CSCI 1470/2470
Spring 2024
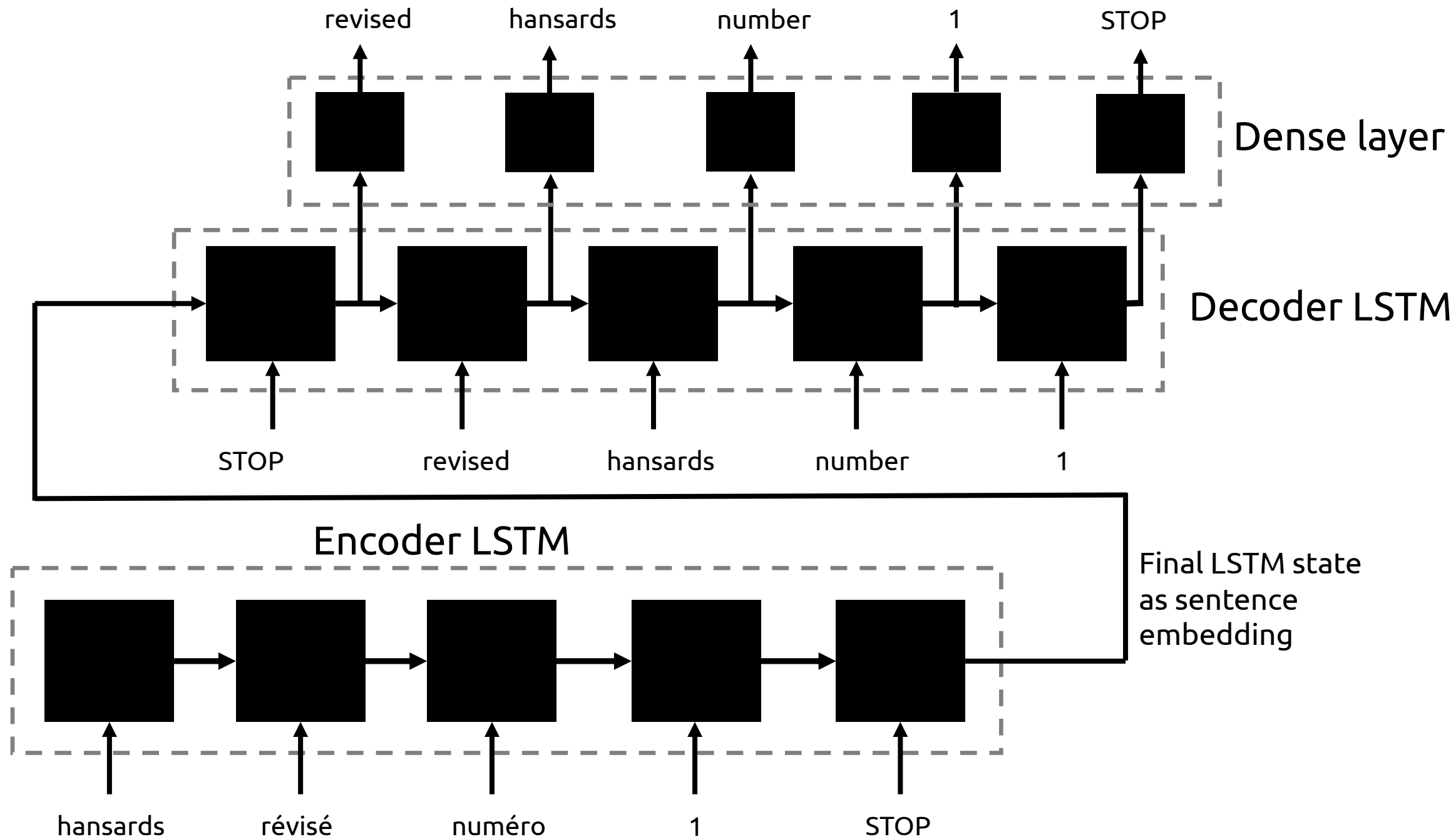
Ritambhara Singh

March 11, 2024

Monday

Transformers

Deep Learning
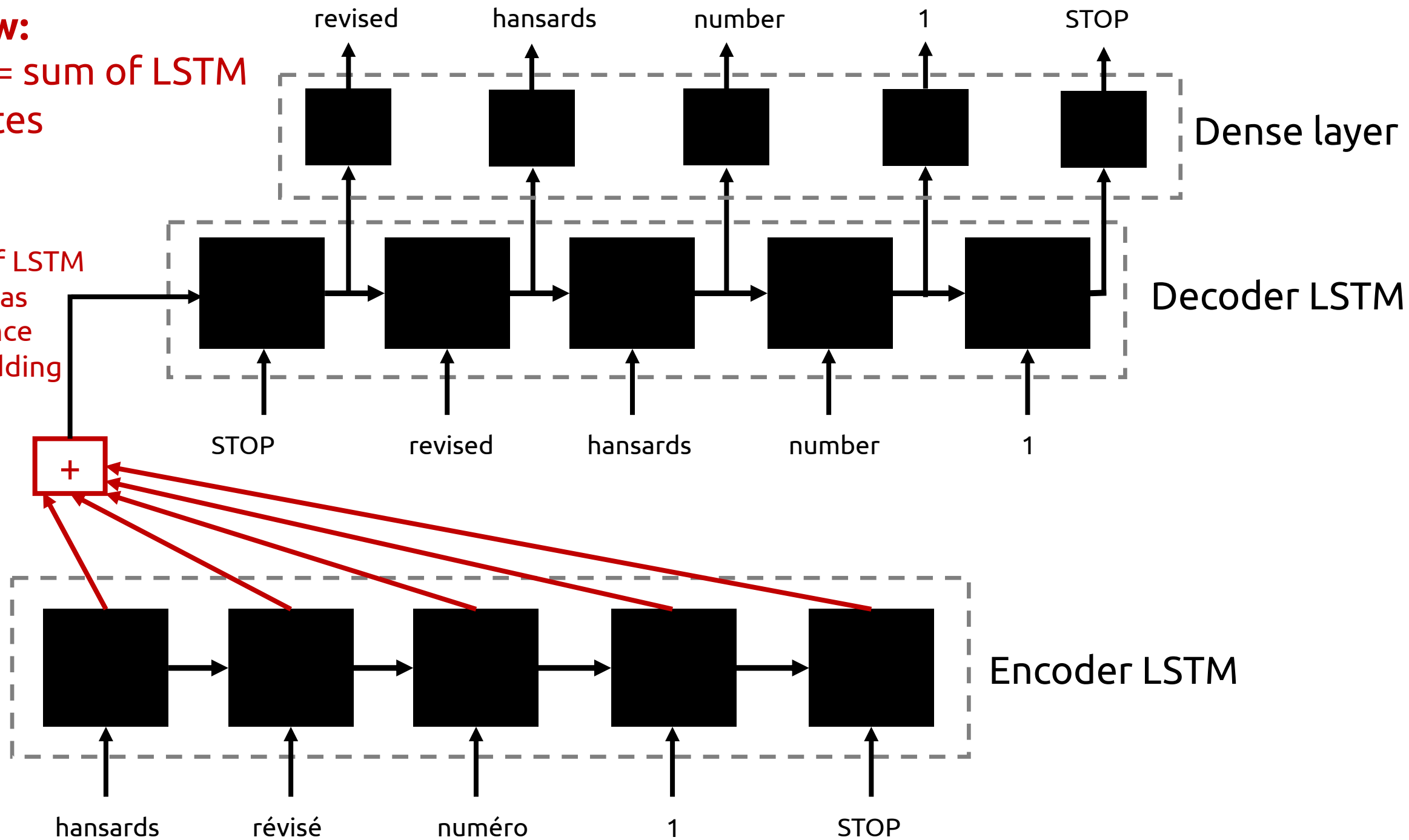
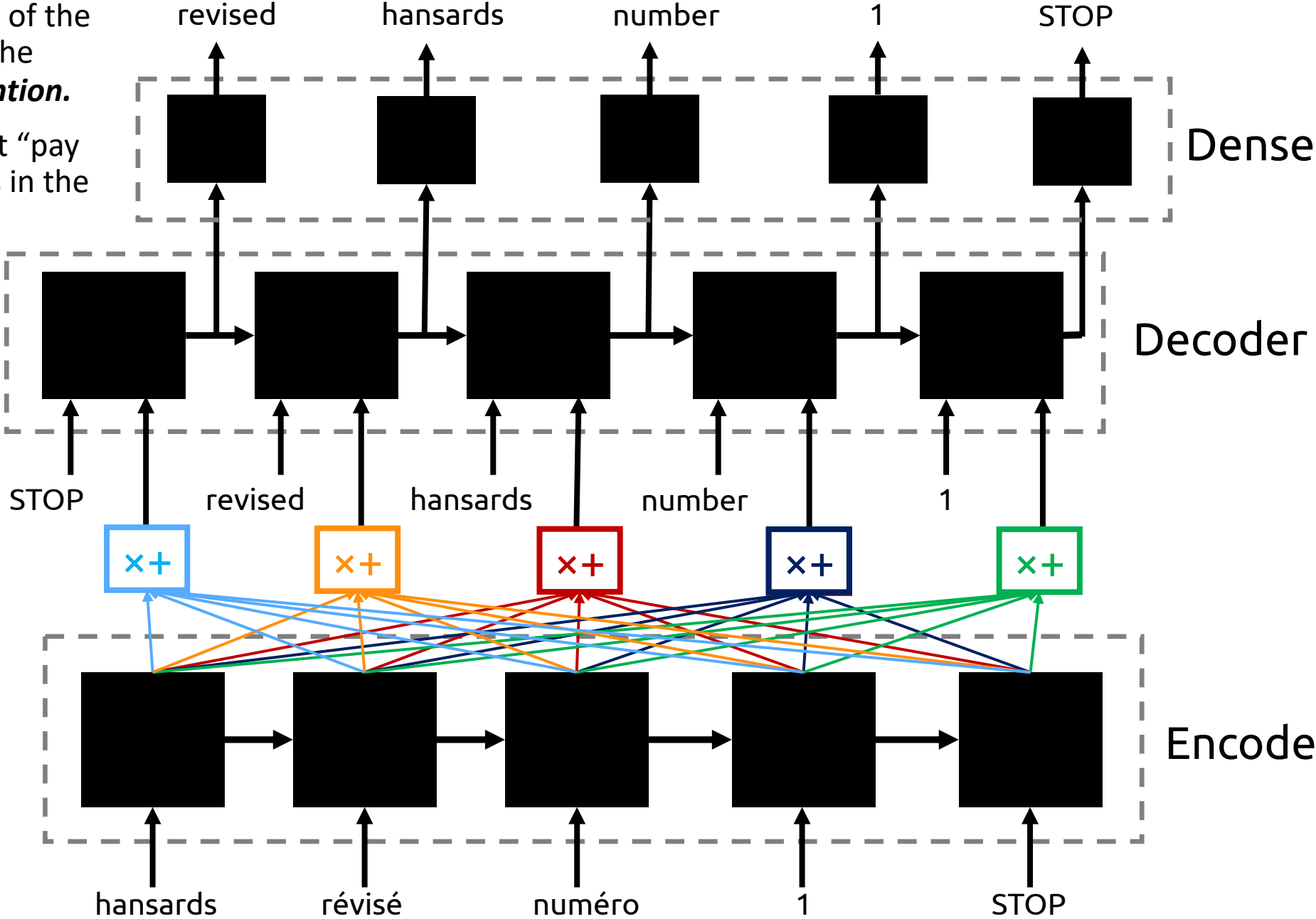ChatGPT prompt "minimalist landscape painting of a deep underwater scene with a blue tang fish in the bottom right corner"

Dense layer

revised · hansards · number · 1 · STOP

Decoder LSTM

STOP · revised · hansards · number · 1

Encoder LSTM

Final LSTM state as sentence embedding

hansards · révisé · numéro · 1 · STOP

**New:**
$E_s$ = sum of LSTM states

revised     hansards     number     1     STOP

Dense layer

Sum of LSTM states as sentence embedding

Decoder LSTM

STOP     revised     hansards     number     1

+

Encoder LSTM

hansards     révisé     numéro     1     STOP

This idea of passing each cell of the decoder a weighted sum of the encoder states is called **attention.**

Different words in the output "pay attention" to different words in the input

# Review: "Attention Is All You Need"

A 2017 paper that introduced the **Transformer** model for machine translation

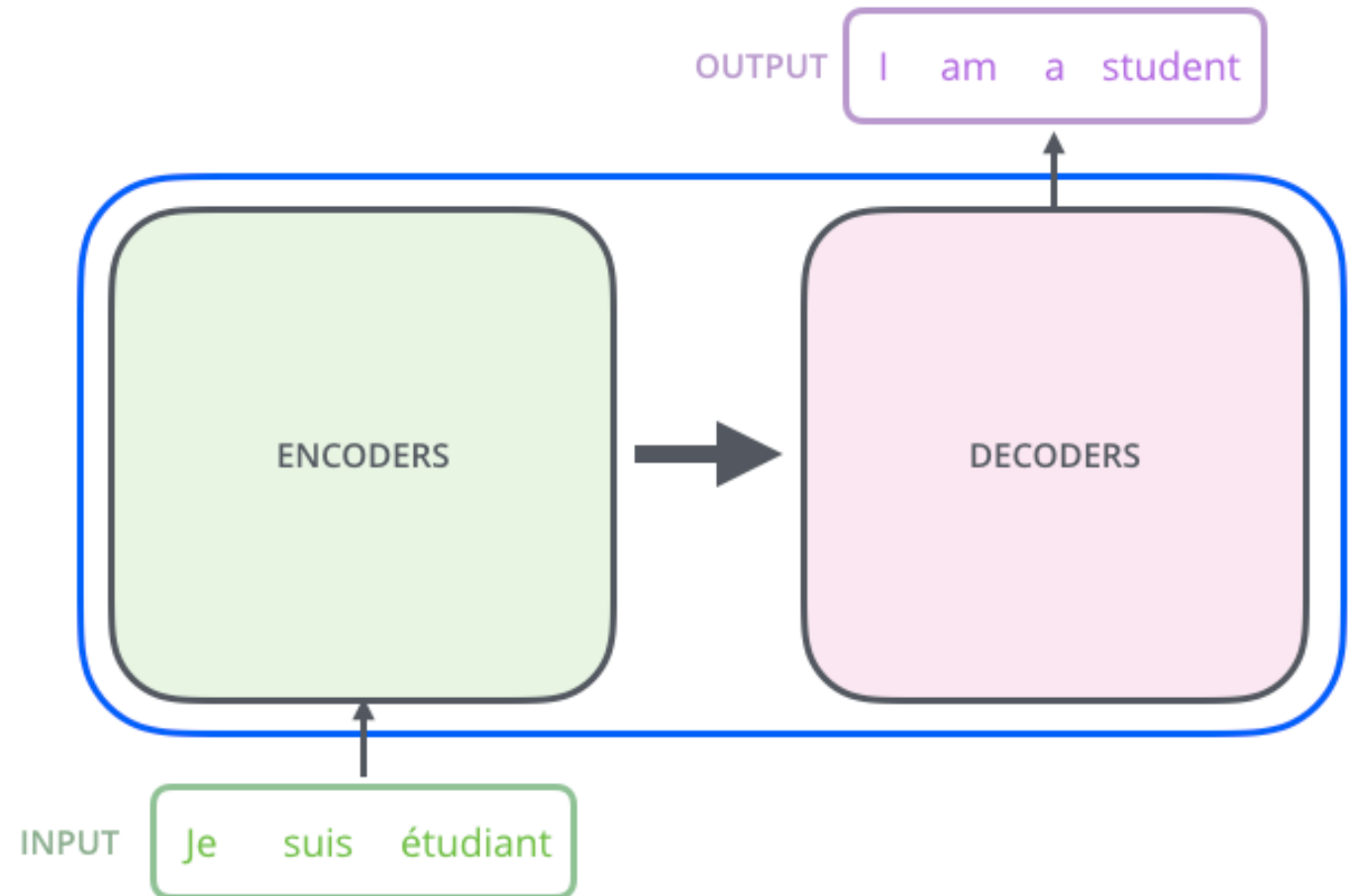- Has no recurrent networks!
- **Only** uses attention

**Motivation:**

- RNN training is hard to parallelize since the previous word must be processed before next word
  - Transformers are trivially parallelizable
- Even with LSTMs / GRUs, preserving important linguistic context over **very** long sequences is difficult
  - Transformers don't even try to remember things (every step looks at a weighted combination of **all** words in the input sentence)
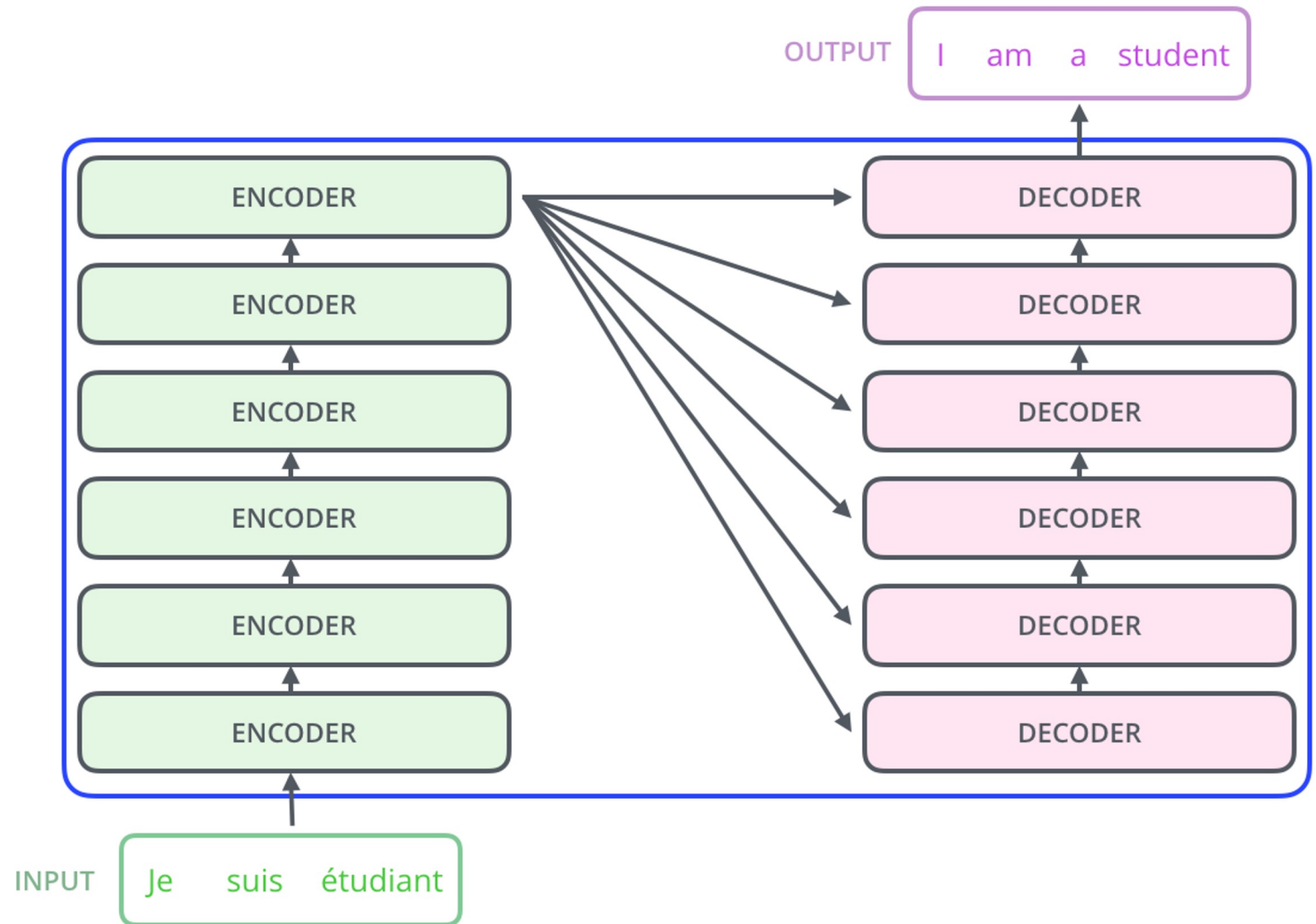
# Review: Transformer Model Overview

- The Transformer model breaks down into Encoder and Decoder blocks.

- At a high level, similar to the seq2seq architecture we've seen already...

- ...but there are no recurrent nets inside the Encoder and Decoder blocks!

Alammar, Jay (2018). The Illustrated Transformer [Blog post]. Retrieved from https://jalammar.github.io/illustrated-transformer/
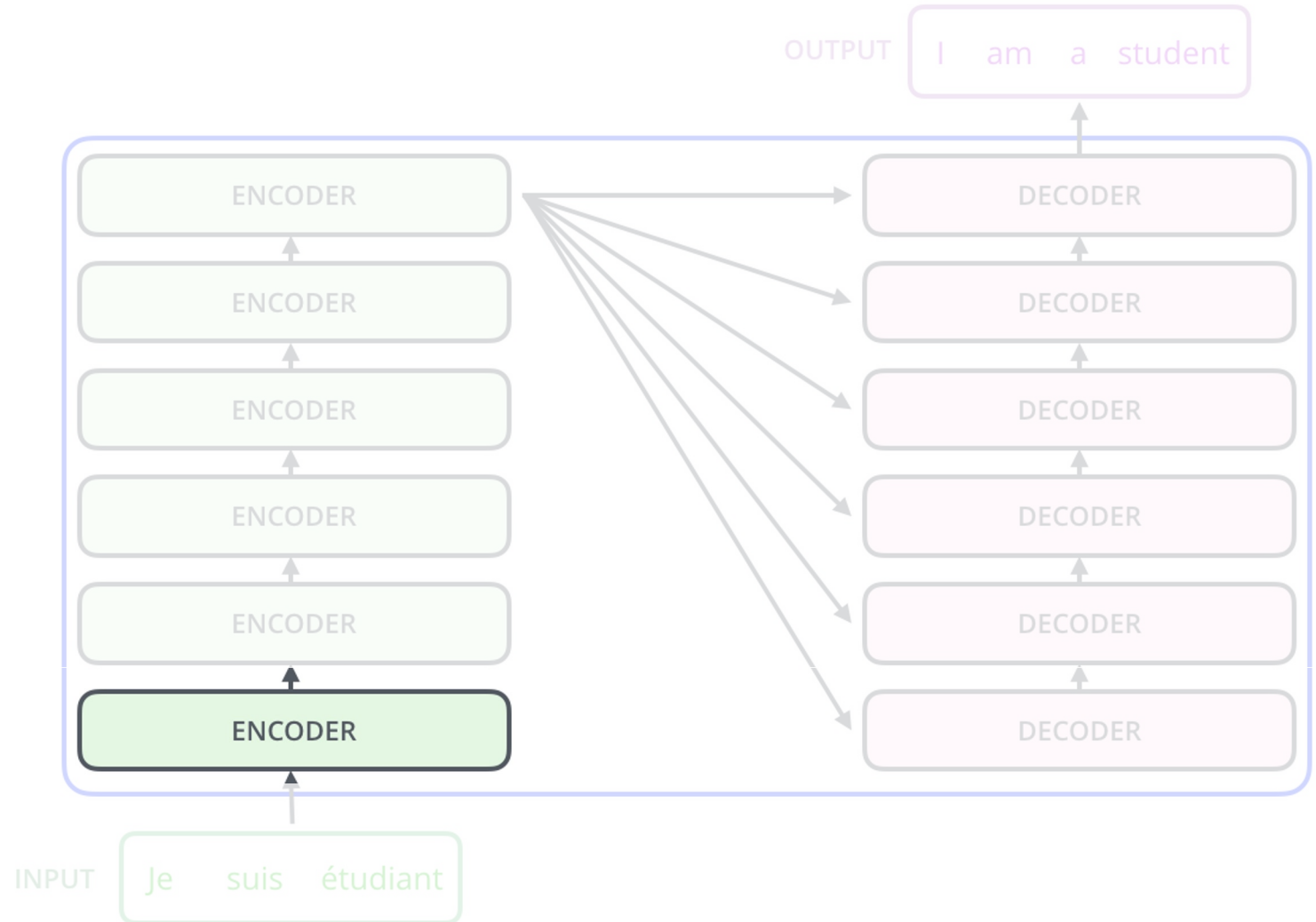
# Review: Transformer Model Overview

- The Transformer model breaks down into Encoder and Decoder blocks.

- At a high level, similar to the seq2seq architecture we've seen already...

- ...but there are no recurrent nets inside the Encoder and Decoder blocks!

- For better performance, often stack multiple Encoder and Decoder blocks (deeper network)

OUTPUT | I   am   a   student

ENCODER

ENCODER

ENCODER

ENCODER

ENCODER

ENCODER

DECODER

DECODER

DECODER

DECODER

DECODER

DECODER

INPUT | Je   suis   étudiant

# Transformer Model Overview

- Let's look at what goes on inside one of these Encoder blocks



Alammar, Jay (2018). The Illustrated Transformer [Blog post]. Retrieved from https://jalammar.github.io/illustrated-transformer/
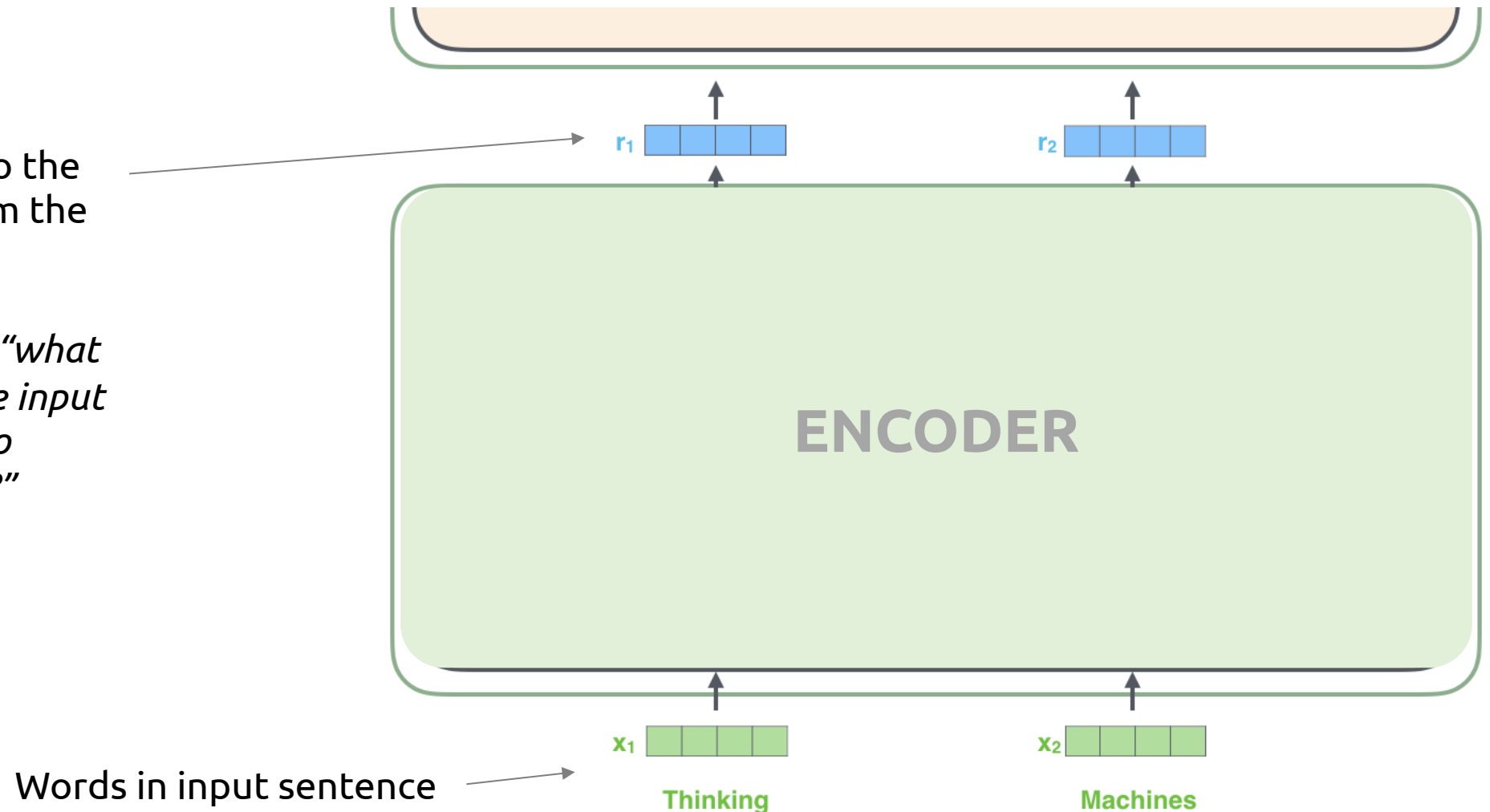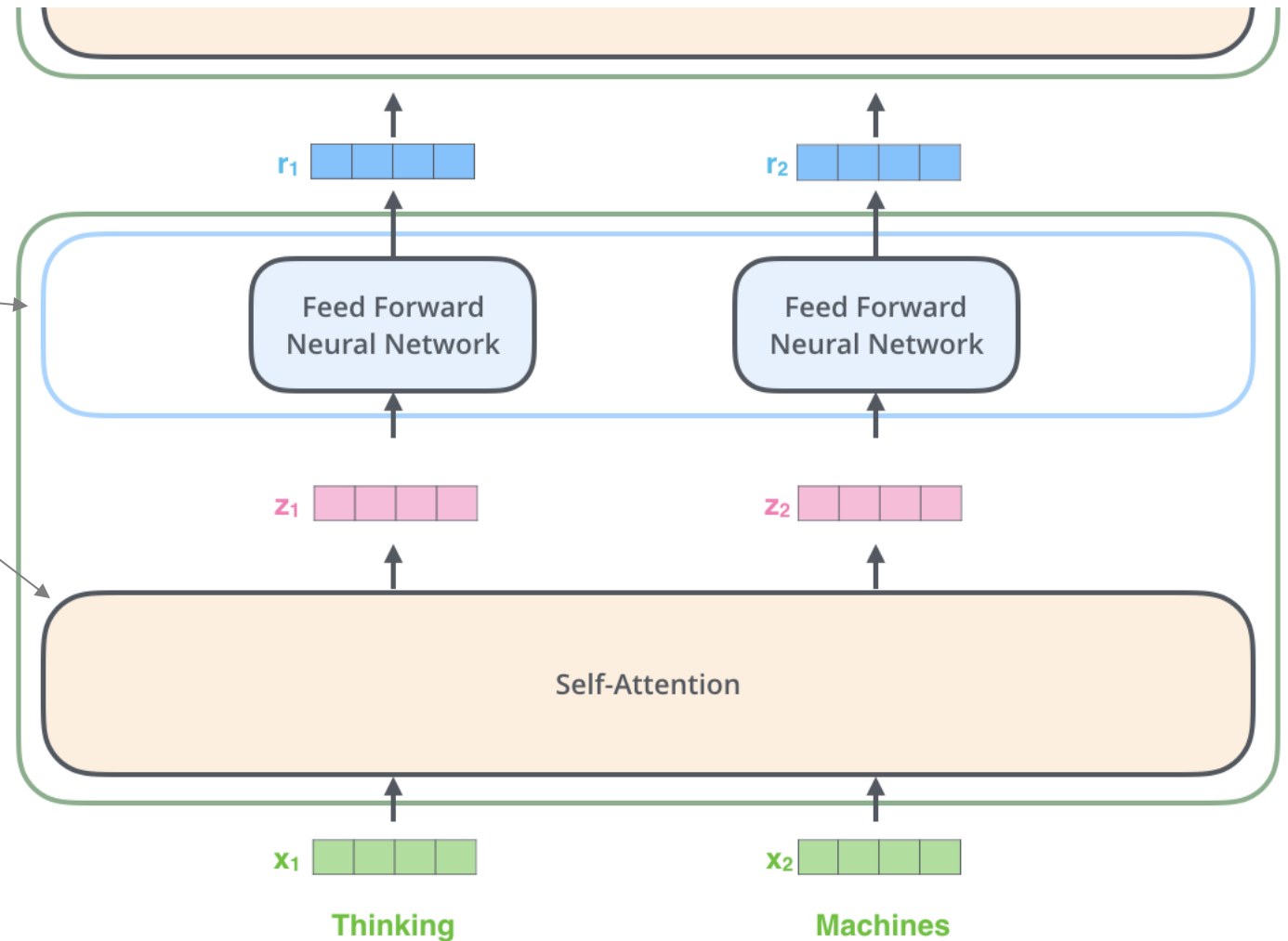
8

# Encoder Block Map

These per-word output vectors are analogous to the LSTM hidden states from the seq2seq2 model

- They should capture *"what information about the input sentence is relevant to translating this word?"*
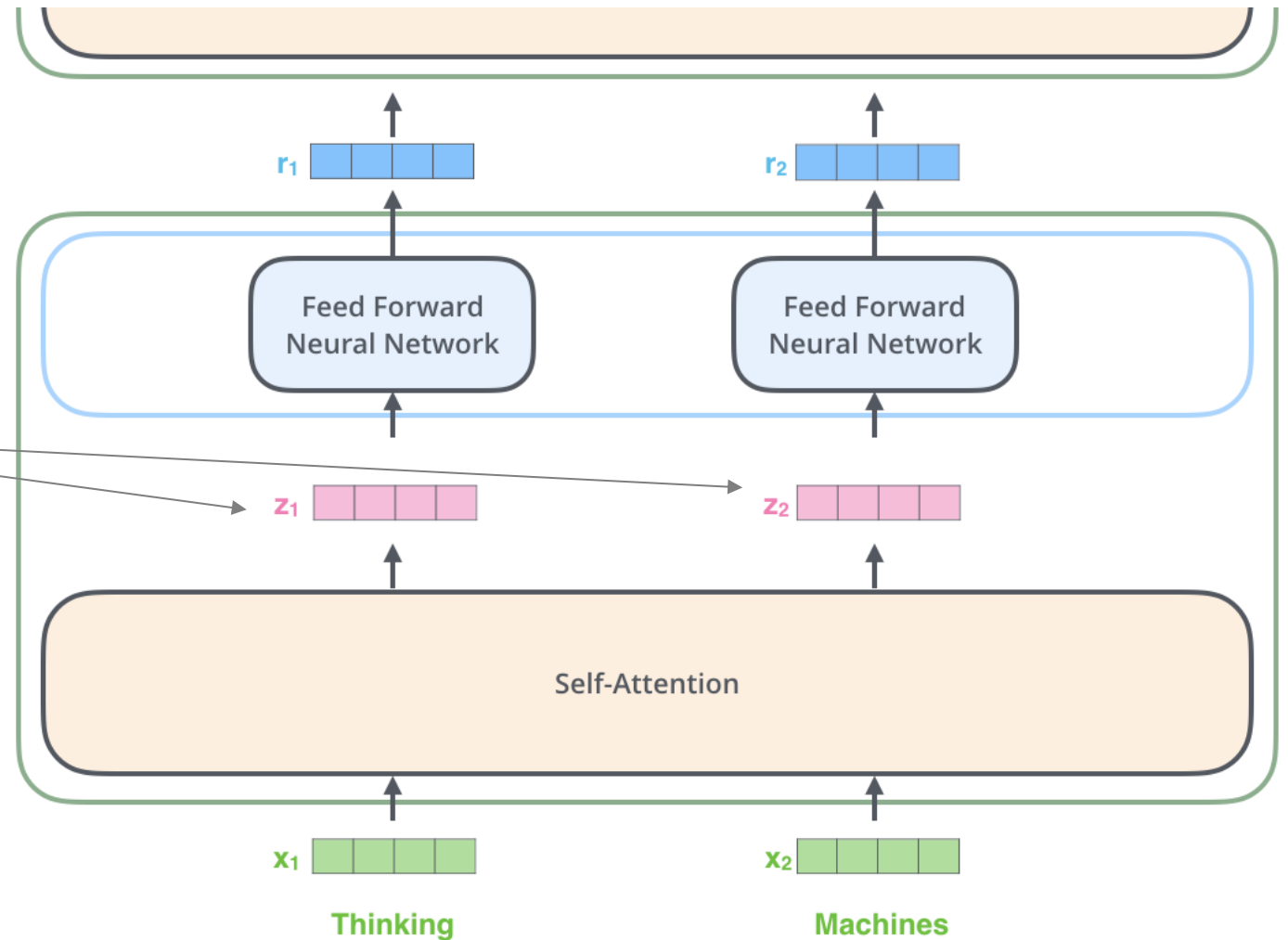
$r_1$ ⬛⬛⬛⬛    $r_2$ ⬛⬛⬛⬛

**ENCODER**

$x_1$ 🟩🟩🟩🟩    $x_2$ 🟩🟩🟩🟩

**Thinking**    **Machines**

Words in input sentence

# Encoder Block Map

- Encoder block breaks down into two main parts: Self-Attention, and Feed Forward layers.



$r_1$ | $r_2$

Feed Forward Neural Network | Feed Forward Neural Network

$z_1$ | $z_2$

Self-Attention

$x_1$ | $x_2$
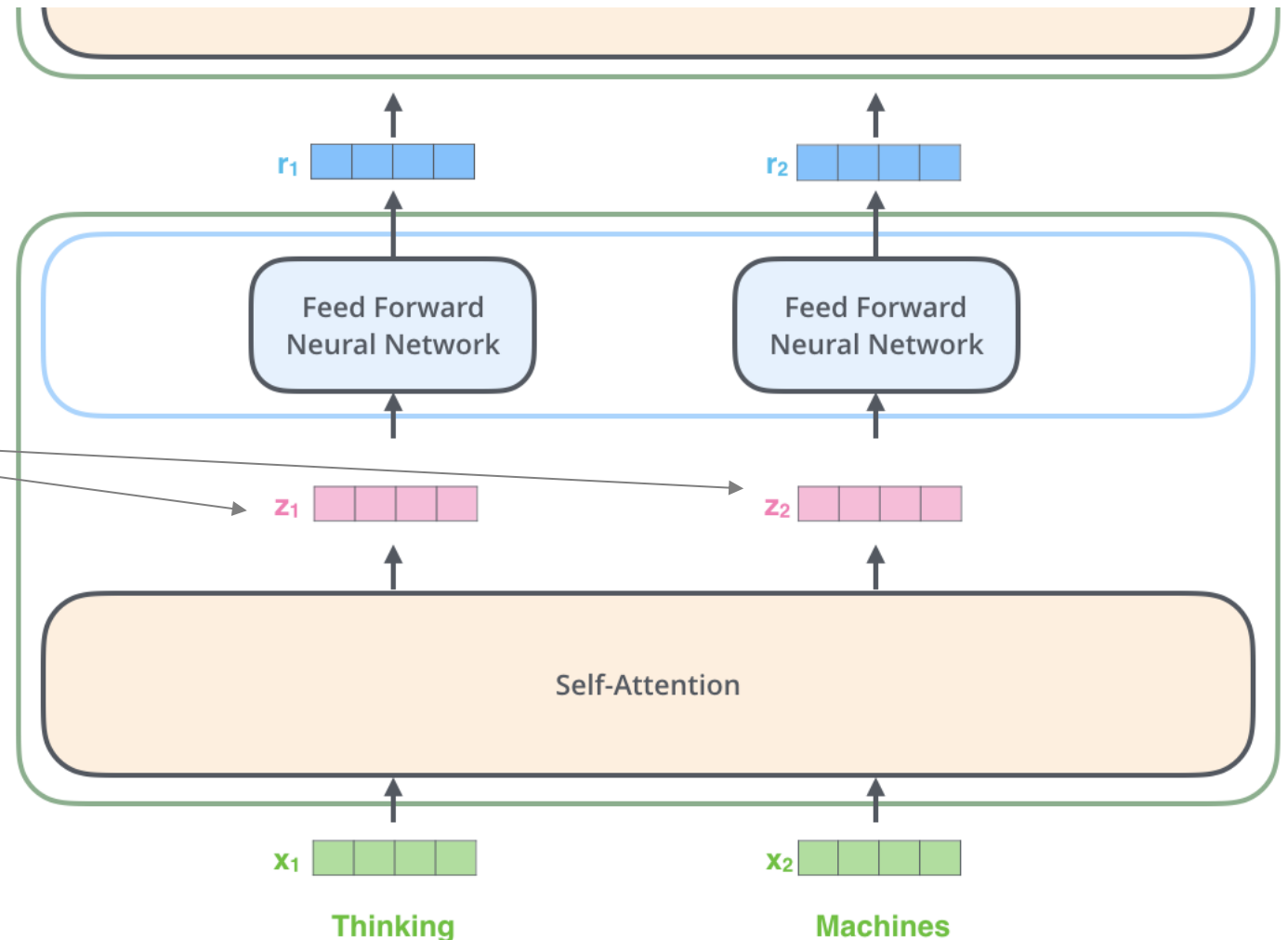
Thinking | Machines

# Encoder Block Map

- Encoder block breaks down into two main parts: Self-Attention, and Feed Forward layers.

- Self-Attention layer is applied to each word individually.

$r_1$

$r_2$

Feed Forward Neural Network

Feed Forward Neural Network

$z_1$

$z_2$

Self-Attention

$x_1$

$x_2$

**Thinking**

**Machines**

# Encoder Block Map

- Encoder block breaks down into two main parts: Self-Attention, and Feed Forward layers.

- Self-Attention layer is applied to each word individually.

*Let's revisit self-attention in detail!*

Alammar, Jay (2018). The Illustrated Transformer [Blog post]. Retrieved from https://jalammar.github.io/illustrated-transformer/

# Review: Attention types

$$\alpha_{t,i} = \text{align}(y_t, x_i) = \frac{\exp(\text{score}(\boldsymbol{s}_{t-1}, \boldsymbol{h}_i))}{\sum_{i'=1}^{n} \exp(\text{score}(\boldsymbol{s}_{t-1}, \boldsymbol{h}_{i'}))}$$

Softmax of some predefined alignment score..

How well two words $y_t$ and $x_i$ are aligned.

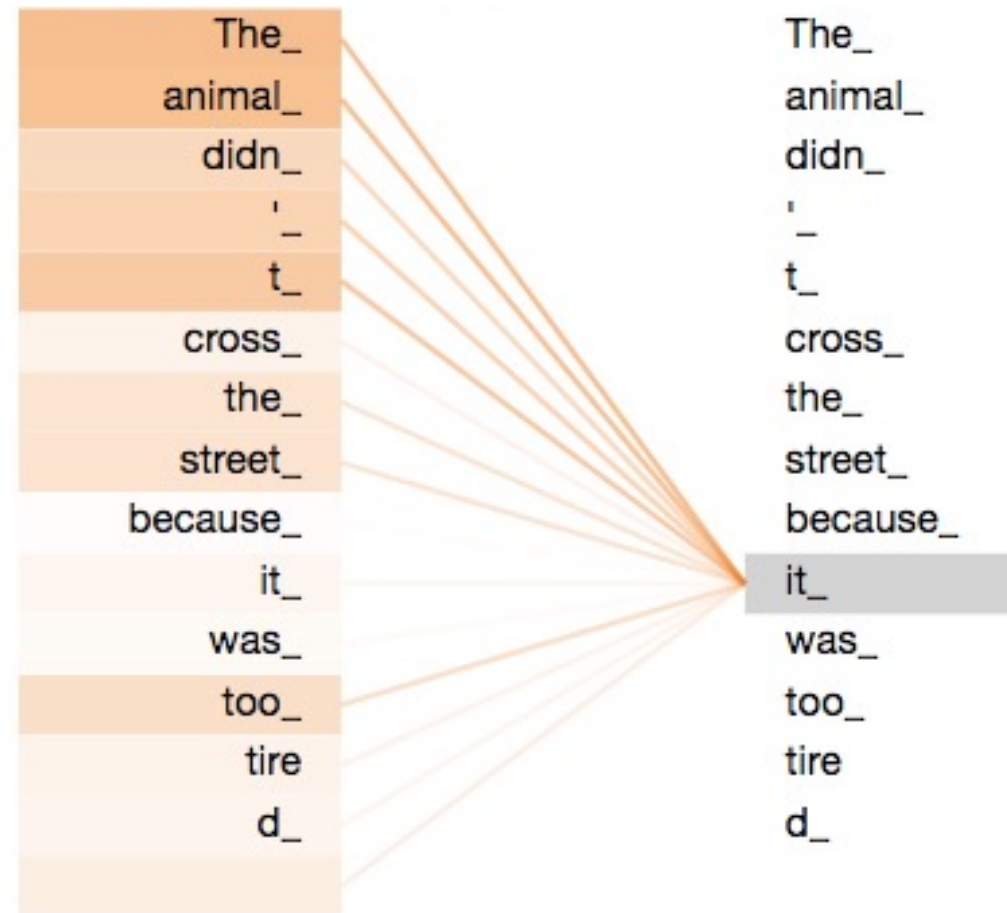| Name | Definition | Citation |
|------|------------|----------|
| Self-Attention(&) | Relating different positions of the same input sequence. Theoretically the self-attention can adopt any score functions above, but just replace the target sequence with the same input sequence. | Cheng2016 |
| Global/Soft | Attending to the entire input state space. | Xu2015 |
| Local/Hard | Attending to the part of input state space; i.e. a patch of the input image. | Xu2015; Luong2015 |

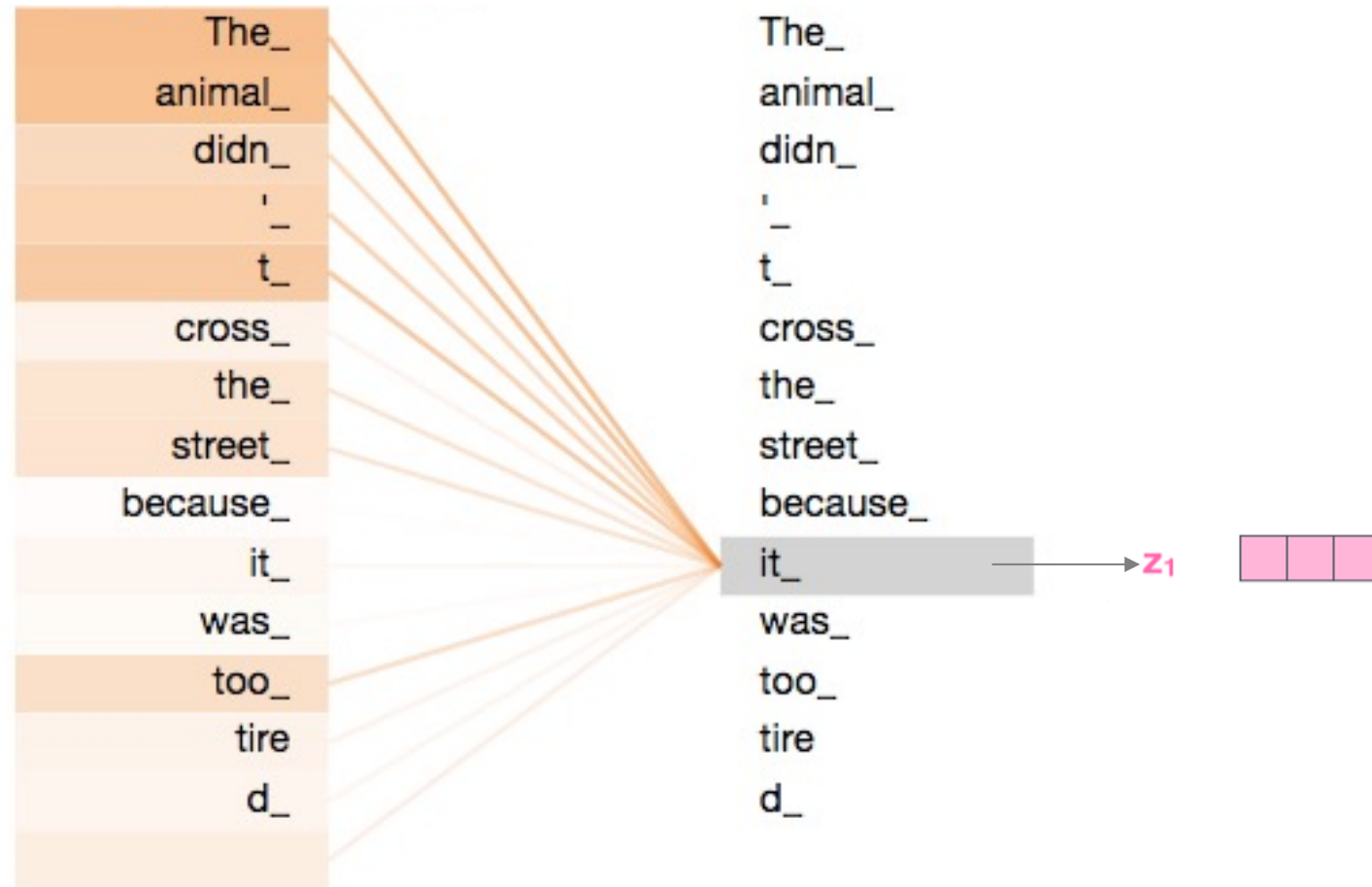# Self-Attention: Input's attention on itself

What do we do next?

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

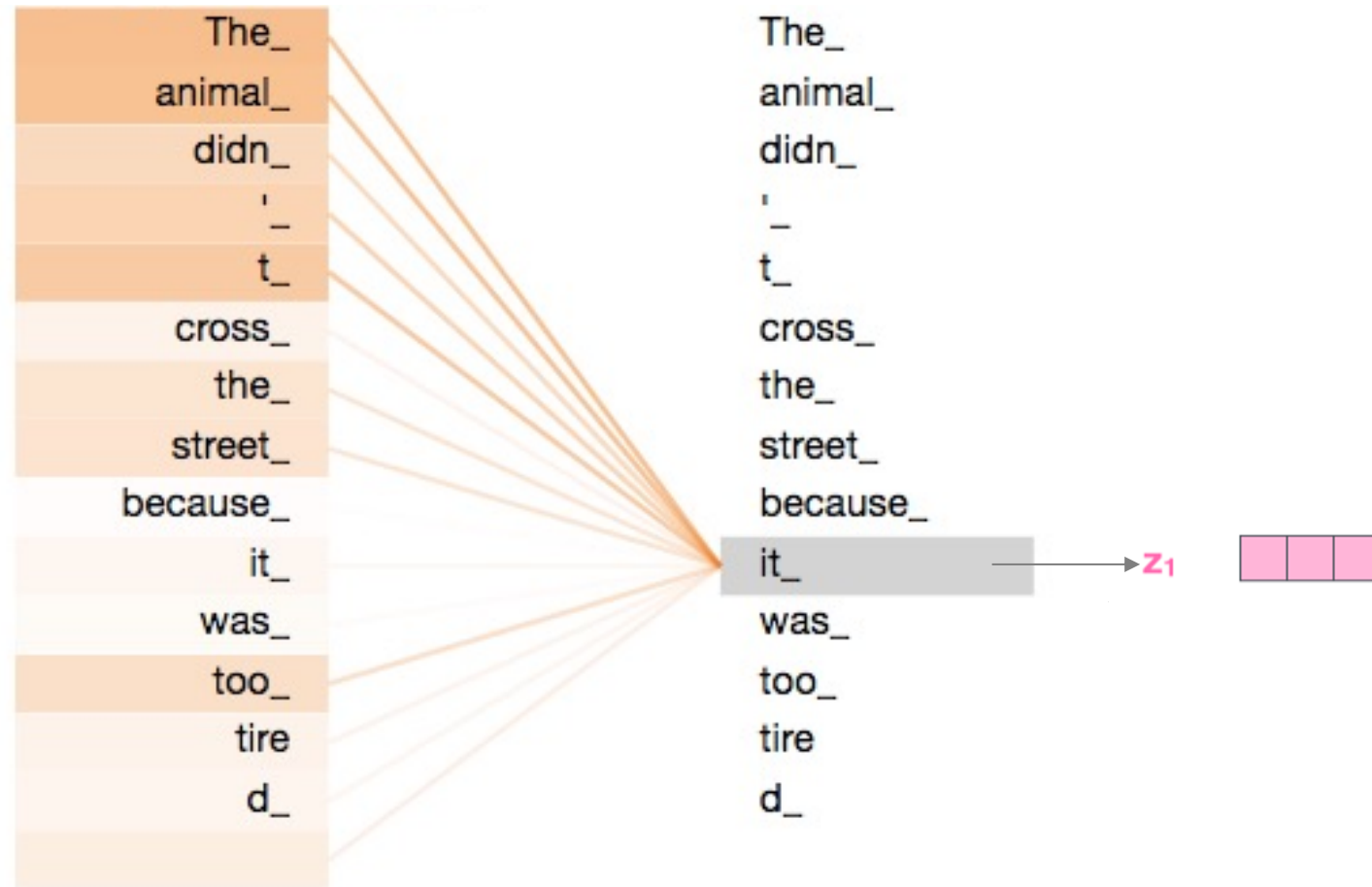# Self-Attention: Input's attention on itself

What do we do next?

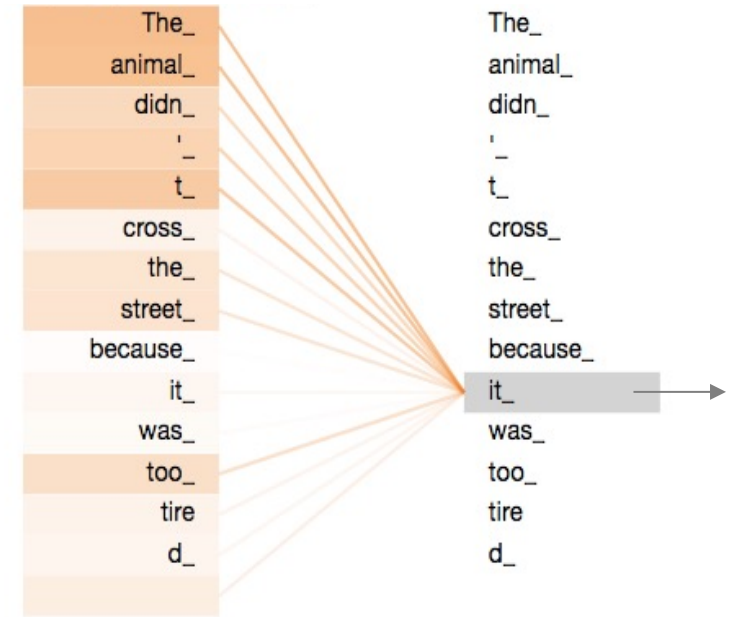# Self-Attention: Input's attention on itself

# Self-Attention: Overview

- **The big idea:**
  Self-attention computes the output vector $z_i$ for each word via a weighted sum of vectors extracted from each word in the input sentence

- Here, self-attention learns that "it" should pay attention to "the animal" (i.e. the entity that "it" refers to)

- Why the name **self**-attention?
  This describes attention that the input sentence pays to itself

Alammar, Jay (2018). The Illustrated Transformer [Blog post]. Retrieved from https://jalammar.github.io/illustrated-transformer/

# Self-Attention: Sketch



Alammar, Jay (2018). The Illustrated Transformer [Blog post]. Retrieved from https://jalammar.github.io/illustrated-transformer/

# Self-Attention: Overview

**Key vectors**

**How it works:**

1. To determine how much attention a word should pay to each other other, we compute
a query vector for the word and compare it to a
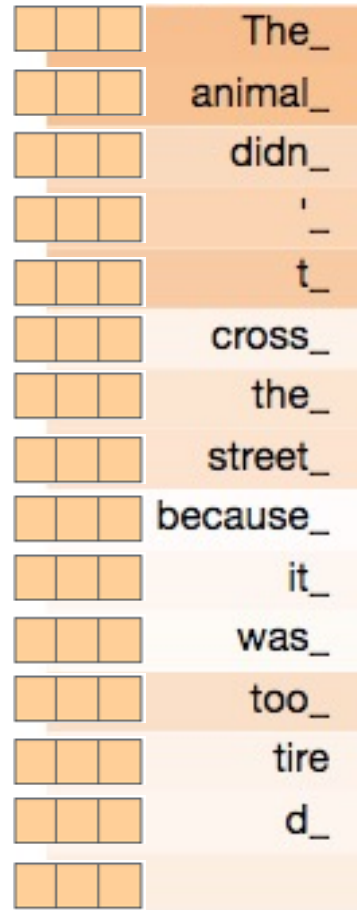key vector for every other word...
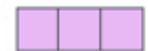
**Query vector**

# Self-Attention: Overview

**How it works:**

1. To determine how much attention a word should pay to each other other, we compute a query vector for the word and compare it to a key vector for every other word...



Key vectors

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

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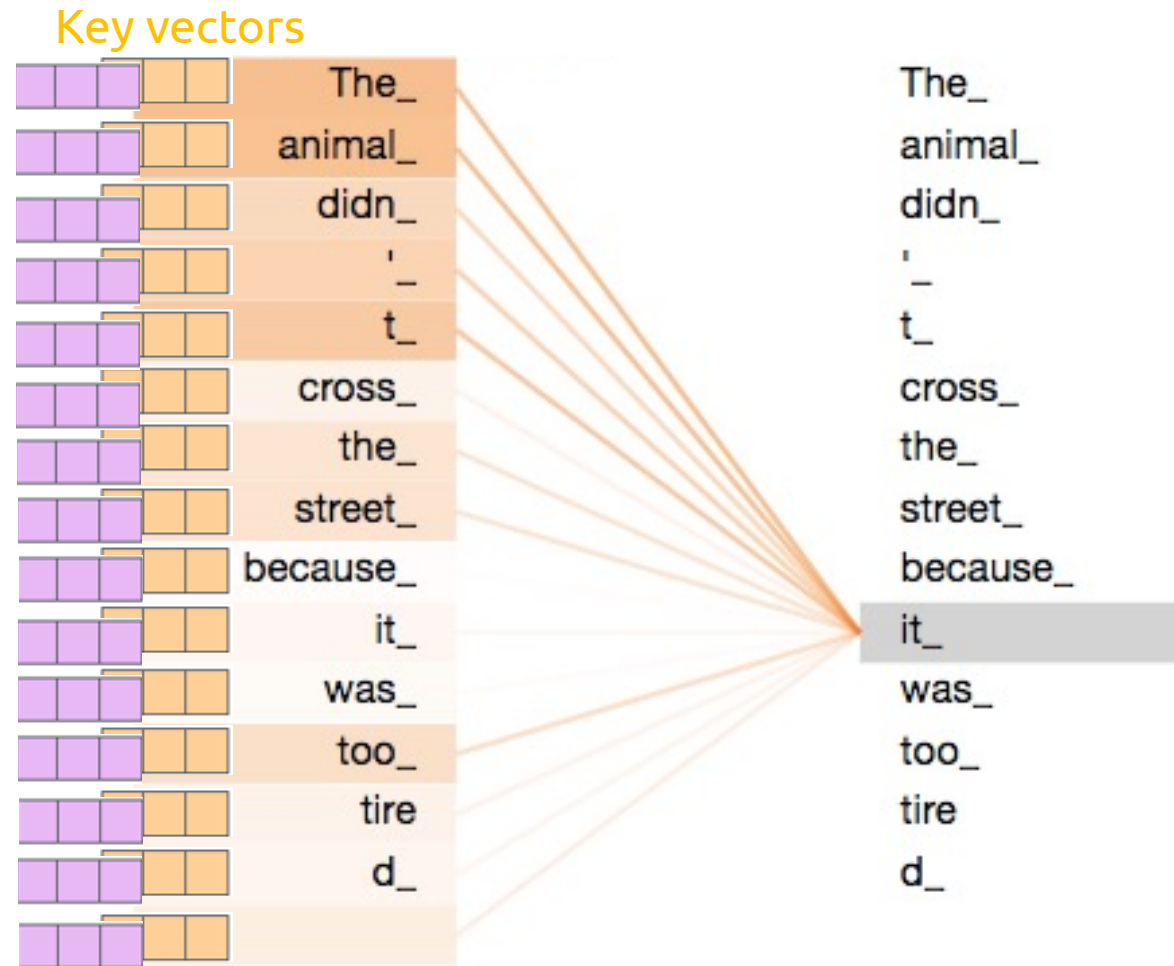didn_
'_
t_
cross_
the_
street_
because_
it_
was_
too_
tire
d_

# Self-Attention: Overview

**How it works:**
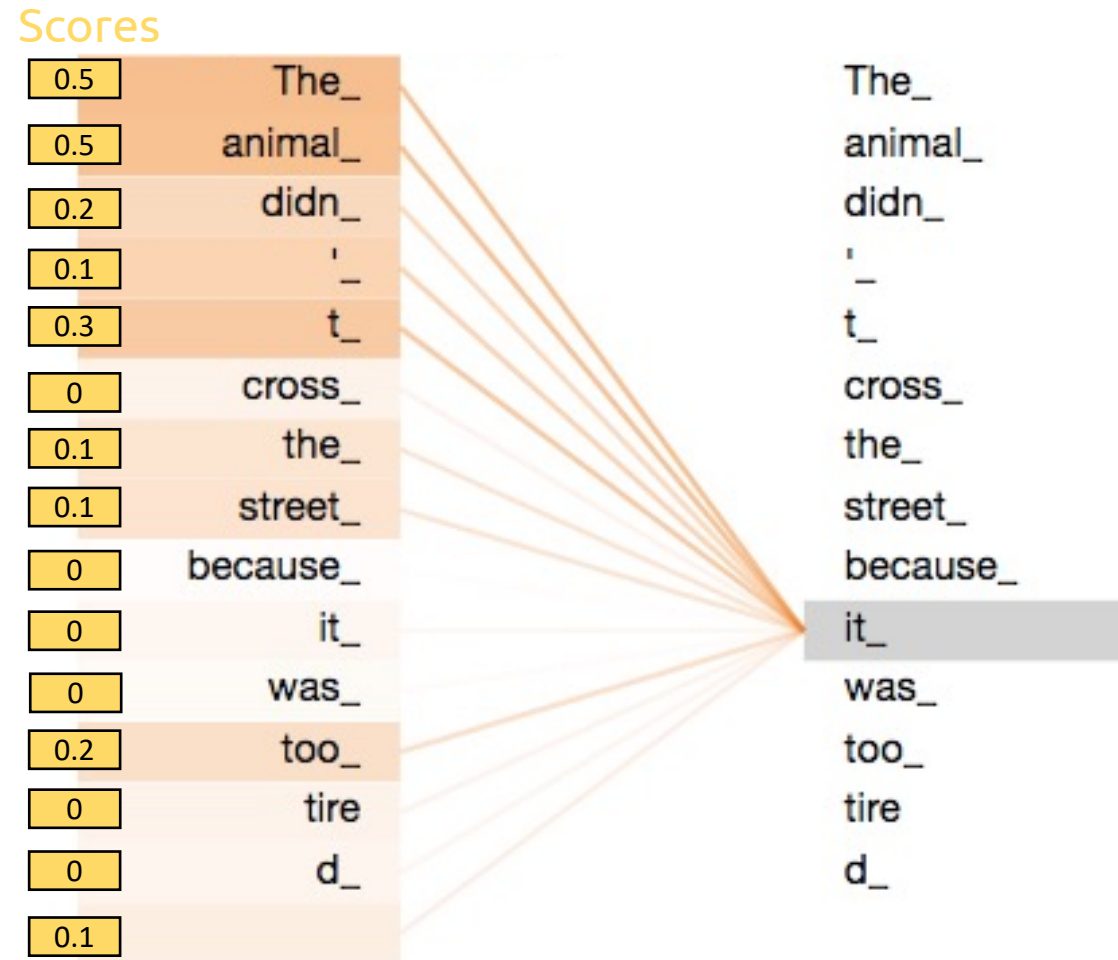
1. To determine how much attention a word should pay to each other other, we compute a query vector for the word and compare it to a key vector for every other word...
**to compute our alignment score**

Scores

| | |
|---|---|
| 0.5 | The_ |
| 0.5 | animal_ |
| 0.2 | didn_ |
| 0.1 | '_ |
| 0.3 | t_ |
| 0 | cross_ |
| 0.1 | the_ |
| 0.1 | street_ |
| 0 | because_ |
| 0 | it_ |
| 0 | was_ |
| 0.2 | too_ |
| 0 | tire |
| 0 | d_ |
| 0.1 | |

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

# Self-Attention: Overview

**How it works:**

1. To determine how much attention a word should pay to each other other, we compute a query vector for the word and compare it to a key vector for every other word... **to compute our alignment score**

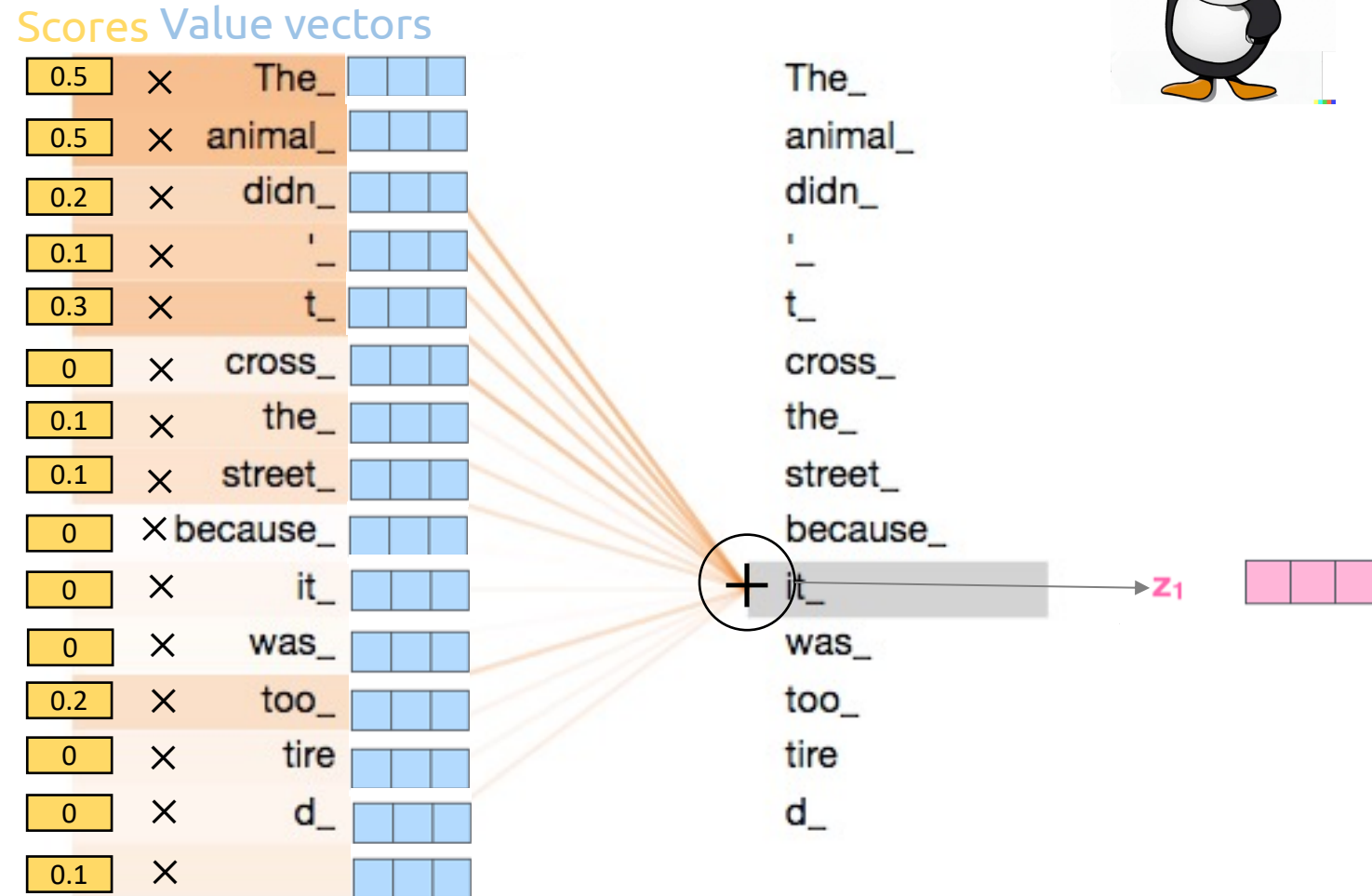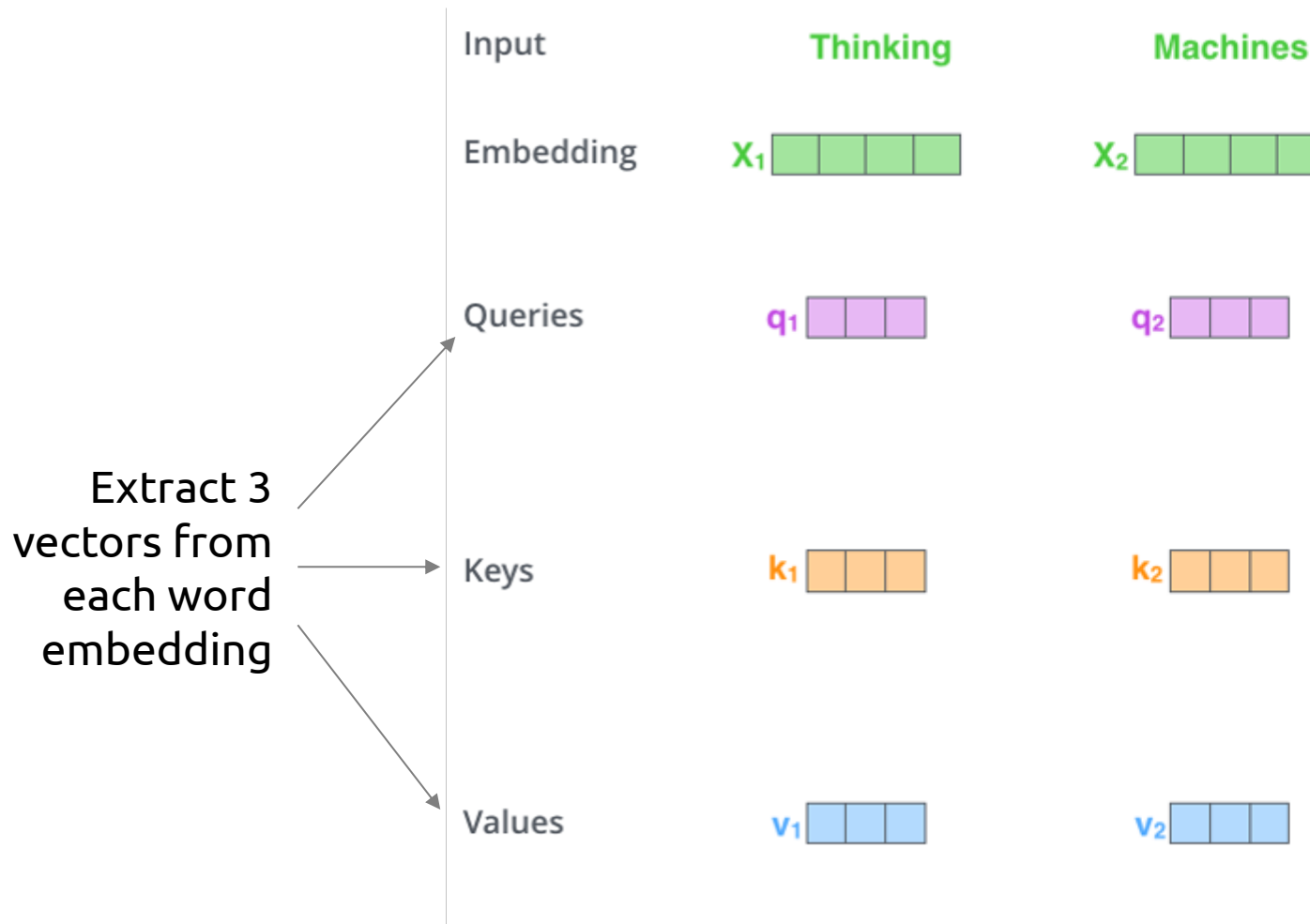2. To produce the output vector, we sum up the value vectors for each word, weighted by the score we computed in step 1
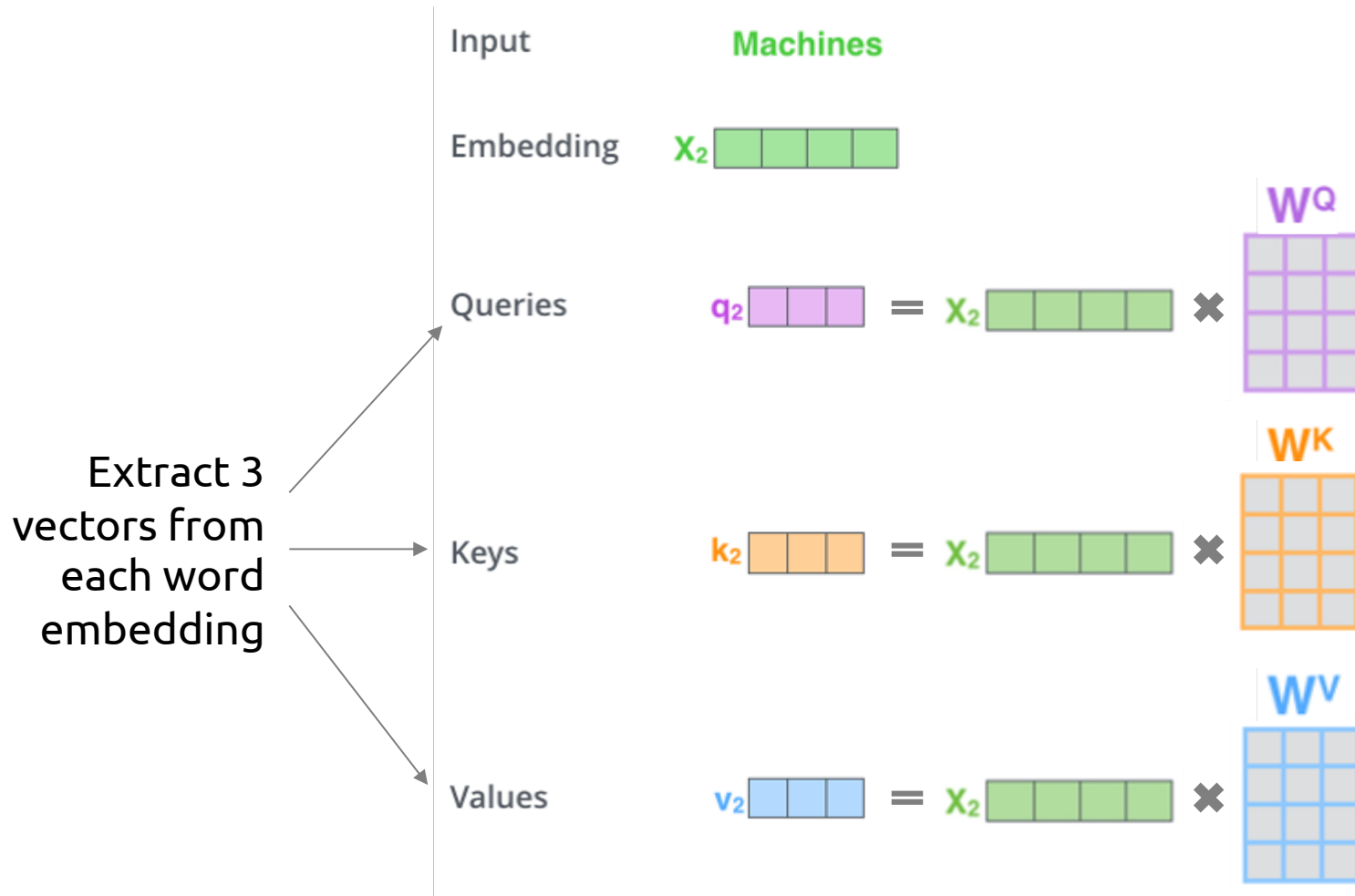


Scores  Value vectors

| 0.5 | × | The_ |
| 0.5 | × | animal_ |
| 0.2 | × | didn_ |
| 0.1 | × | '_ |
| 0.3 | × | t_ |
| 0 | × | cross_ |
| 0.1 | × | the_ |
| 0.1 | × | street_ |
| 0 | × | because_ |
| 0 | × | it_ |
| 0 | × | was_ |
| 0.2 | × | too_ |
| 0 | × | tire |
| 0 | × | d_ |
| 0.1 | × | |

The_
animal_
didn_
'_
t_
cross_
the_
street_
because_
it_  $z_1$
was_
too_
tire
d_

Alammar, Jay (2018). The Illustrated Transformer [Blog post]. Retrieved from https://jalammar.github.io/illustrated-transformer/

# Self-Attention: Details

Alammar, Jay (2018). The Illustrated Transformer [Blog post]. Retrieved from https://jalammar.github.io/illustrated-transformer/

# Self-Attention: Details

Extract 3 vectors from each word embedding

# Self-Attention: Details



Input: **Machines**

Embedding: $X_2$

Queries: $q_2 = X_2 \times W^Q$

Keys: $k_2 = X_2 \times W^K$

Values: $v_2 = X_2 \times W^V$

Extract 3 vectors from each word embedding

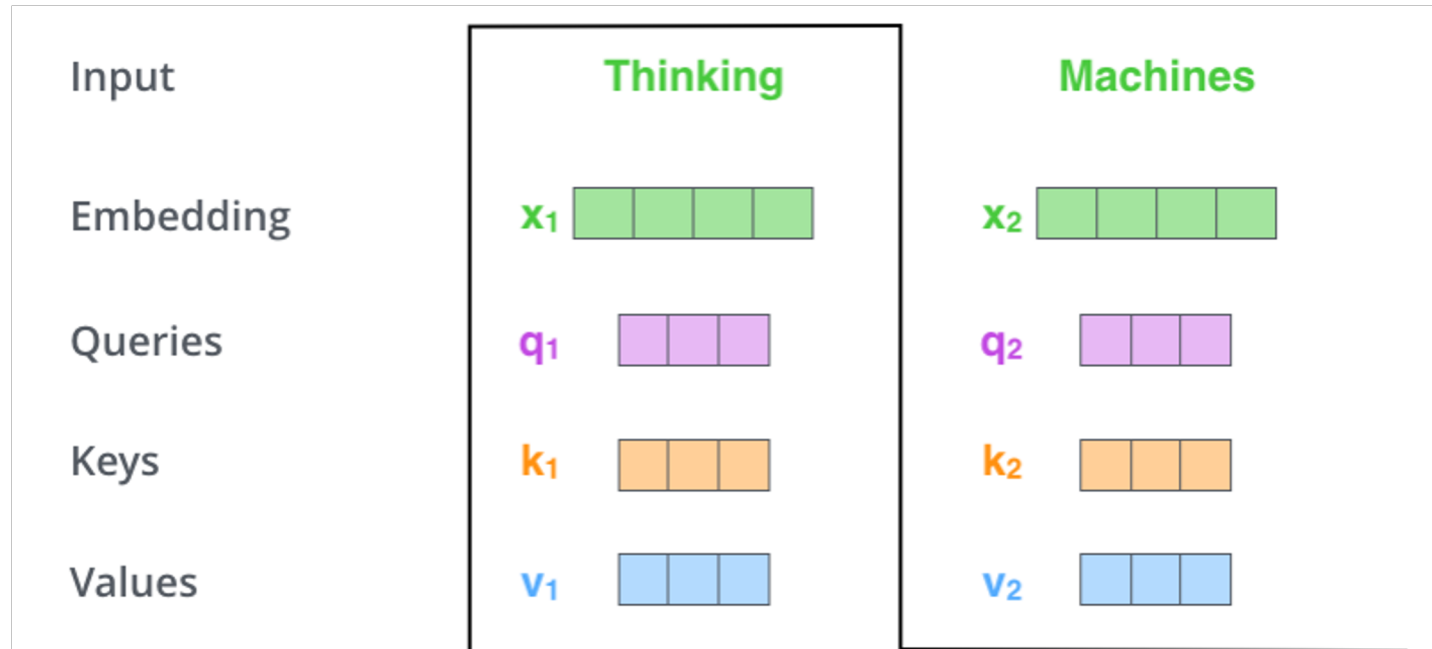Each vector is obtained by multiplying the embedding with the respective weight matrix.

How do we get these weight matrices?

These matrices are the **_trainable parameters_** of the network

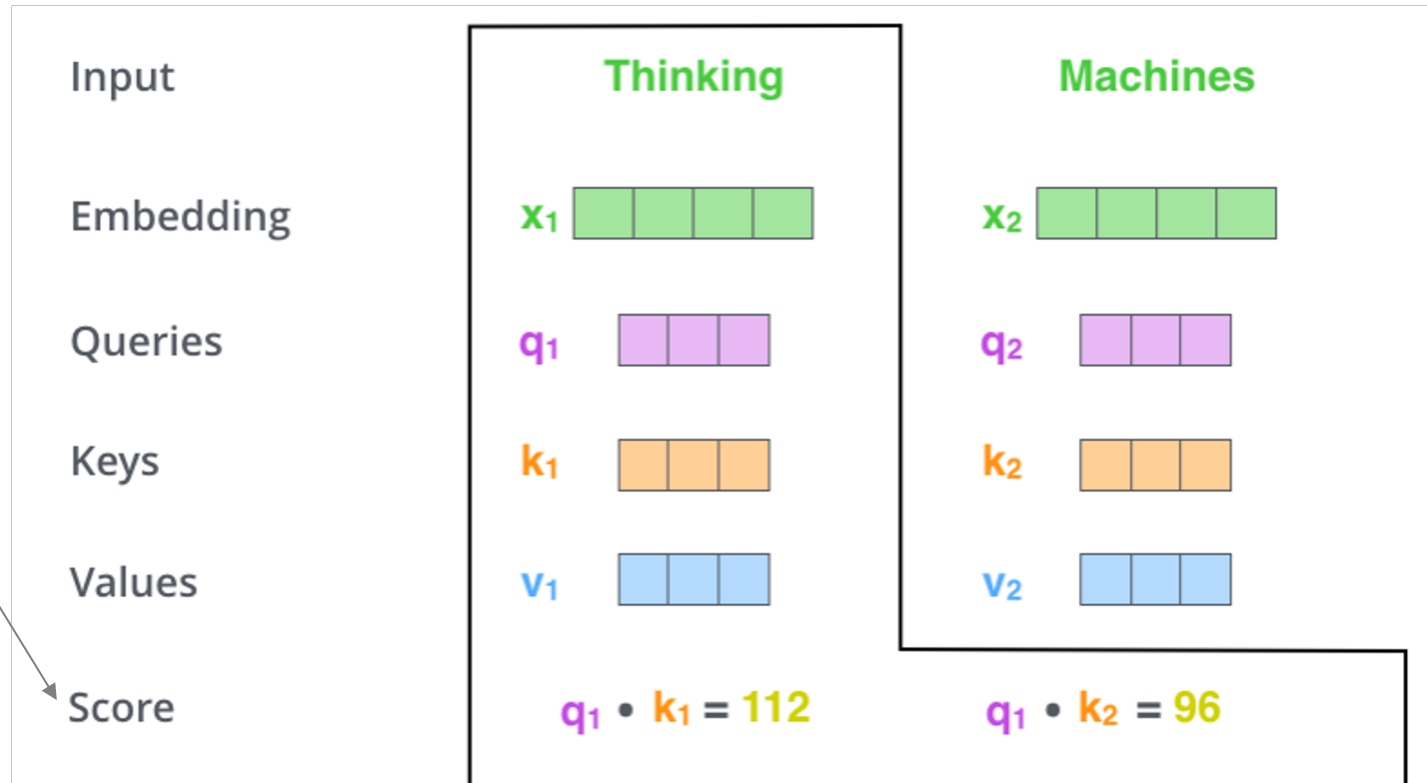# Self-Attention: Details

Computing self-attention for "Thinking"



What do we calculate next?

# Self-Attention: Details

Computing self-attention for "Thinking"

1. **Score:** Dot product query vector for "Thinking" ($q_1$) with the key vectors of each word in the sentence ($k_{1,2,\ldots,n}$).
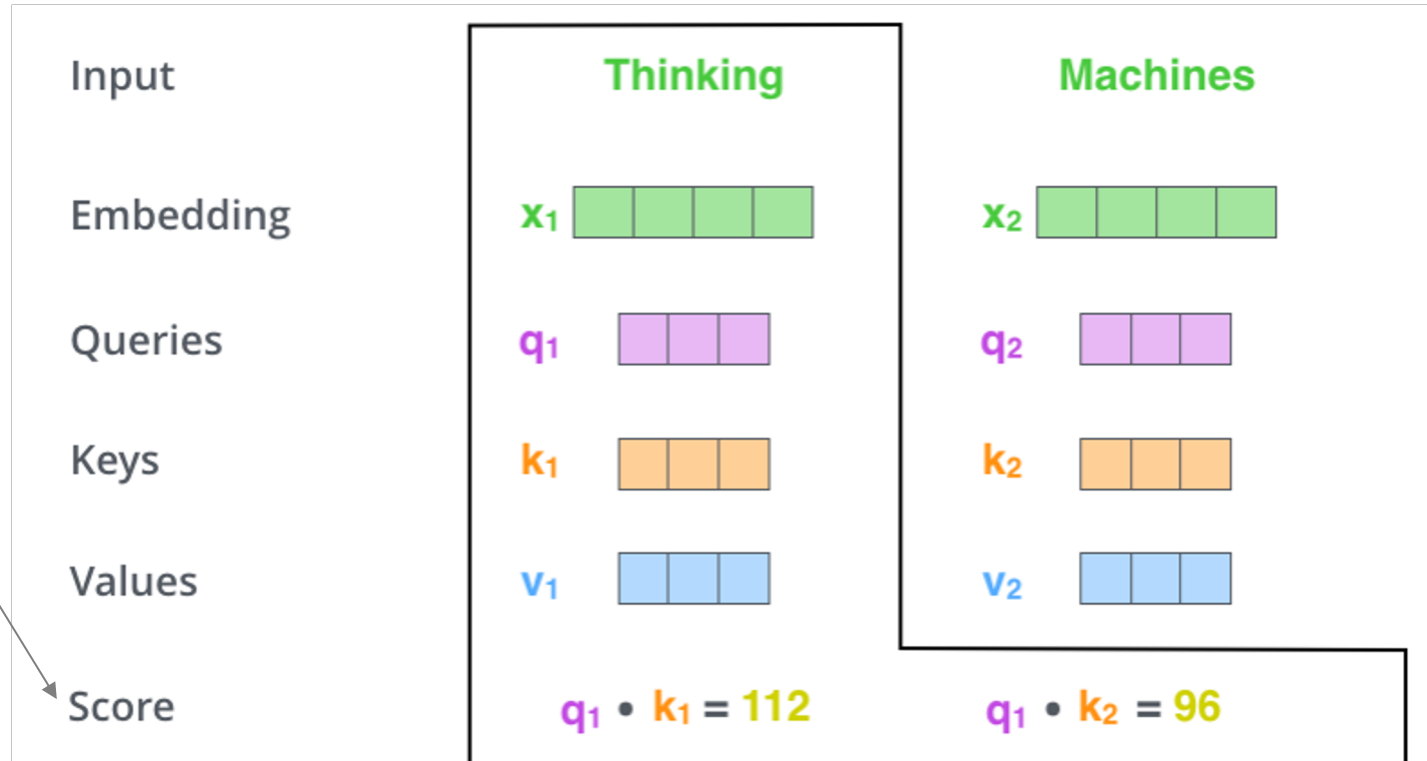
# Self-Attention: Details

Computing self-attention for "Thinking"

1. **Score:** Dot product query vector for "Thinking" ($q_1$) with the key vectors of each word in the sentence ($k_{1,2,...,n}$).

What this is essentially asking is: *How much should "Thinking" pay attention to each other word in the sequence?*

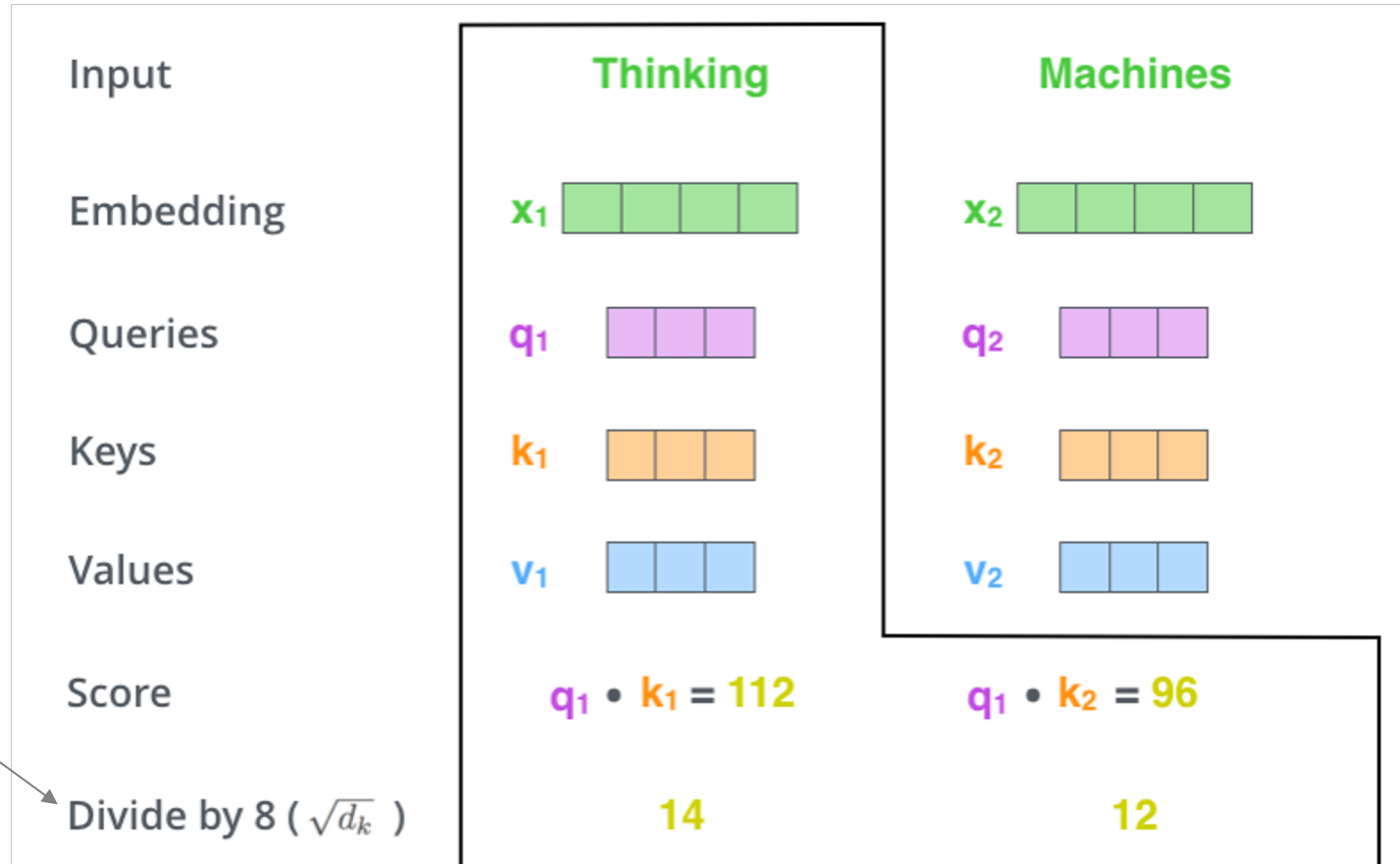Query vectors are asking the question and key vectors respond.



| | Thinking | Machines |
|---|---|---|
| Input | | |
| Embedding | $x_1$ | $x_2$ |
| Queries | $q_1$ | $q_2$ |
| Keys | $k_1$ | $k_2$ |
| Values | $v_1$ | $v_2$ |
| Score | $q_1 \cdot k_1 = 112$ | $q_1 \cdot k_2 = 96$ |

What do we calculate next?

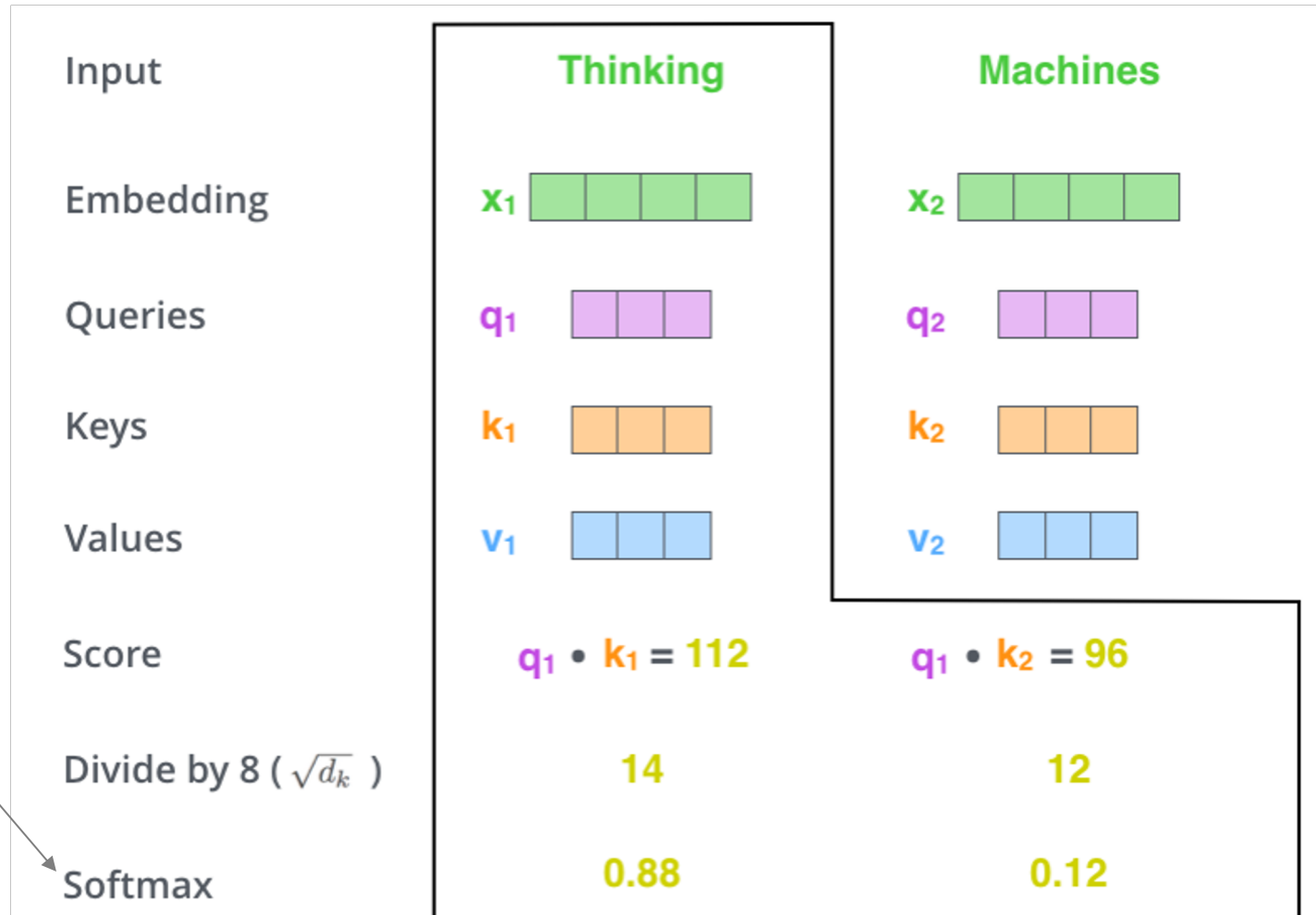# Self-Attention: Details

Computing self-attention for "Thinking"

1. **Score:** Dot product query vector for "Thinking" ($q_1$) with the key vectors of each word in the sentence ($k_{1,2,...,n}$).
2. **Scale:** Divide each score by square root of key vector dimensionality. Results in more stable gradients.



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

Alammar, Jay (2018). The Illustrated Transformer [Blog post]. Retrieved from https://jalammar.github.io/illustrated-transformer/
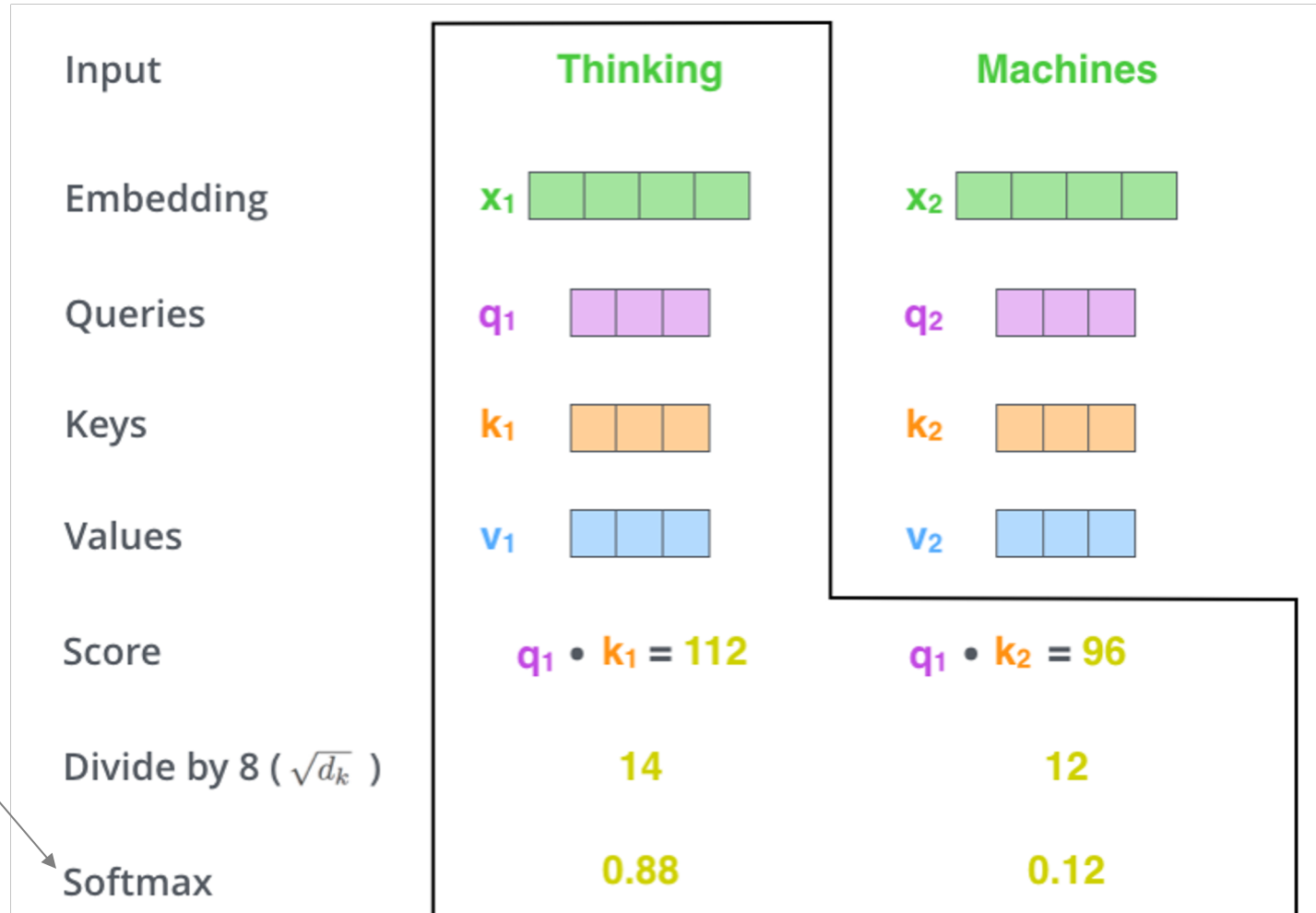
# Self-Attention: Details

Computing self-attention for "Thinking"

1. **Score:** Dot product query vector for "Thinking" ($q_1$) with the key vectors of each word in the sentence ($k_{1,2,...,n}$).
2. **Scale:** Divide each score by square root of key vector dimensionality. Results in more stable gradients.
3. **Softmax:** Apply softmax layer.

| | 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: Details

Computing self-attention for "Thinking"

1. **Score:** Dot product query vector for "Thinking" ($q_1$) with the key vectors of each word in the sentence ($k_{1,2,...,n}$).
2. **Scale:** Divide each score by square root of key vector dimensionality. Results in more stable gradients.
3. **Softmax:** Apply softmax layer.

By applying softmax, we transform the scores into attention weights.

What do we calculate next?



| | 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: Details

Computing self-attention for "Thinking"

1. **Score:** Dot product query vector for "Thinking" ($q_1$) with the key vectors of each word in the sentence ($k_{1,2,...,n}$).
2. **Scale:** Divide each score by square root of key vector dimensionality. Results in more stable gradients.
3. **Softmax:** Apply softmax layer.
4. **Weighting:** Multiply value vector of each word in the sentence ($v_{1,2,...,n}$) with the respective softmax values.

| | 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 | $v_1$ | $v_2$ |

# Self-Attention: Details

Computing self-attention for "Thinking"

1. **Score:** Dot product query vector for "Thinking" ($q_1$) with the key vectors of each word in the sentence ($k_{1,2,...,n}$).
2. **Scale:** Divide each score by square root of key vector dimensionality. Results in more stable gradients.
3. **Softmax:** Apply softmax layer.
4. **Weighting:** Multiply value vector of each word in the sentence ($v_{1,2,...,n}$) with the respective softmax values.

The idea here is that the value vectors store the contextual information that each word provides.

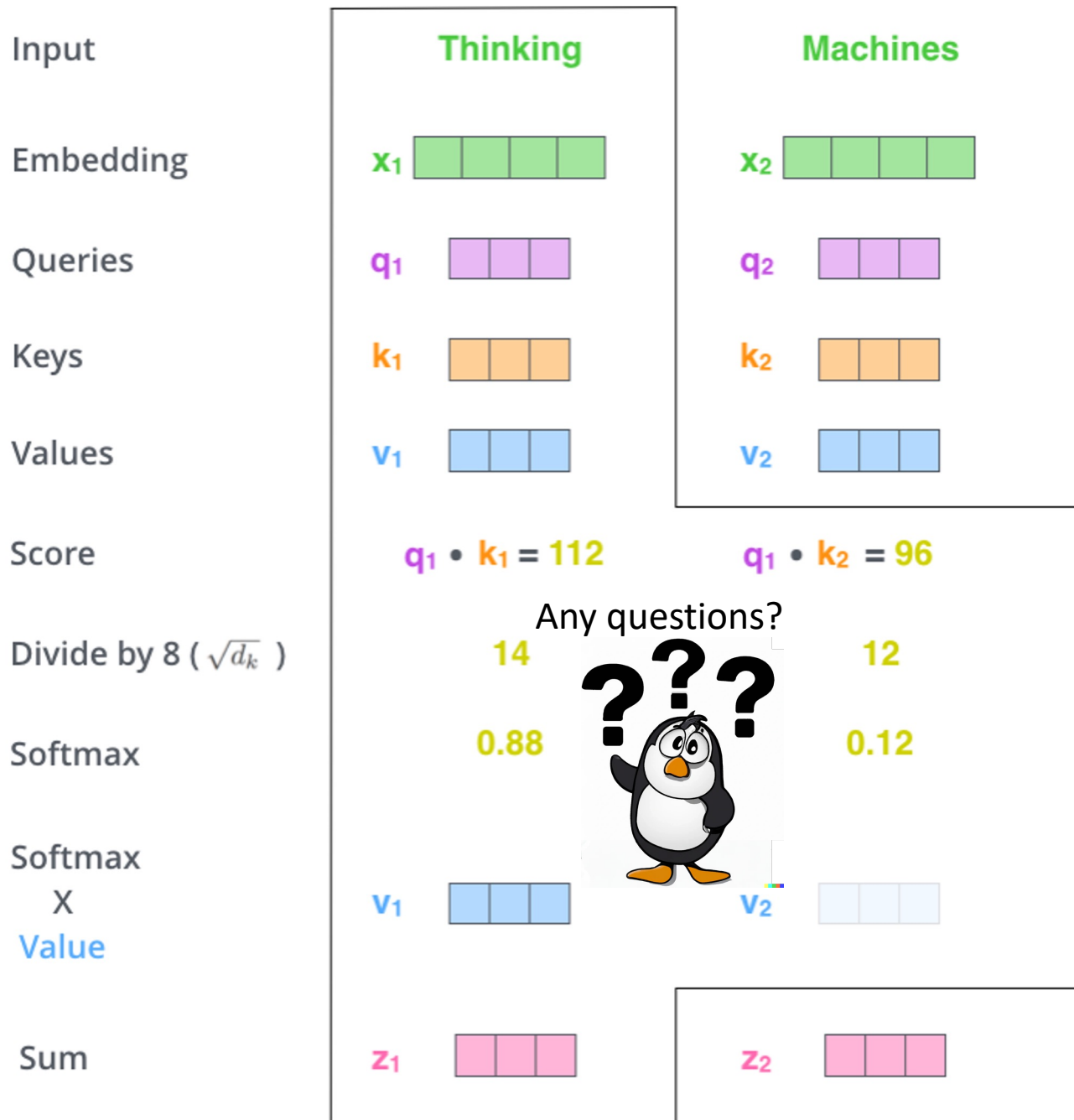| | Thinking | Machines |
|---|---|---|
| Input | | |
| Embedding | $x_1$ | $x_2$ |
| Queries | $q_1$ | $q_2$ |
| Keys | $k_1$ | $k_2$ |
| Values | $v_1$ | $v_2$ |
| Score | $q_1 \cdot k_1 = 112$ | $q_1 \cdot k_2 = 96$ |
| Divide by 8 ( $\sqrt{d_k}$ ) | 14 | 12 |
| Softmax | 0.88 | 0.12 |
| Softmax X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

# Self-Attention: Details

Computing self-attention for "Thinking"

1. **Score:** Dot product query vector for "Thinking" ($q_1$) with the key vectors of each word in the sentence ($k_{1,2,...,n}$).
2. **Scale:** Divide each score by square root of key vector dimensionality. Results in more stable gradients.
3. **Softmax:** Apply softmax layer.
4. **Weighting:** Multiply value vector of each word in the sentence ($v_{1,2,...,n}$) with the respective softmax values.
5. **Sum:** Sum up weighted value vectors ($v_{1,2,...,n}$) into one final self-attention vector for "Thinking" ($z_1$)
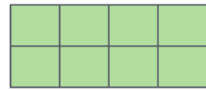
| | Thinking | Machines |
|---|---|---|
| Input | | |
| Embedding | $x_1$ | $x_2$ |
| Queries | $q_1$ | $q_2$ |
| Keys | $k_1$ | $k_2$ |
| Values | $v_1$ | $v_2$ |
| Score | $q_1 \bullet k_1 = 112$ | $q_1 \bullet k_2 = 96$ |
| Divide by 8 ( $\sqrt{d_k}$ ) | 14 | 12 |
| Softmax | 0.88 | 0.12 |
| Softmax X Value | $v_1$ | $v_2$ |
| Sum | $z_1$ | $z_2$ |

# Self-Attention: Details

Computing self-attention for "Thinking"

1. **Score:** Dot product query vector for "Thinking" ($q_1$) with the key vectors of each word in the sentence ($k_{1,2,...,n}$).
2. **Scale:** Divide each score by square root of key vector dimensionality. Results in more stable gradients.
3. **Softmax:** Apply softmax layer.
4. **Weighting:** Multiply value vector of each word in the sentence ($v_{1,2,...,n}$) with the respective softmax values.
5. **Sum:** Sum up weighted value vectors ($v_{1,2,...,n}$) into one final self-attention vector for "Thinking" ($z_1$)

We are weighting the *context* provided by each word by the amount of *attention* we should pay.

| | 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 | $z_1$ | $z_2$ |

Any questions?

???

Alammar, Jay (2018). The Illustrated Transformer [Blog post]. Retrieved from https://jalammar.github.io/illustrated-transformer/

# Self-Attention as a Matrix Computation

**X**

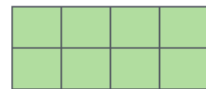Each row of **X** is a word embedding of a word in our sentence.

**X**

Get your pens/papers or tablets ready!

**X**

What would be the dimensions of the weight matrices to calculate the query, key, and value?

What would be the dimensions of the query, key, and value matrices?

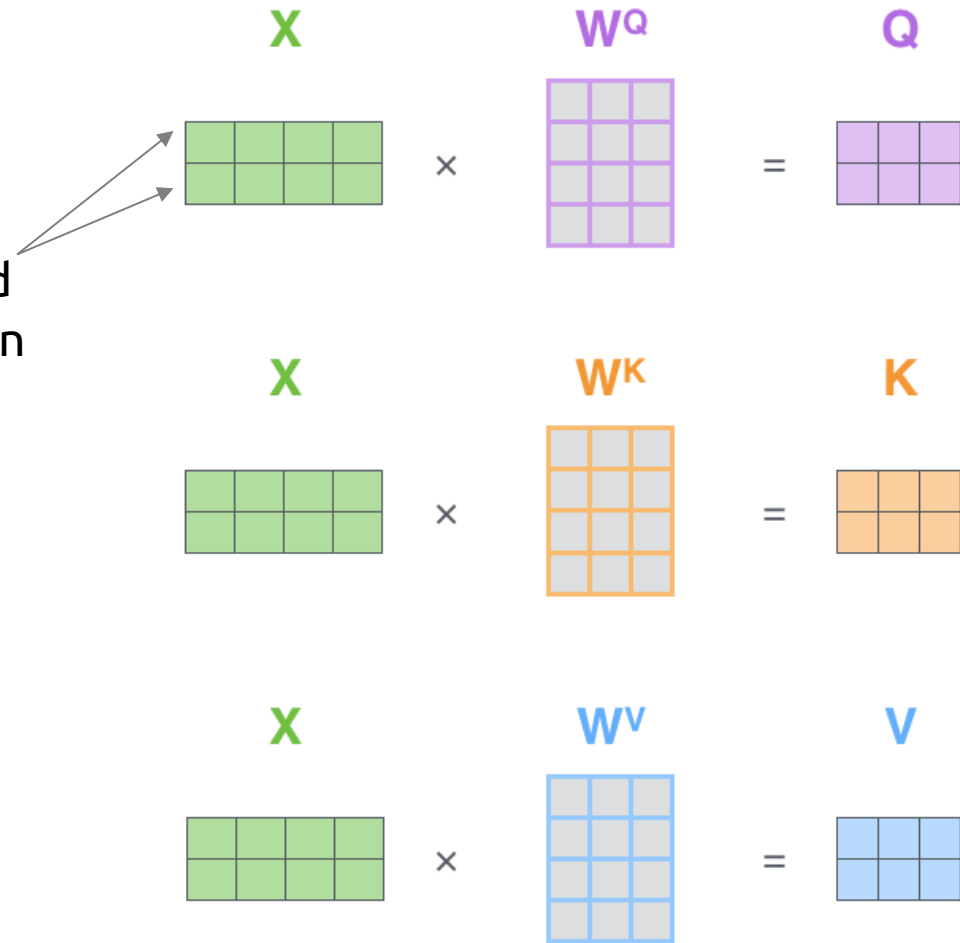Apply the steps of calculating attention weights on the query and key matrices.

What is the dimension of attention weight matrix?

Multiply the attention weight matrix to value matrix produce the output matrix.

What are the dimensions of output matrix?

# Self-Attention as a Matrix Computation



Each row of **X** is a word embedding of a word in our sentence.

Alammar, Jay (2018). The Illustrated Transformer [Blog post]. Retrieved from https://jalammar.github.io/illustrated-transformer/

# Self-Attention as a Matrix Computation



Each row of **X** is a word embedding of a word in our sentence.

Matrix multiplication with $W^Q$, $W^K$, and $W^V$ gives us matrices **Q**, **K**, and **V**, where the $i^{th}$ rows of each matrix represent the vectors $q_i$, $k_i$, and $v_i$.

# Self-Attention as a Matrix Computation

Matrix multiplying **Q** and the transpose of **K** calculates all the "score" values.



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

Dividing by $\sqrt{d_k}$ correctly scales values.

# Self-Attention as a Matrix Computation

Matrix multiplying **Q** and the transpose of **K** calculates all the "score" values.

The result is a **Z** matrix where the i[th] row represents the self-attention vector $z_i$

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

**Q** **K**[T] **V** **Z**

Any questions?

Dividing by $\sqrt{d_k}$ correctly scales values.

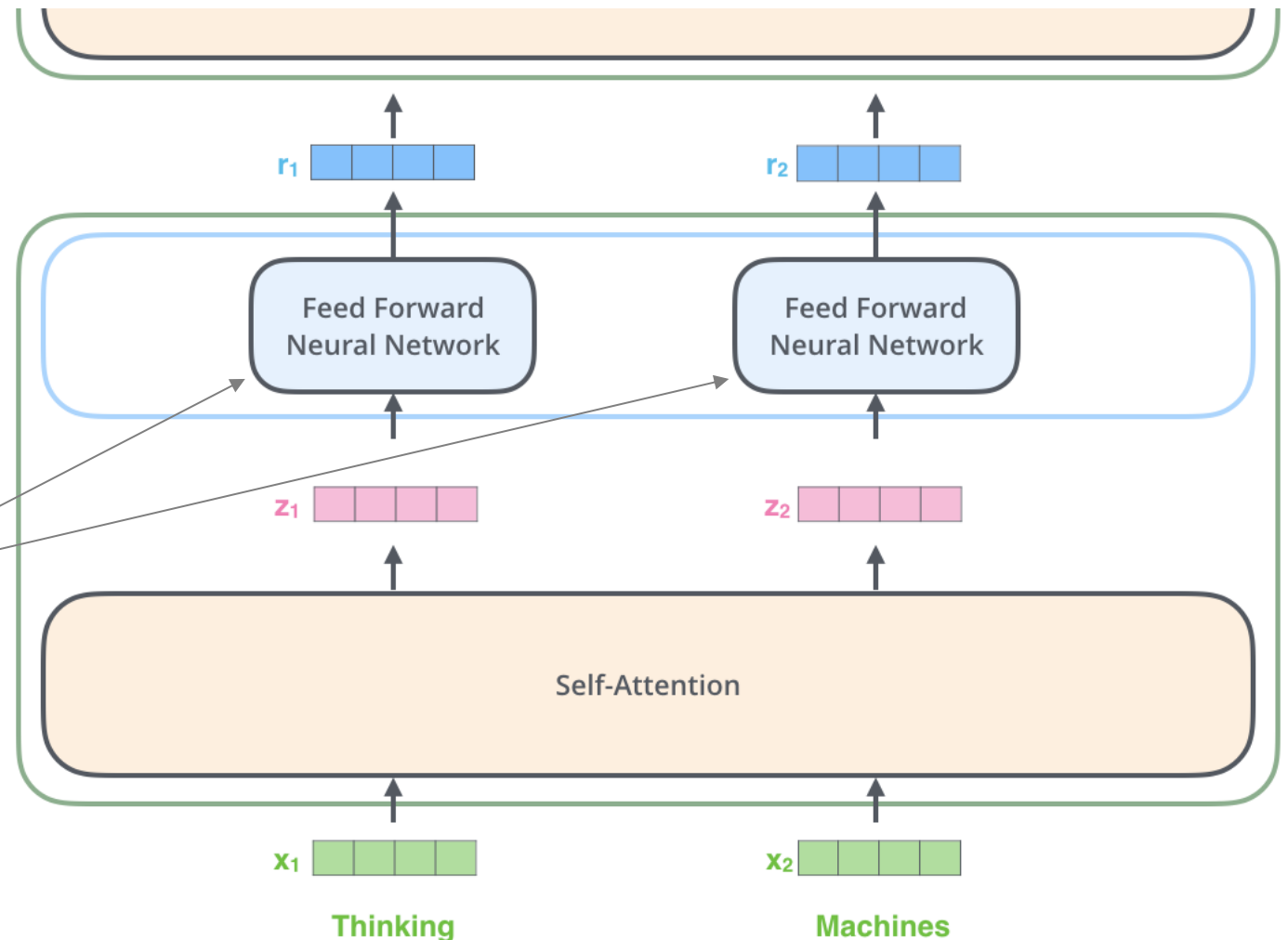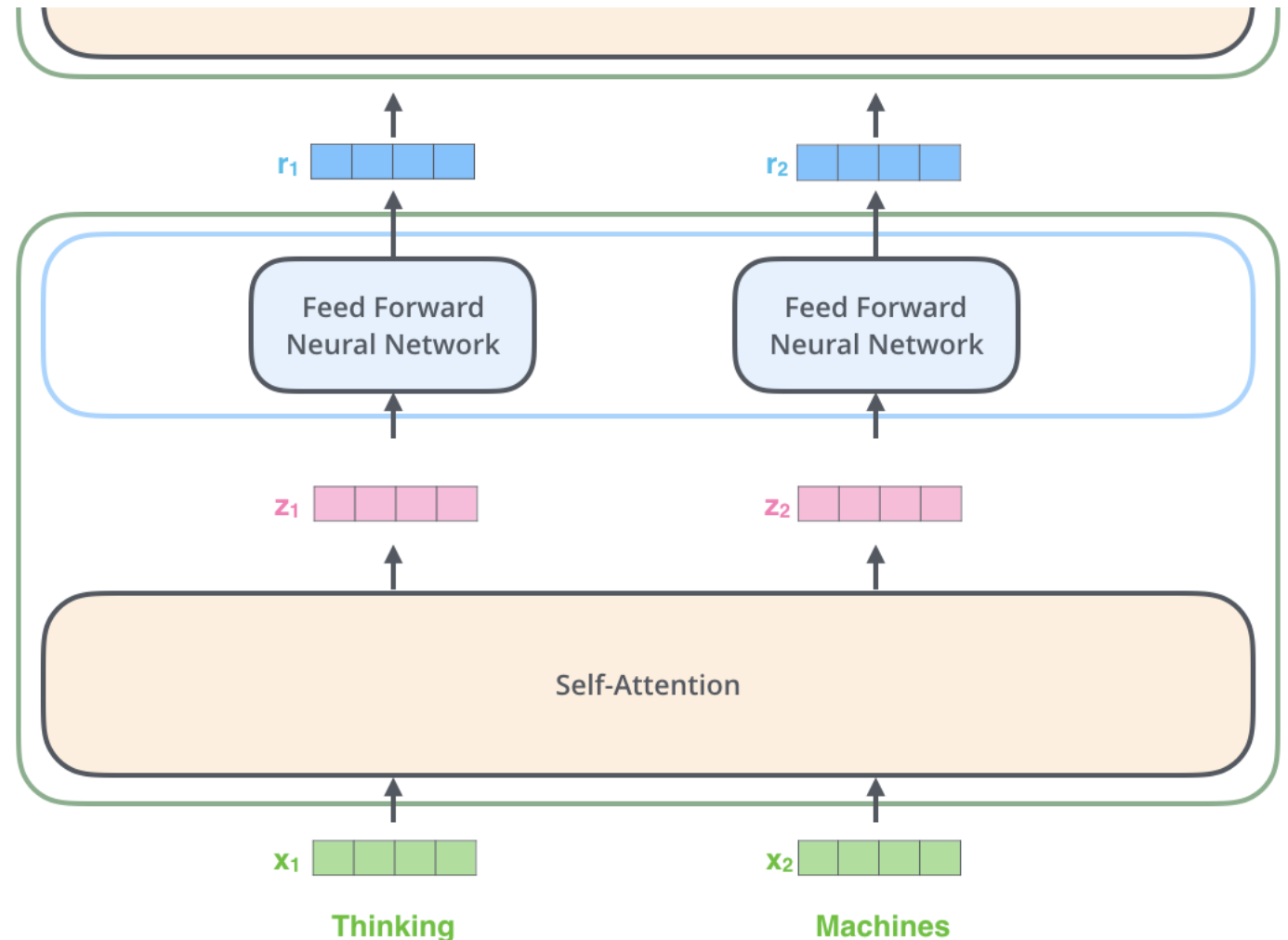Multiplying the resulting vector with **V** properly weighs the $v_i$ vectors.

40

# Encoder Block Map

- Encoder block breaks down into two main parts: Self-Attention, and Feed Forward layers.

- Self-Attention layer is applied to each word individually.

# Encoder Block Map

- Encoder block breaks down into two main parts: Self-Attention, and Feed Forward layers.

- Self-Attention layer is applied to each word individually.

- Feed Forward layer is applied to each word individually.

Alammar, Jay (2018). The Illustrated Transformer [Blog post]. Retrieved from https://jalammar.github.io/illustrated-transformer/

# Encoder Block Map

- Encoder block breaks down into two main parts: Self-Attention, and Feed Forward layers.

- Self-Attention layer is applied to each word individually.

- Feed Forward layer is applied to each word individually.

- The outputs of the feed forward layer are then passed as the inputs of the next encoder block.

# Encoder Block Map

- Encoder block breaks down into two main parts: Self-Attention, and Feed Forward layers.

- Self-Attention layer is applied to each word individually.

- Feed Forward layer is applied to each word individually.

- The outputs of the feed forward layer are then passed as the inputs of the next encoder block.
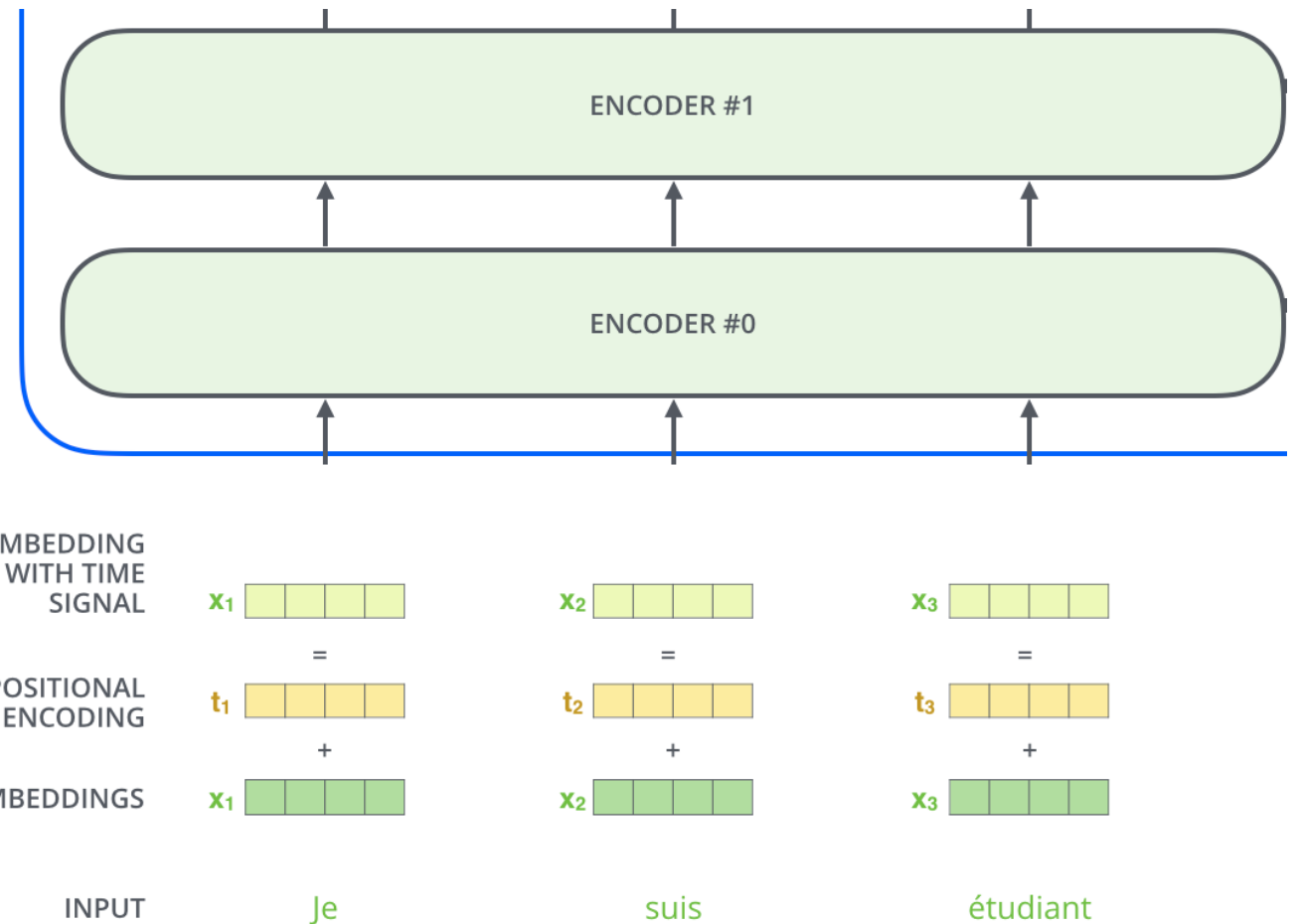
- *But we forgot about something...*

# What are we missing?

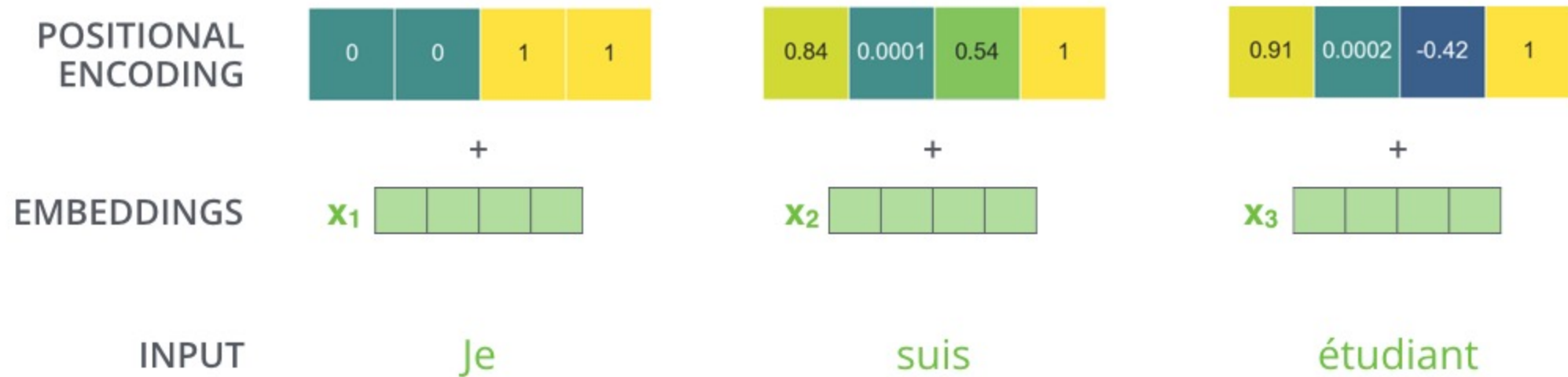*Hint: Remember – we are not using RNNs anymore.*

*Have we neglected/lost any information about the original input sequence?*
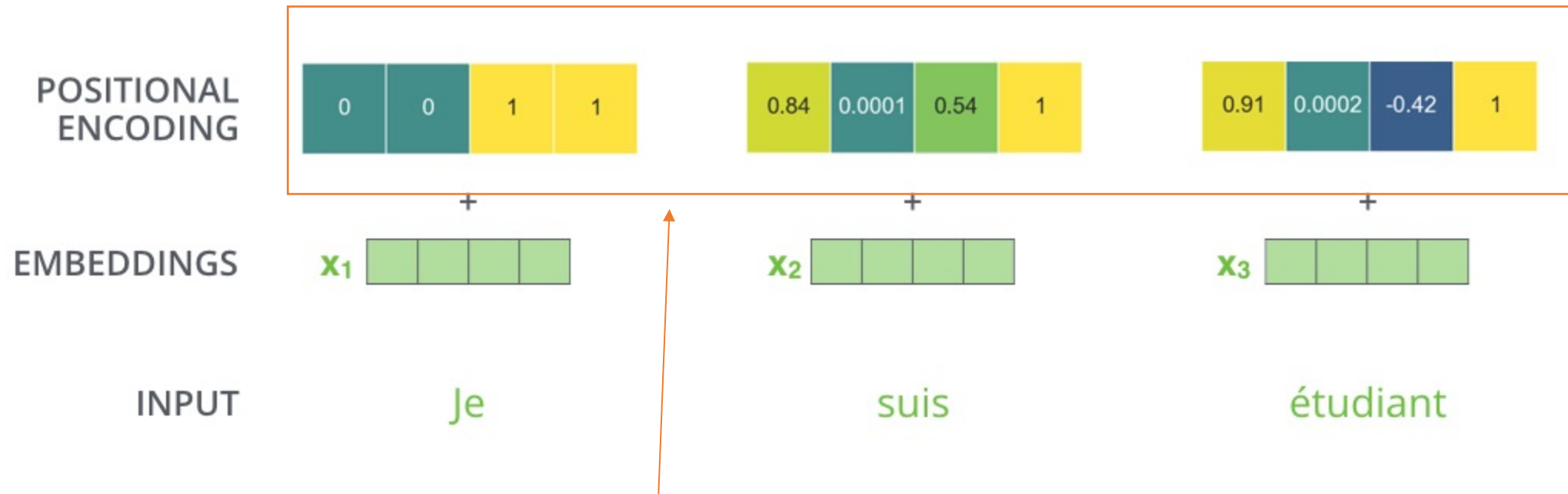
# Positional Encodings

- Instead of passing Embedding vector to encoder, we pass ***Embedding with Time Signal*** vector.

- Positional Encoding is `embedding_size` vector that encodes information about the position of a word in a sequence.

- Positional Encodings can be learned or defined by a fixed function.

- We add the Positional Encoding to the Embedding to get our Embedding with Time Signal vector.

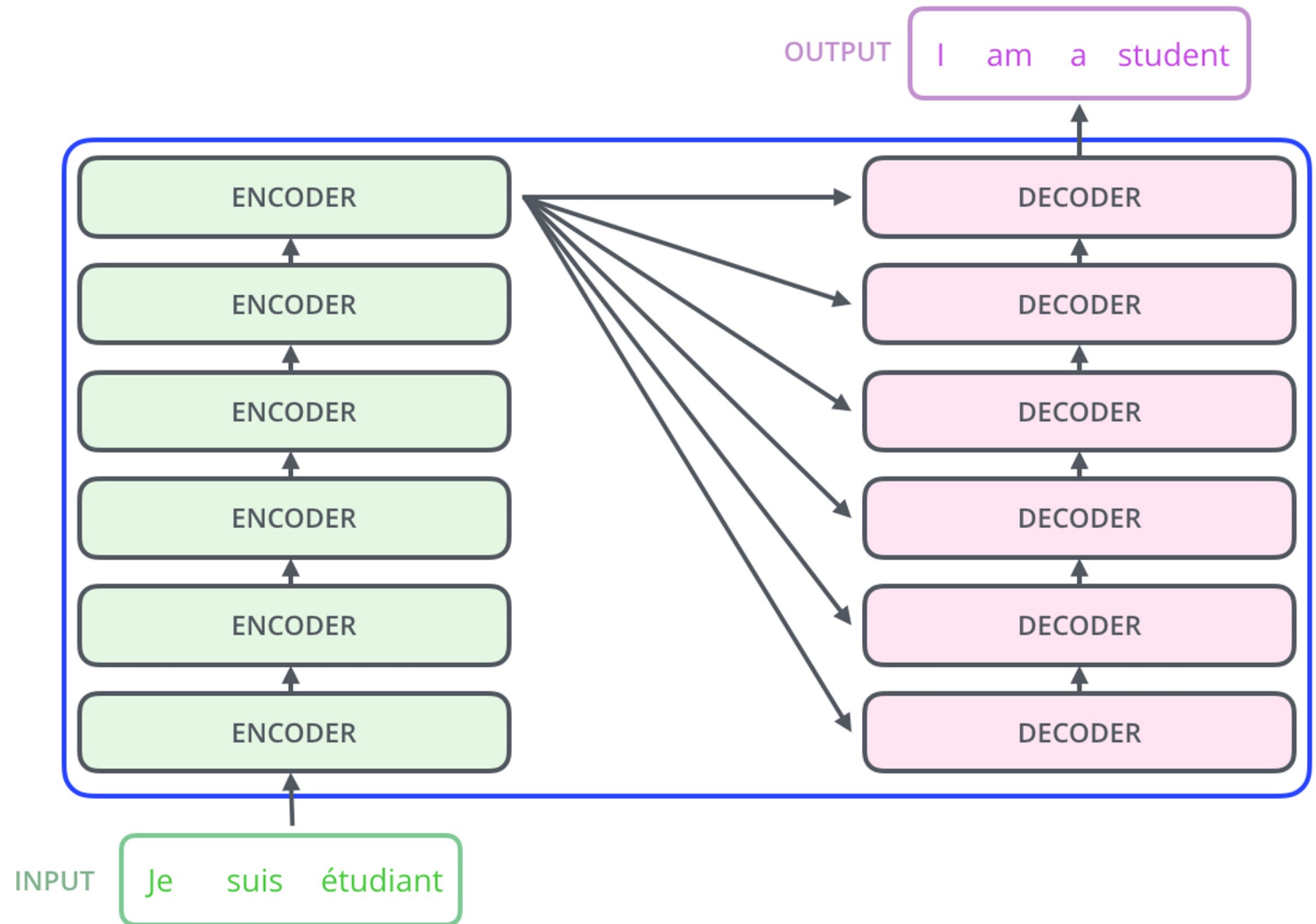# Positional Encodings

# Positional Encodings



**Where do these numbers come from?**
- Carefully-chosen sinusoidal patterns such that when we add them to the embedding vectors, their dot products w/ each other reflect the distance between them in the sentence.
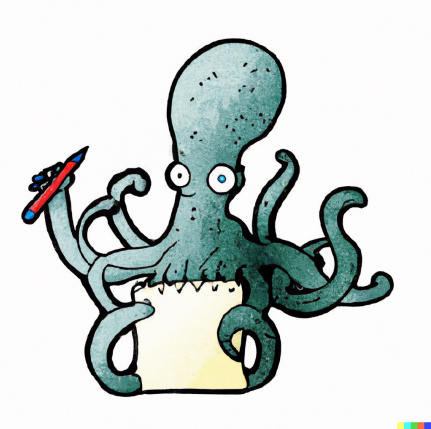
# More to come on Transformer!

- **Multi-headed attention**

- **Modifications for efficiency**

- **Decoder**

Alammar, Jay (2018). The Illustrated Transformer [Blog post]. Retrieved from https://jalammar.github.io/illustrated-transformer/

# Recap

Seq-to-seq using transformers

RNNs cannot be parallelized

Can forget information

Transformers – Encoder-Decoder with just attention

Encoder module

Self-attention

Fully connected layers

Positional encodings