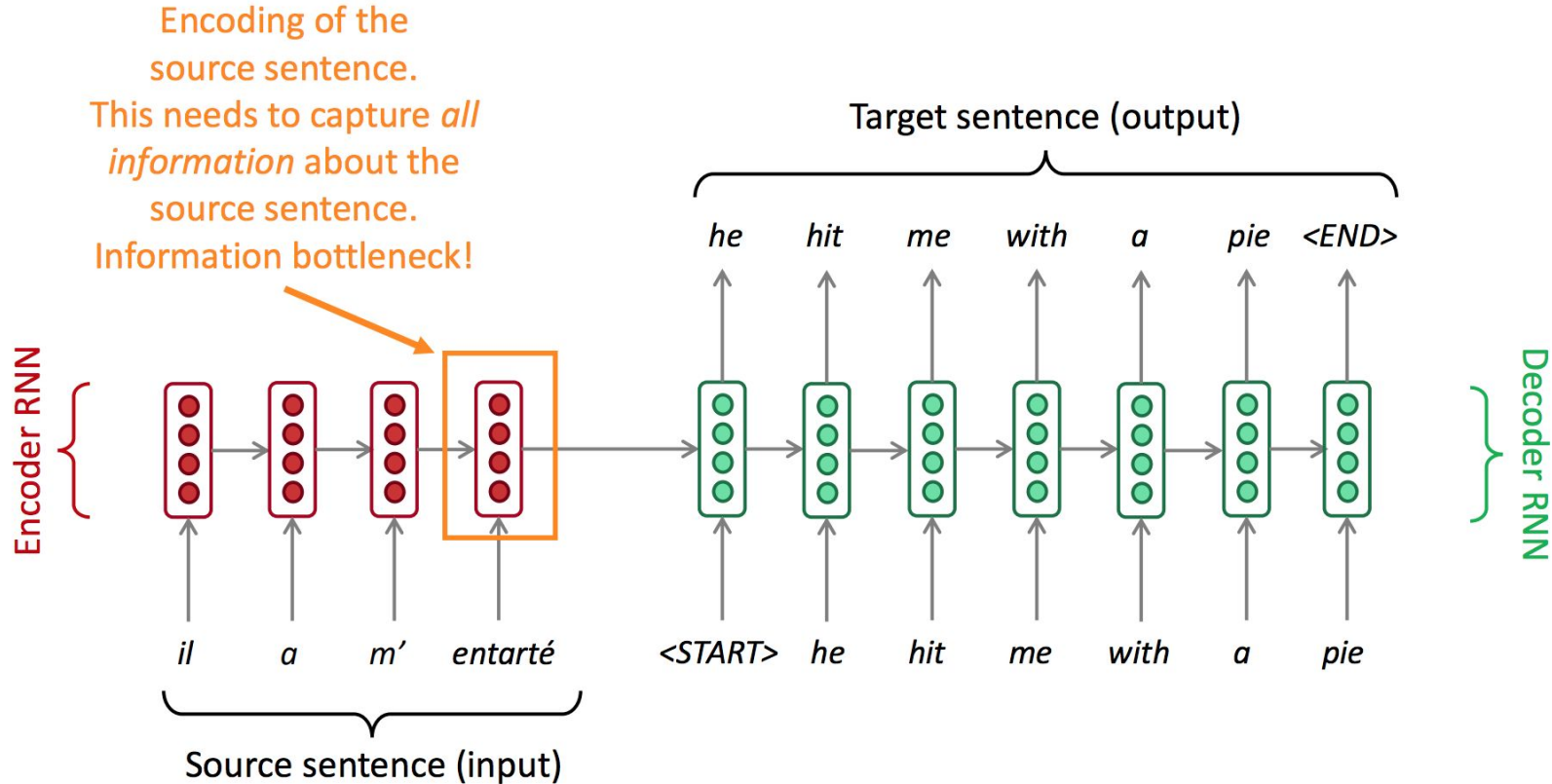
**Radoslav Neychev**

MIPT
27.09.2019, Moscow

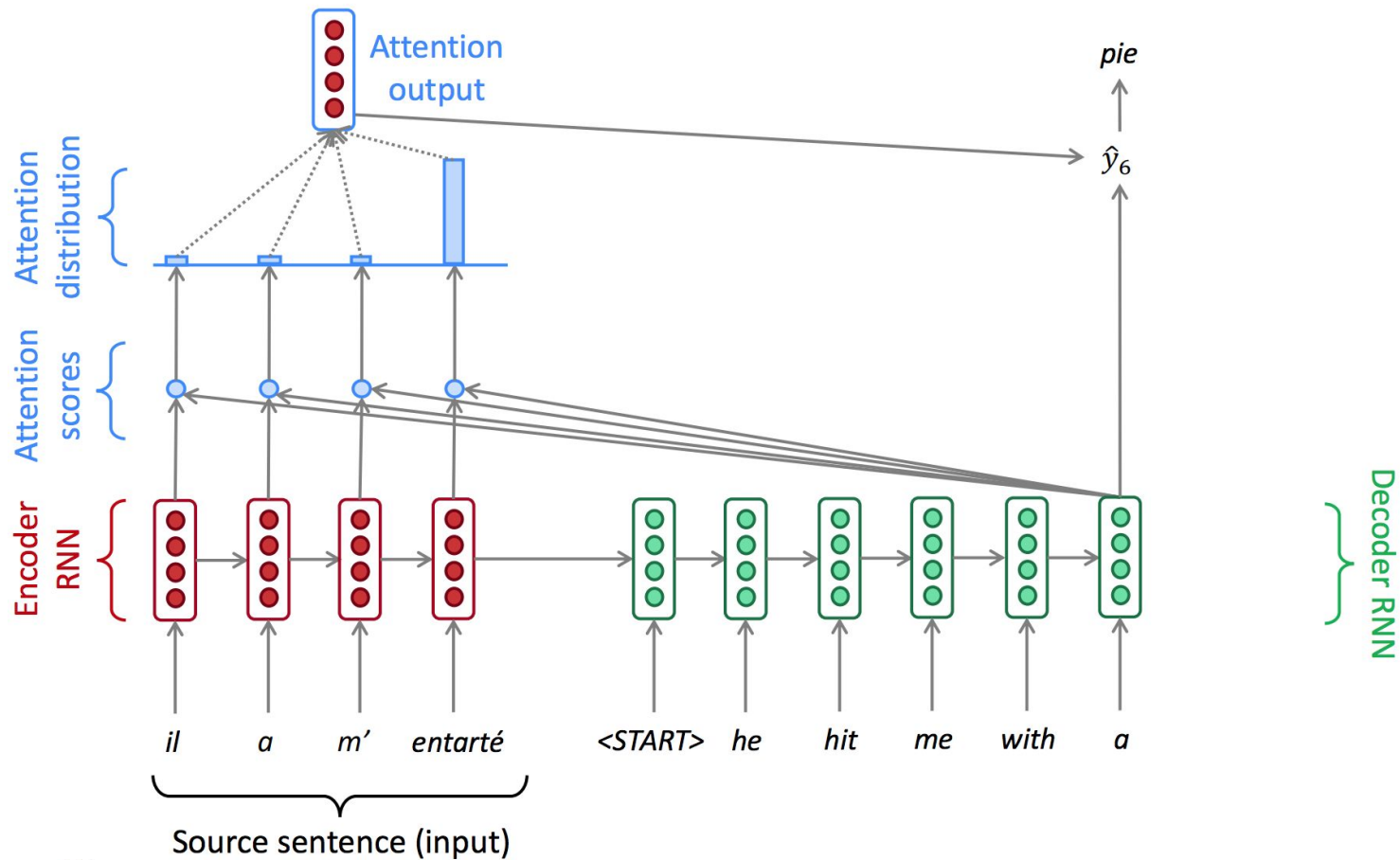# Outline

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

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

Encoding of the source sentence. This needs to capture *all information* about the source sentence. Information bottleneck!

Target sentence (output)

he    hit    me    with    a    pie    <END>

Encoder RNN

Decoder RNN

il    a    m'    entarté

<START>    he    hit    me    with    a    pie

Source sentence (input)

# recap: Attention in seq2seq

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

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

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

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

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

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

We use $\alpha^t$ 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
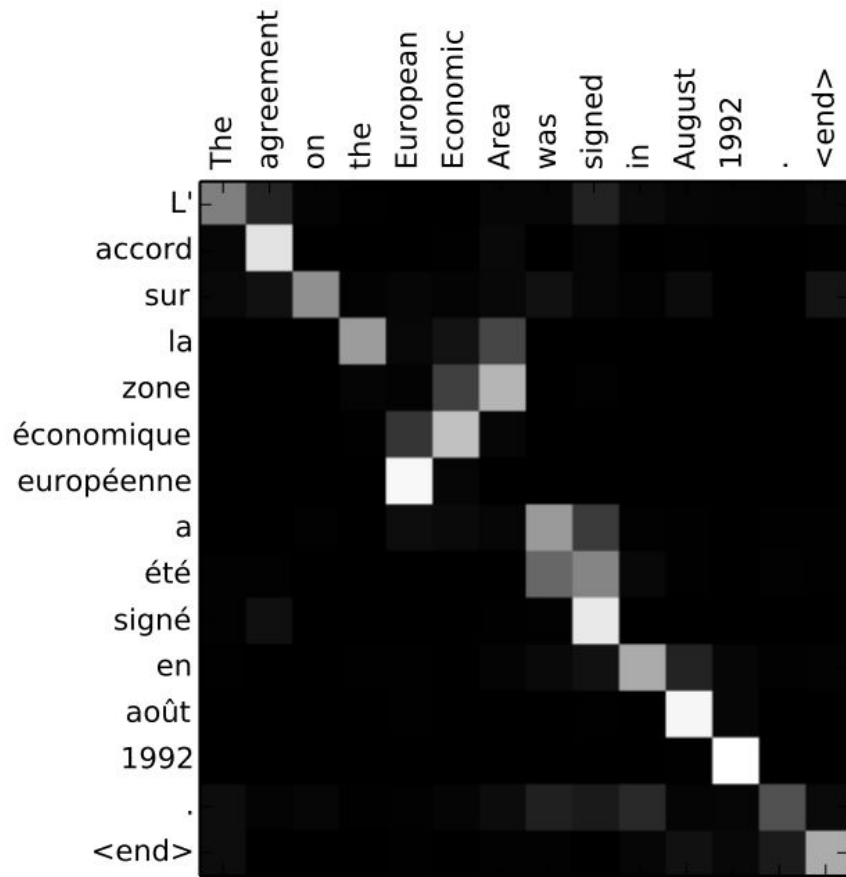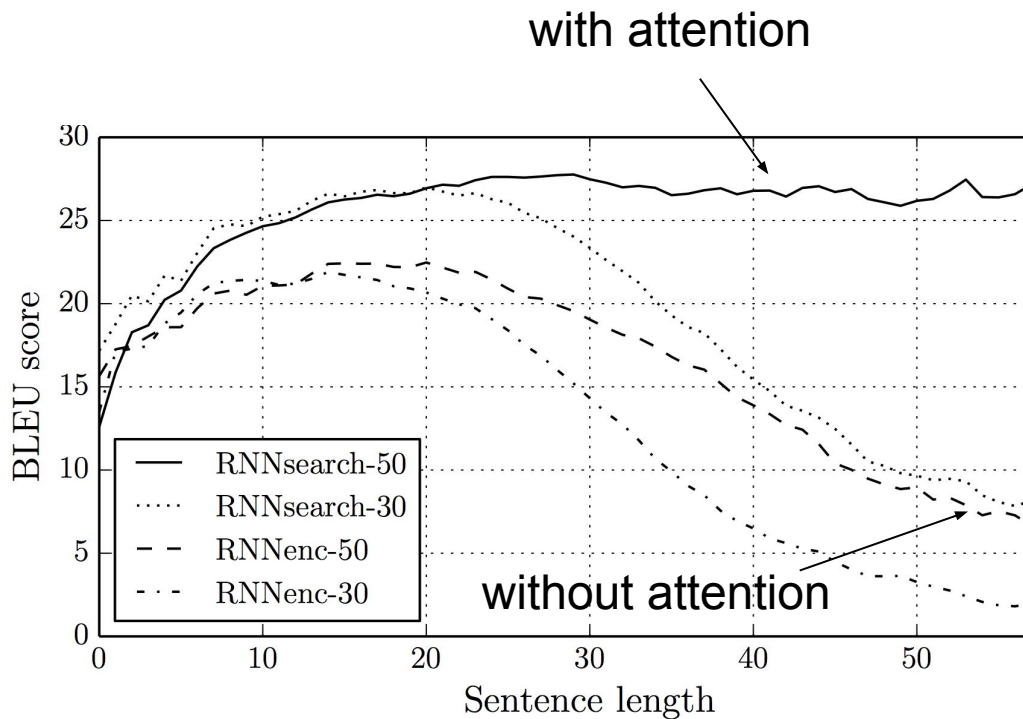
$$[a_t; s_t] \in \mathbb{R}^{2h}$$

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

# Attention advantages

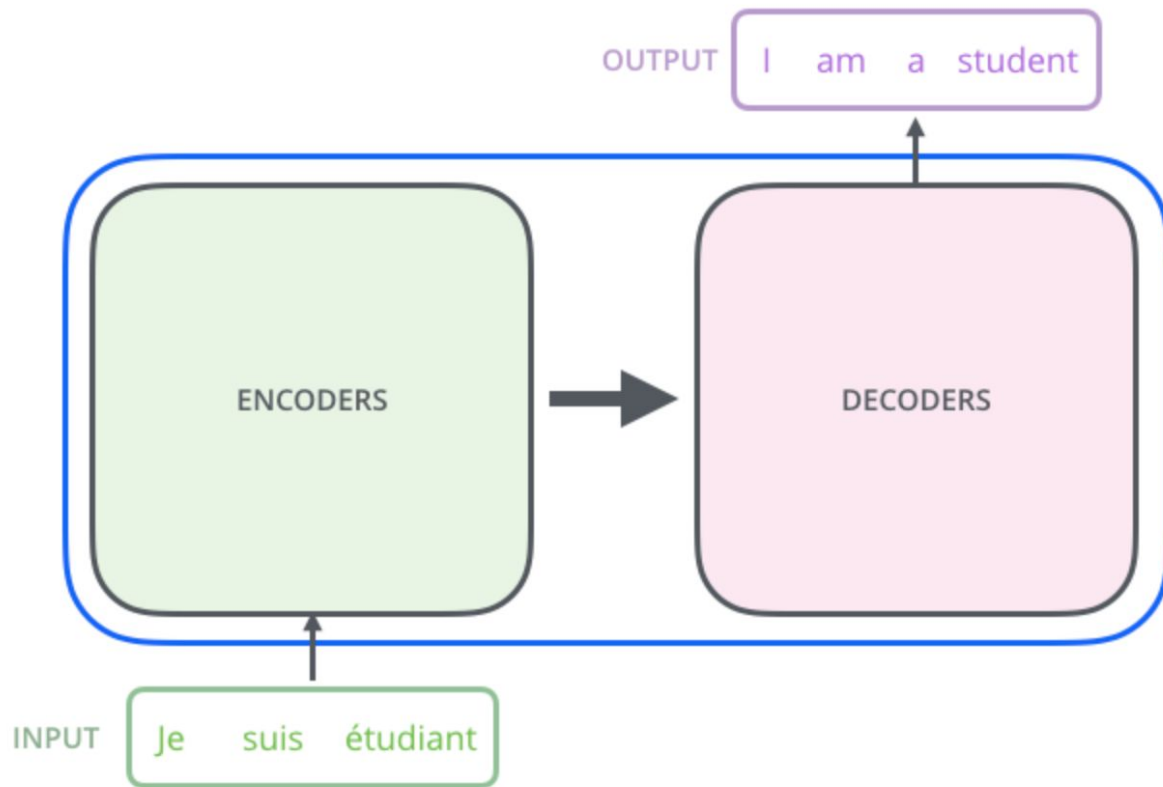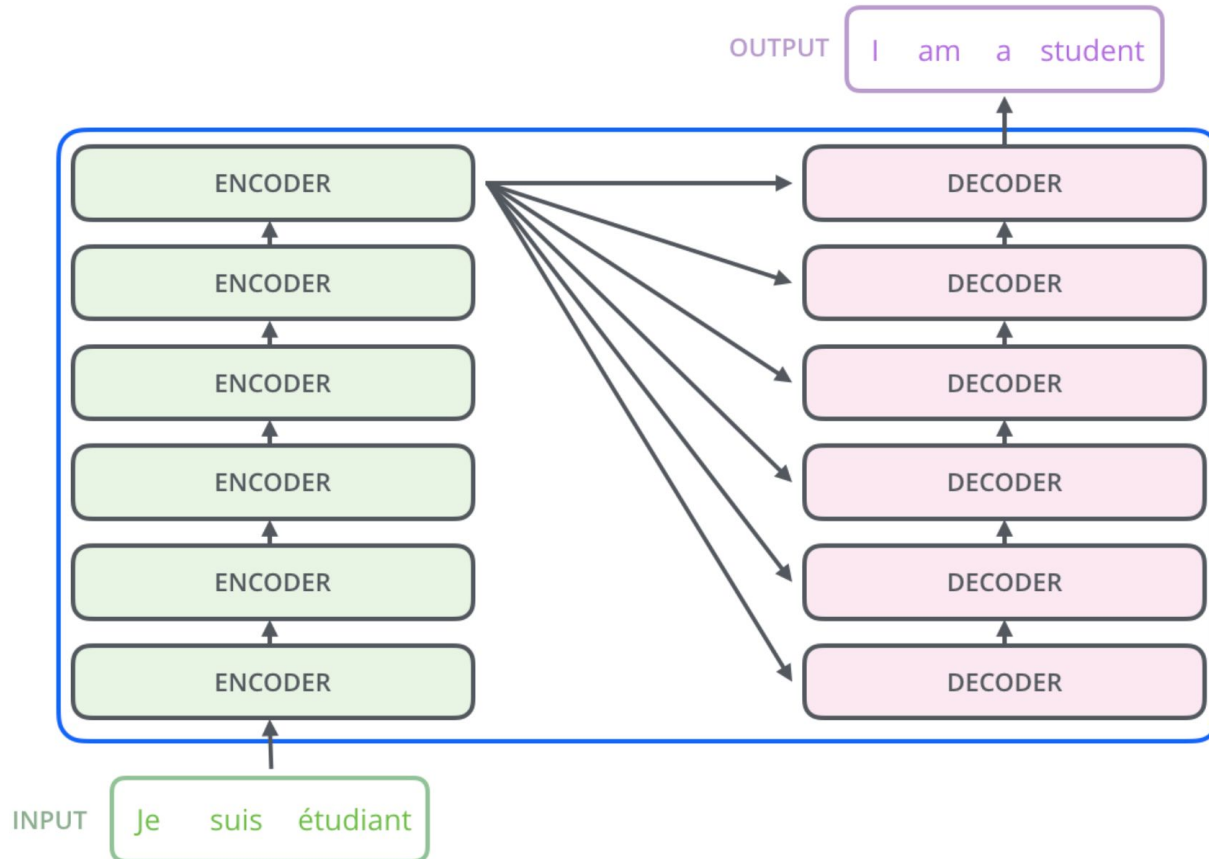- "Free" word alignment
- Better results on long sequences



with attention

without attention
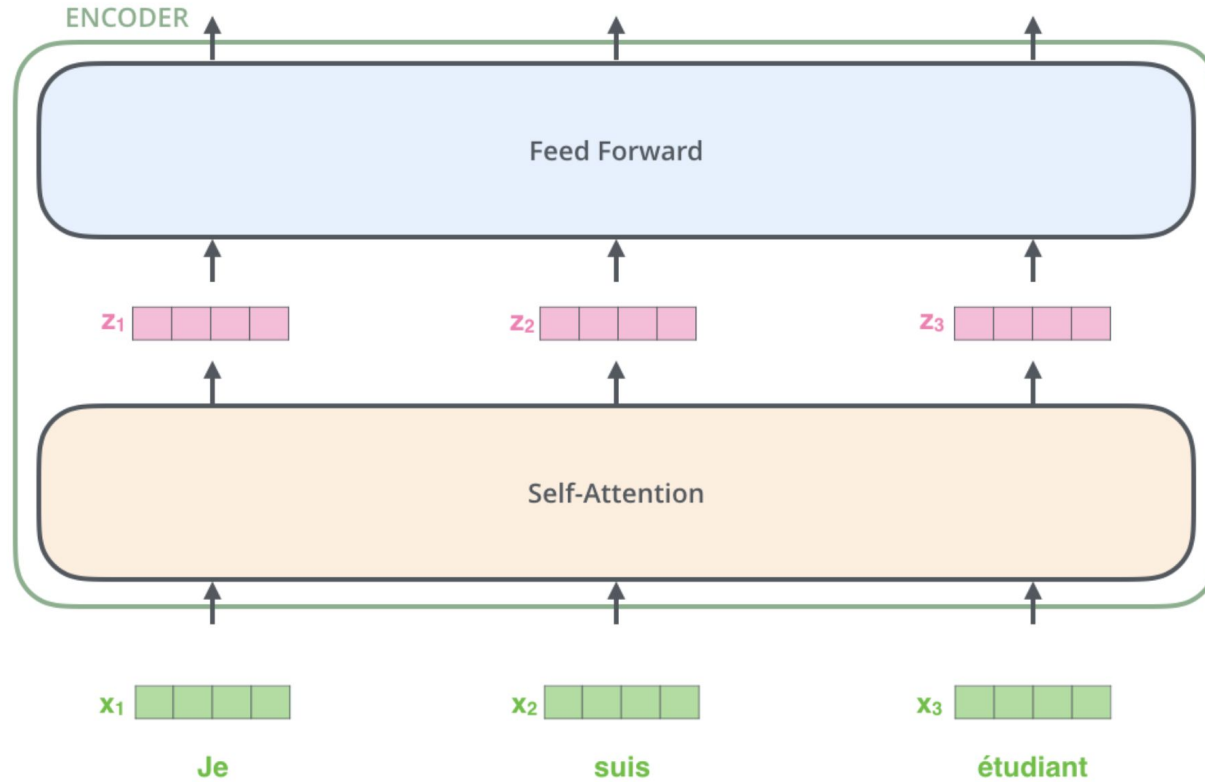
RNNsearch-50
RNNsearch-30
RNNenc-50
RNNenc-30

# The Transformer

# The Transformer

ENCODER

Feed Forward

$z_1$ | $z_2$ | $z_3$

Self-Attention

$x_1$ Je

$x_2$ suis

$x_3$ étudiant

# The Encoder Side

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

ENCODER

Can be parallelized →

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

# The Encoder Side

Can be parallelized →

Feed Forward Neural Network

Feed Forward Neural Network

$r_1$

$r_2$

$z_1$

$z_2$

Self-Attention

$x_1$ **Thinking**

$x_2$ **Machines**

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

15

# 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

"The animal didn't cross the street because it was too tired"
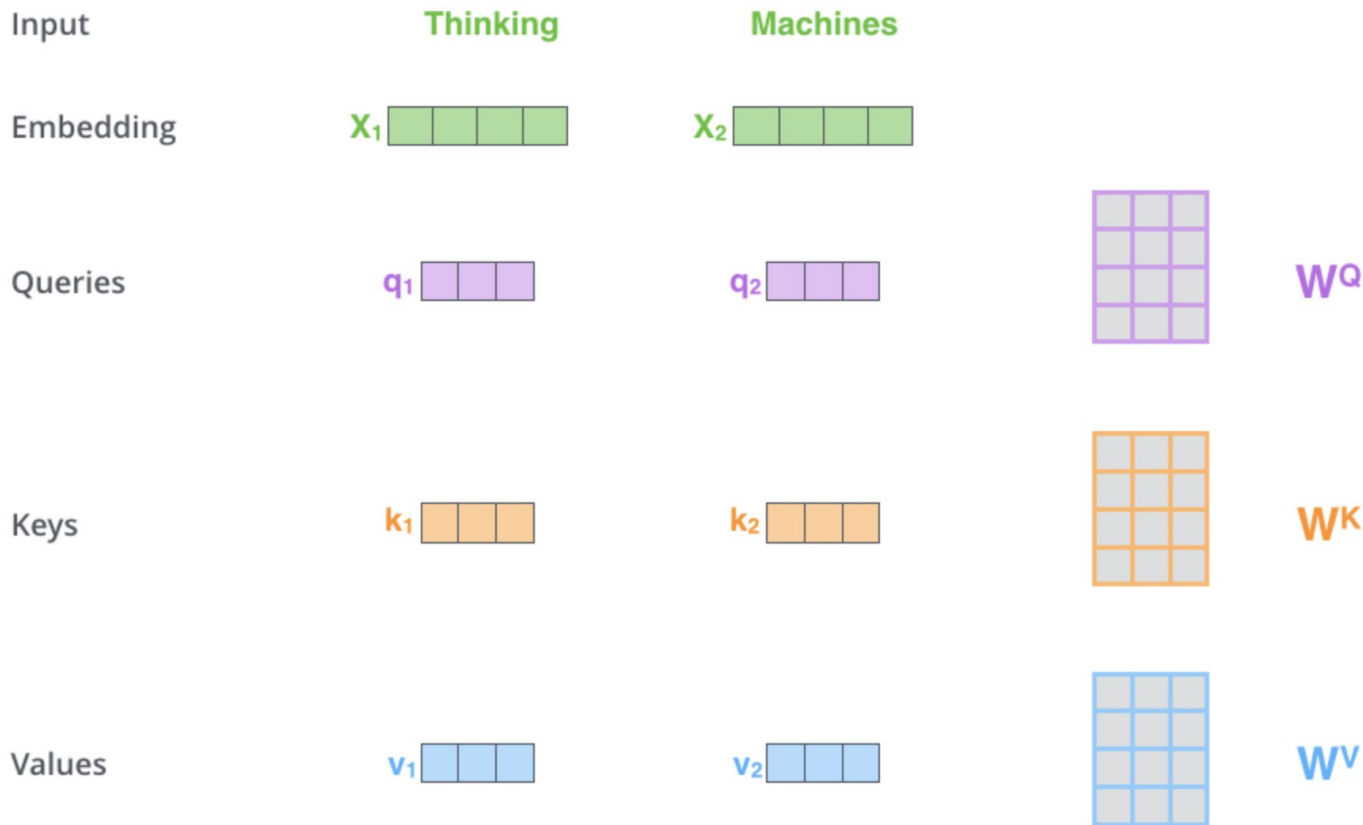
- What does "it" in this sentence refer to?

"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

# Self-Attention: detailed explanation

# Self-Attention: detailed explanation

## STEP 1:

create 3 vectors
(Query, Key, Value)

from each of the encoder's
input vectors

| Input | | Thinking | Machines | |
|---|---|---|---|---|
| Embedding | | $X_1$ | $X_2$ | |
| Queries | | $q_1$ | $q_2$ | $W^Q$ |
| Keys | | $k_1$ | $k_2$ | $W^K$ |
| Values | | $v_1$ | $v_2$ | $W^V$ |

# Self-Attention: detailed explanation

What are the "query", "key", and "value" vectors?

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)
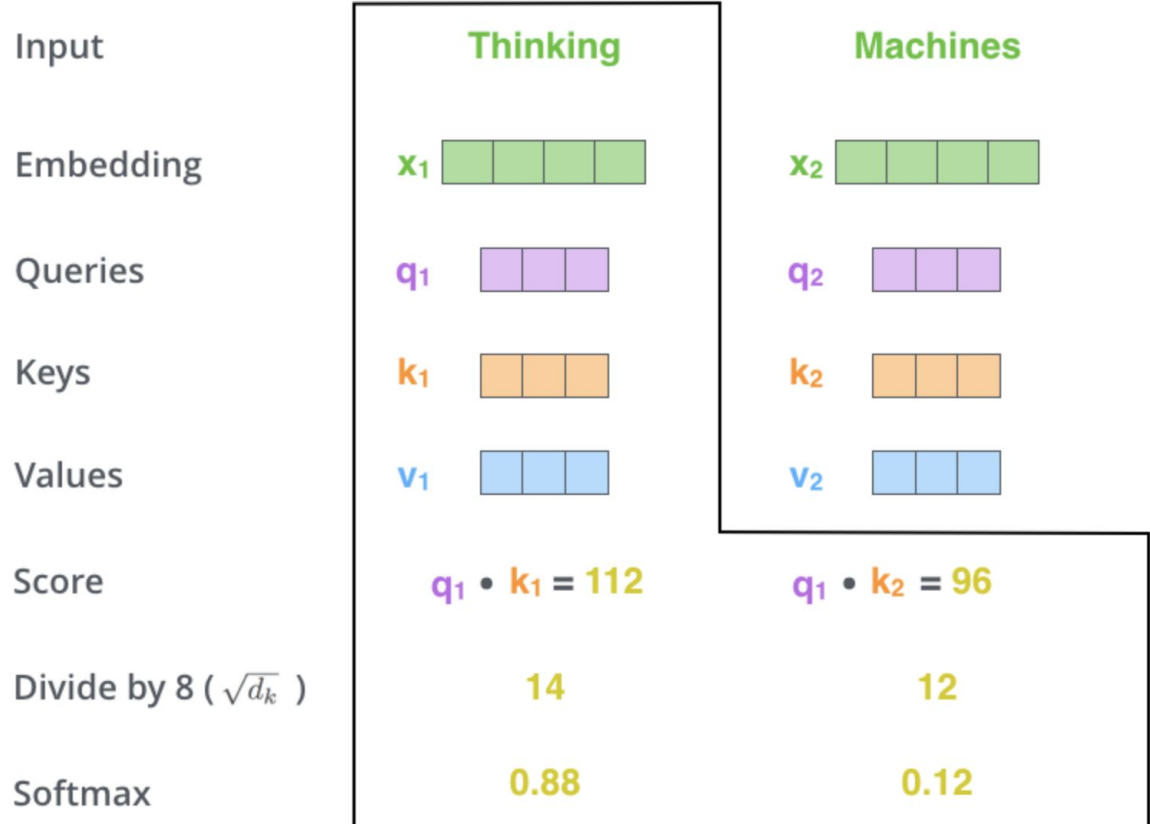
# Self-Attention: detailed explanation

**STEP 3:**

divide the scores by 8

(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 |

# 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

**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 (**Wk, Wq, Wv**)

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

$$= Z$$

# Multi-Head Attention

X

Thinking
Machines

Calculating attention separately in
eight different attention heads

ATTENTION
HEAD #0

ATTENTION
HEAD #1

...

ATTENTION
HEAD #7

$Z_0$

$Z_1$

$Z_7$

# 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
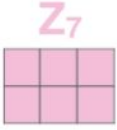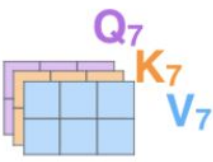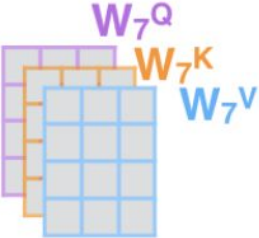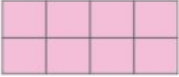
Thinking Machines

$X$

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

$Q_0$
$K_0$
$V_0$

$Z_0$

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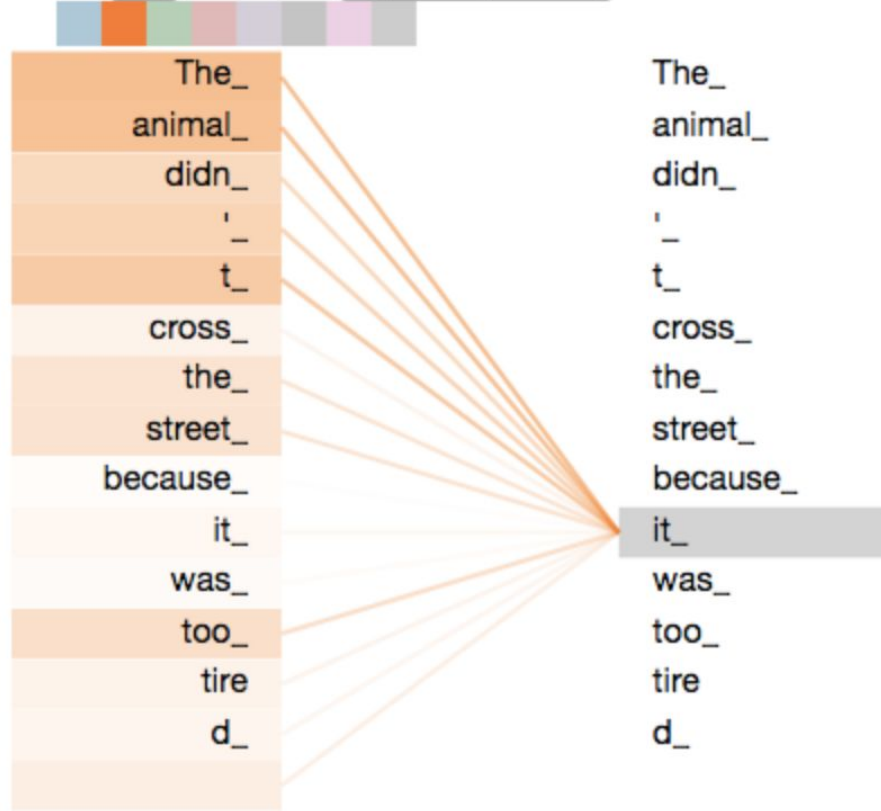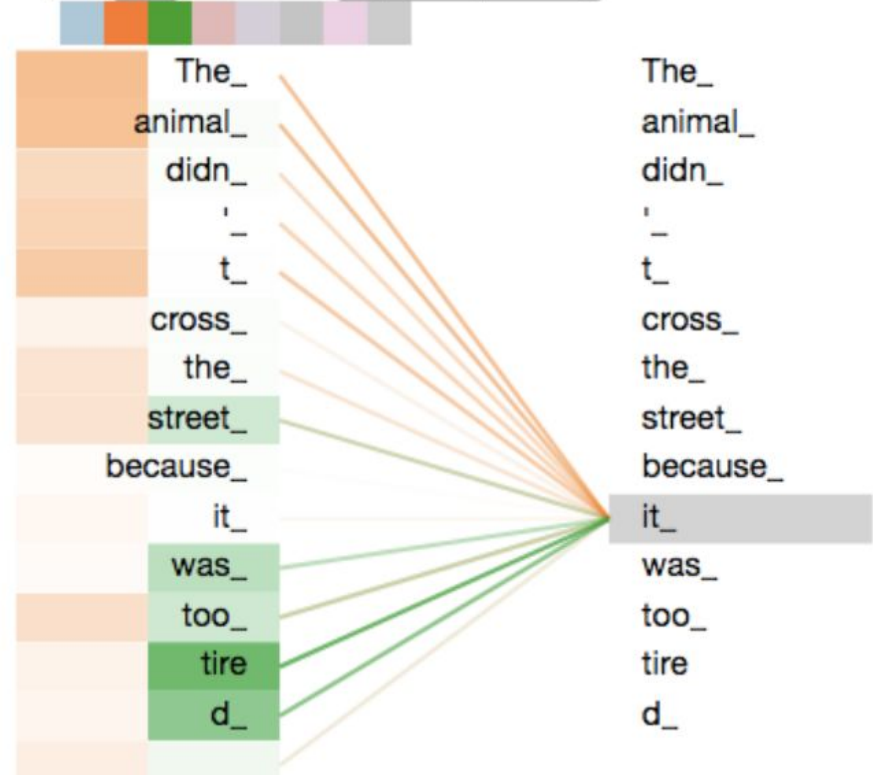
* In all encoders other than #0, we don't need embedding. We start directly with the output of the encoder right below this one

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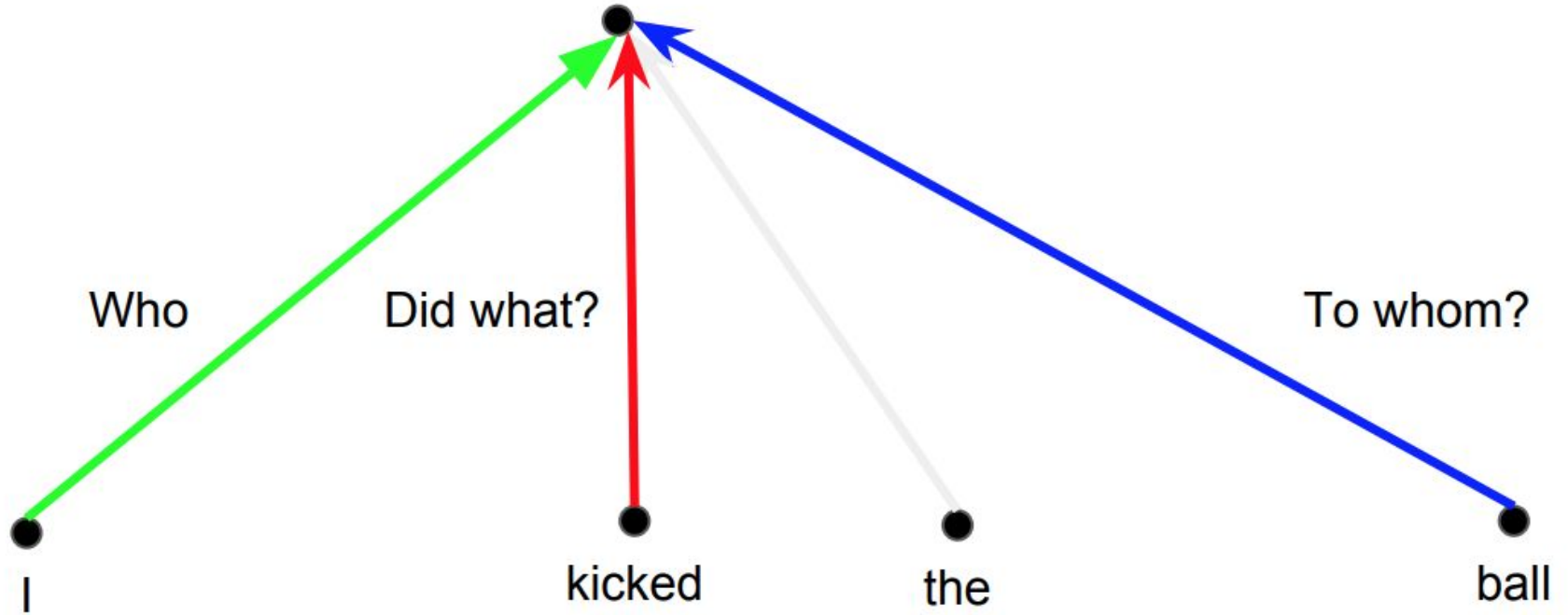$W_1^V$

$Q_1$
$K_1$
$V_1$

$Z_1$

$Z$

$R$

...

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

...

$Q_7$
$K_7$
$V_7$

...

$Z_7$

# Multi-Head Attention

Who

Did what?

To whom?

I

kicked

the

ball

# 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.

- 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

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

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

- *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 2pi to 2pi * 10 000

# Positional Encoding

| POSITIONAL ENCODING | 0 | 0 | 1 | 1 |
|---|---|---|---|---|

| | 0.84 | 0.0001 | 0.54 | 1 |
|---|---|---|---|---|

| | 0.91 | 0.0002 | -0.42 | 1 |
|---|---|---|---|---|

+      +      +

EMBEDDINGS $x_1$ [ ] [ ] [ ] [ ]    $x_2$ [ ] [ ] [ ] [ ]    $x_3$ [ ] [ ] [ ] [ ]

INPUT     Je     suis     étudiant

# Layer Normalization

The Transformer: recap

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

Nx

Nx

Positional Encoding

Positional Encoding

Input Embedding

Output Embedding

Inputs

Outputs (shifted right)

50

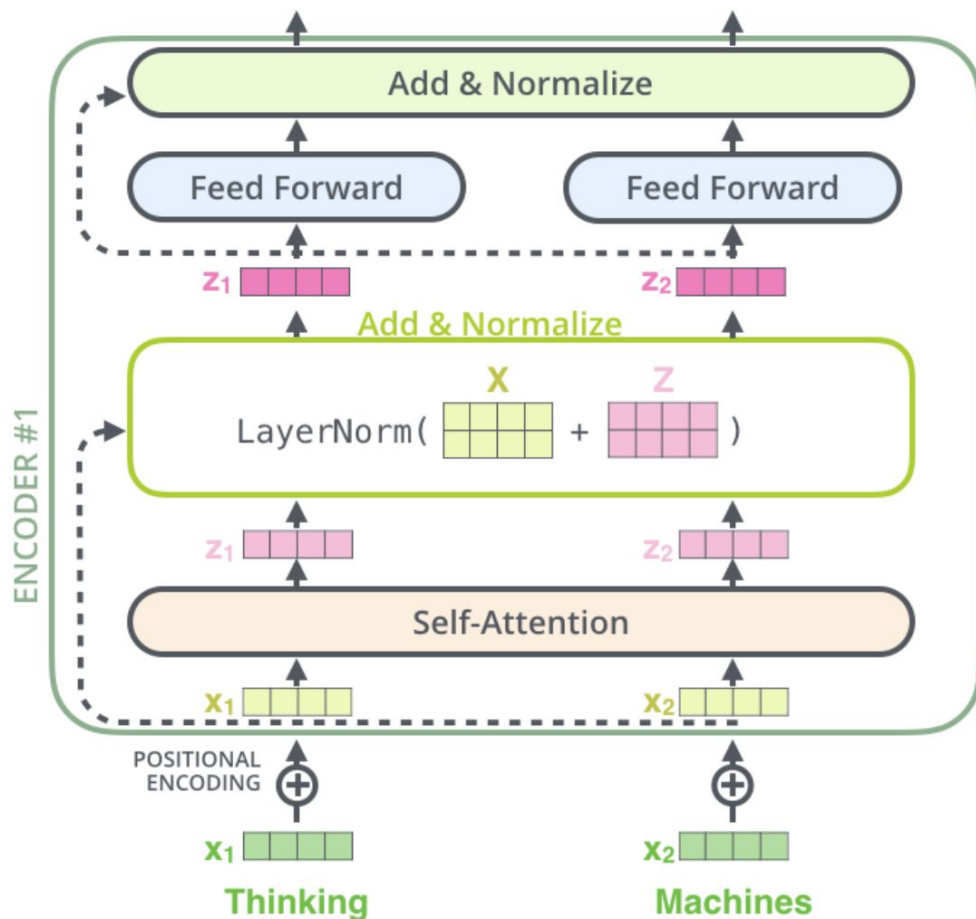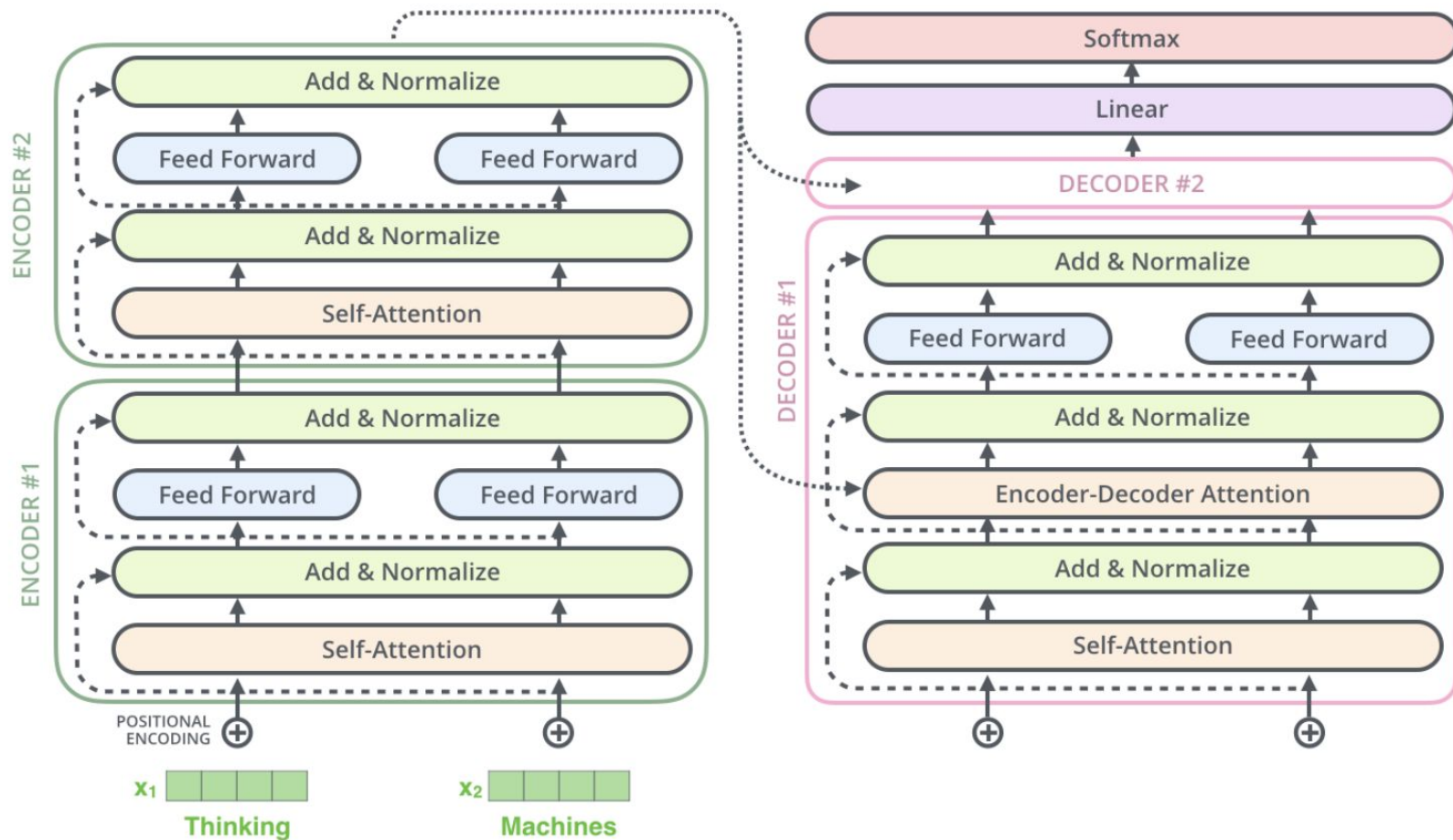# Layer Normalization

# Layer Normalization

# Layer Normalization

# The Decoder

# The Decoder Side



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

Linear + Softmax

ENCODER

ENCODER

DECODER

DECODER

EMBEDDING WITH TIME SIGNAL

EMBEDDINGS

INPUT    Je    suis    étudiant

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

Linear + Softmax

ENCODER                    DECODER

ENCODER                    DECODER

EMBEDDING
WITH TIME
SIGNAL

*Here comes
the mask*

EMBEDDINGS

INPUT        Je        suis      étudiant

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

OUTPUT

Linear + Softmax

ENCODER

ENCODER

DECODER

EMBEDDING WITH TIME SIGNAL

EMBEDDINGS

INPUT        Je        suis      étudiant

The masked decoder input

# The Decoder Side



Decoding time step: 1 ② 3 4 5 6    OUTPUT    I

$K_{encdec}$    $V_{encdec}$    Linear + Softmax

ENCODERS    DECODERS

EMBEDDING WITH TIME SIGNAL

EMBEDDINGS

INPUT    Je    suis    étudiant    PREVIOUS OUTPUTS    I

# 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

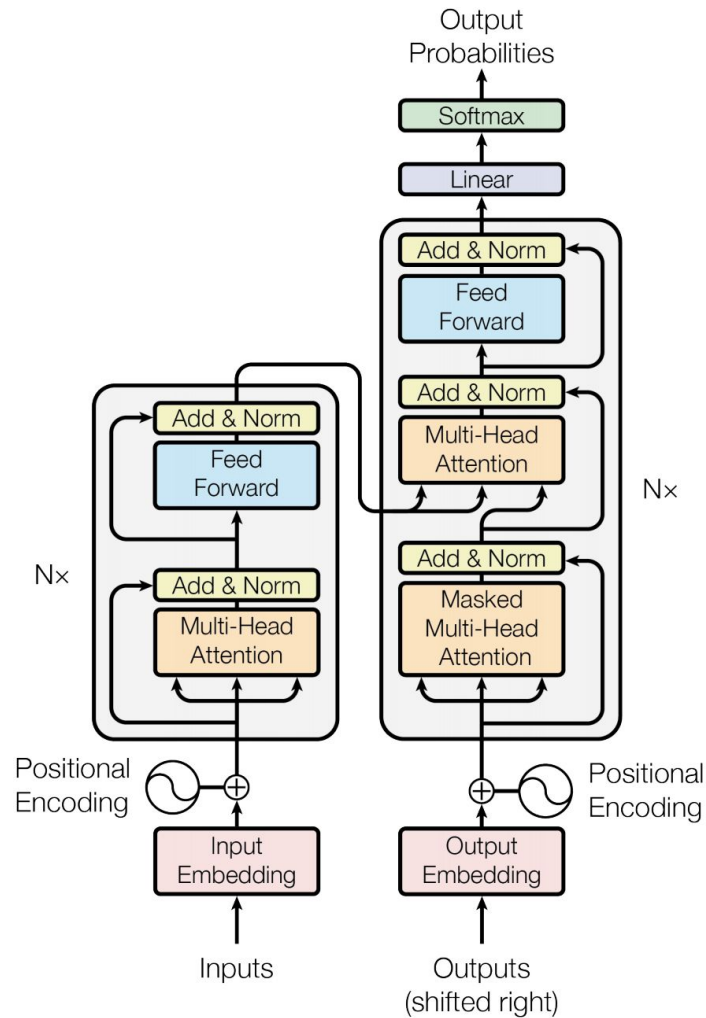0 1 2 3 4 5 … vocab_size

**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

Add & Norm

Multi-Head
Attention

Masked
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