# Transformer
## (Attention is All You Need)

# Transformer (A High-Level Look)
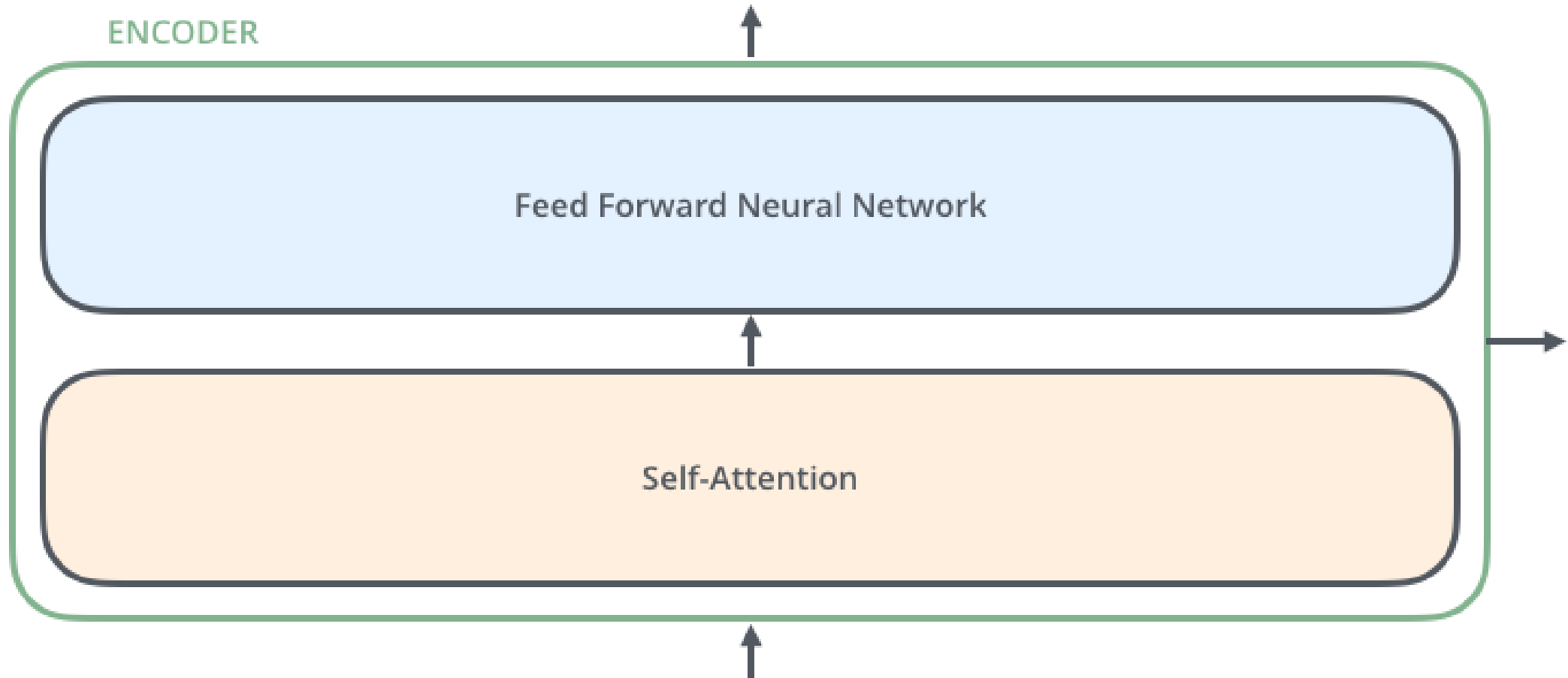
INPUT

Je   suis   étudiant

THE
TRANSFORMER

OUTPUT

I   am   a   student

# Transformer (A High-Level Look)

OUTPUT | I  am  a  student

ENCODERS → DECODERS

INPUT | Je  suis  étudiant

OUTPUT | I am a student

ENCODER

ENCODER

ENCODER

ENCODER

ENCODER

ENCODER

DECODER

DECODER

DECODER

DECODER

DECODER
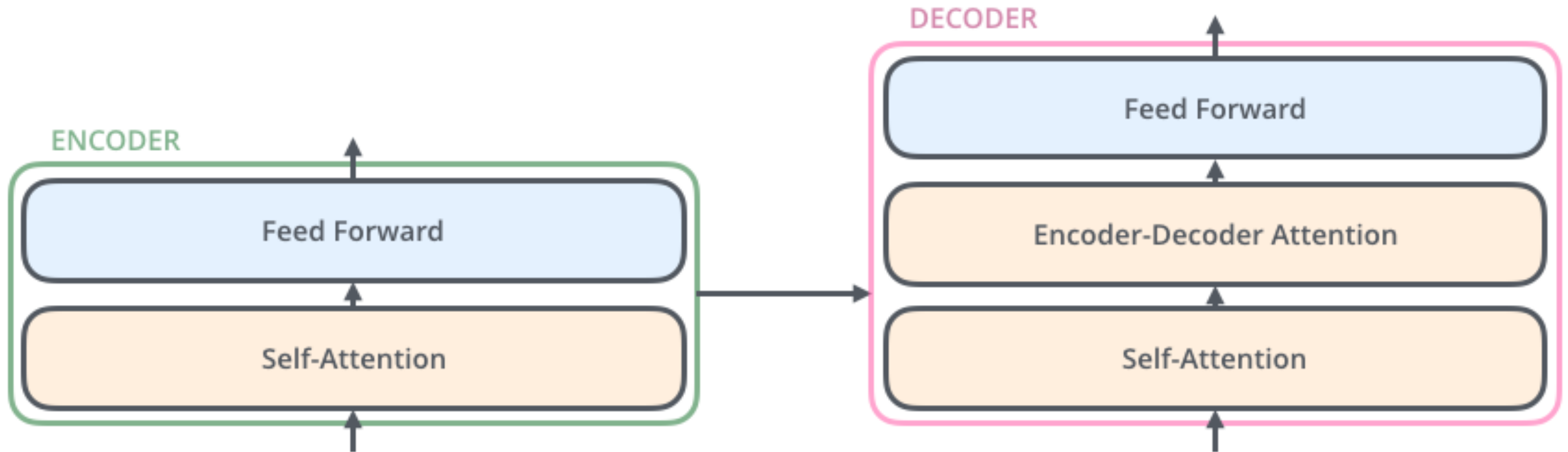
DECODER

INPUT | Je suis étudiant

The encoders are all identical in structure. Each one is broken down into two sub-layers:
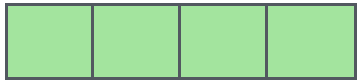
The encoders are all identical in structure. Each one is broken down into two sub-layers:
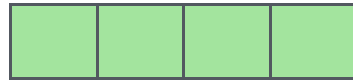
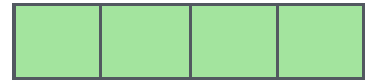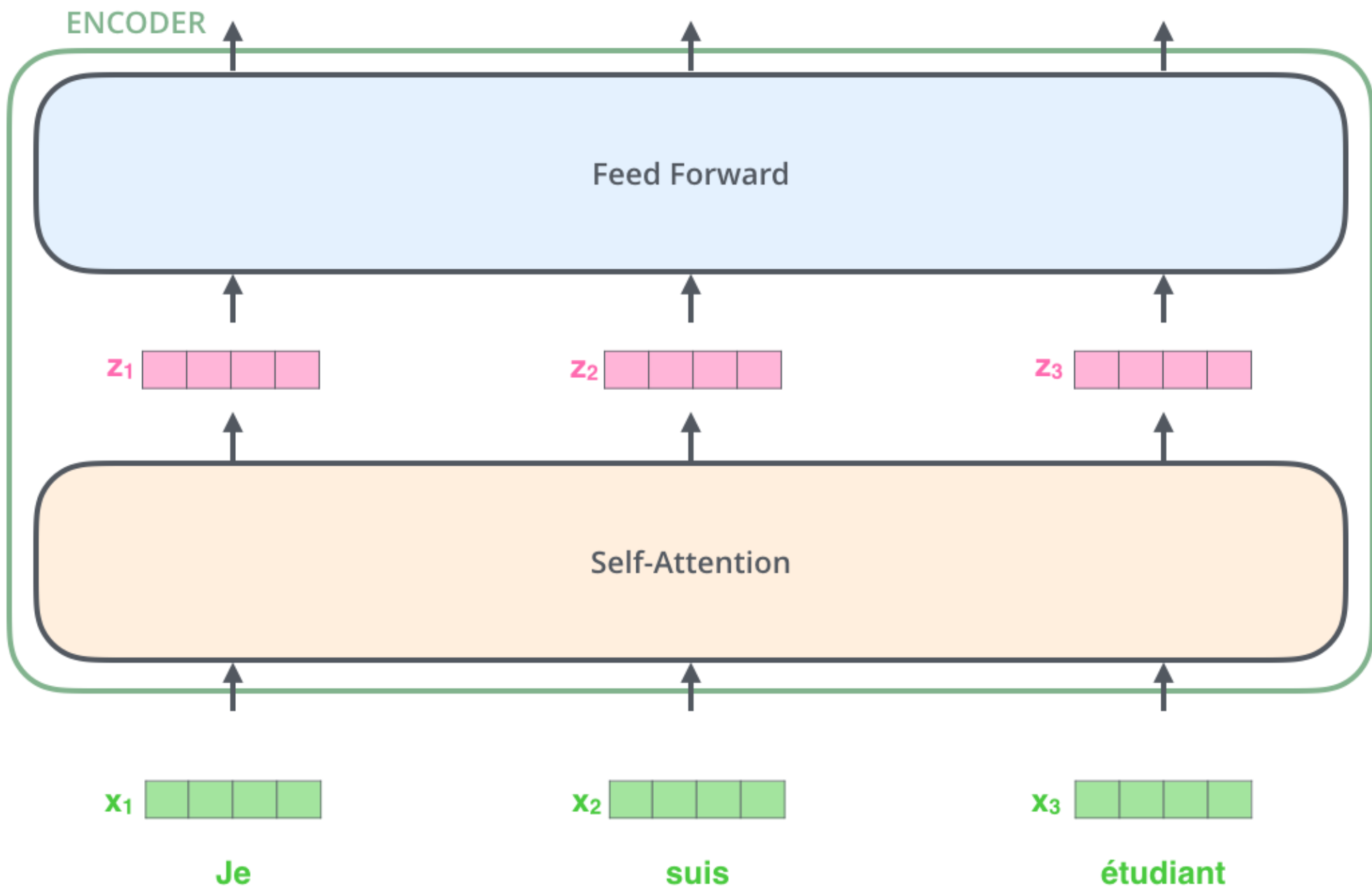Each word is embedded into a vector of size 512. We'll represent those vectors with these simple boxes.

$x_1$ [green box] $x_2$ [green box] $x_3$ [green box]

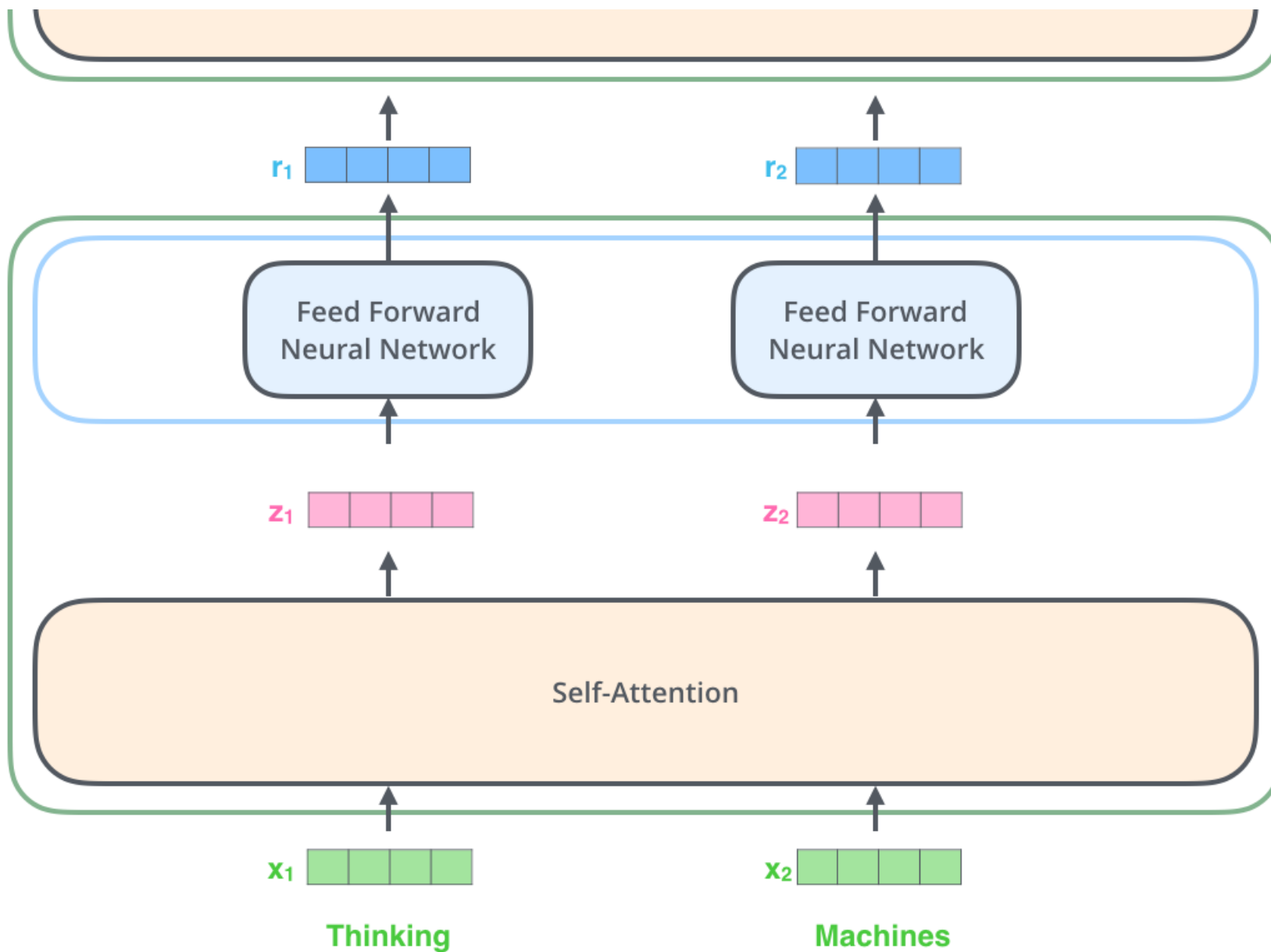Je       suis       étudiant

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

Layer: 5 ⬍ Attention: Input - Input ⬍

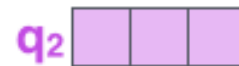| | |
|---|---|
| The_ | The_ |
| animal_ | animal_ |
| didn_ | didn_ |
| '_ | '_ |
| t_ | t_ |
| cross_ | cross_ |
| the_ | the_ |
| street_ | street_ |
| because_ | because_ |
| it_ | it_ |
| was_ | was_ |
| too_ | too_ |
| tire | tire |
| d_ | d_ |

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

| 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 |

| | 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 | $z_1$ | $z_2$ |

# The self-attention calculation in matrix form

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

$$= Z$$

# The Beast With Many Heads
## Multi Head Attention

- It expands the model's ability to focus on different positions

- It gives the attention layer multiple "representation subspaces"

**X**

Thinking
Machines

ATTENTION HEAD #0

ATTENTION HEAD #1

$Q_0$

$W_0^Q$

$Q_1$

$W_1^Q$

$K_0$

$W_0^K$

$K_1$

$W_1^K$

$V_0$

$W_0^V$

$V_1$

$W_1^V$

**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

R

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

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

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

...

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

$Q_0$  $K_0$  $V_0$

$Q_1$  $K_1$  $V_1$

...

$Q_7$  $K_7$  $V_7$

$Z_0$

$Z_1$

...

$Z_7$

$W^O$

Z

# Multi Head Attention

# Multi Head Attention

# Representing The Order of The Sequence Using Positional Encoding

# Representing The Order of The Sequence Using 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

# Positional Embedding

$$e'_w = e_w + \left[ sin\left(\frac{pos}{10000^0}\right), cos\left(\frac{pos}{10000^0}\right), sin\left(\frac{pos}{10000^{2/4}}\right), cos\left(\frac{pos}{10000^{2/4}}\right) \right]$$

$$= e_w + \left[ sin(pos), cos(pos), sin\left(\frac{pos}{100}\right), cos\left(\frac{pos}{100}\right) \right]$$

$$PE(pos, 2i) = sin\left(\frac{pos}{10000^{2i/d_{model}}}\right),$$

$$PE(pos, 2i+1) = cos\left(\frac{pos}{10000^{2i/d_{model}}}\right).$$

# Positional Embedding

# Positional Embedding

# The Decoder Side

# Attention in Decoder

- The self attention layers in the decoder operate in a slightly different way than the one in the encoder

- In the decoder, the self-attention layer is only allowed to attend to earlier positions in the output sequence. This is done by masking future positions (setting them to -inf) before the softmax step in the self-attention calculation

- The "Encoder-Decoder Attention" layer works just like multiheaded self-attention, except it creates its Queries matrix from the layer below it, and takes the Keys and Values matrix from the output of the encoder stack

# The 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

# Recap Of Training

Output Vocabulary

| WORD | a | am | I | thanks | student | <eos> |
|------|---|----|----|--------|---------|-------|
| INDEX | 0 | 1 | 2 | 3 | 4 | 5 |

One-hot encoding of the word "am"

| 0.0 | 1.0 | 0.0 | 0.0 | 0.0 | 0.0 |

## Untrained Model Output

| 0.2 | 0.2 | 0.1 | 0.2 | 0.2 | 0.1 |
|-----|-----|-----|-----|-----|-----|

## Correct and desired output

| 0.0 | 0.0 | 0.0 | 1.0 | 0.0 | 0.0 |
|-----|-----|-----|-----|-----|-----|
| a | am | I | thanks | student | <eos> |

# The Loss Function

- Each probability distribution is represented by a vector of width vocab_size (6 in our toy example, but more realistically a number like 3,000 or 10,000)
- The first probability distribution has the highest probability at the cell associated with the word "i"
- The second probability distribution has the highest probability at the cell associated with the word "am"
- And so on, until the fifth output distribution indicates

# Trained Model Outputs

Beam Search

Output Vocabulary: a   am   I   thanks   student   <eos>

| | a | am | I | thanks | student | <eos> |
|---|---|---|---|---|---|---|
| position #1 | 0.01 | 0.02 | 0.93 | 0.01 | 0.03 | 0.01 |
| position #2 | 0.01 | 0.8 | 0.1 | 0.05 | 0.01 | 0.03 |
| position #3 | 0.99 | 0.001 | 0.001 | 0.001 | 0.002 | 0.001 |
| position #4 | 0.001 | 0.002 | 0.001 | 0.02 | 0.94 | 0.01 |
| position #5 | 0.01 | 0.01 | 0.001 | 0.001 | 0.001 | 0.98 |

a   am   I   thanks   student   <eos>

# References

- [Jay Alammar](https://jalammar.github.io/illustrated-transformer/)  https://jalammar.github.io/illustrated-transformer/

- Ashish Vaswani, Noam Shazeer, Niki Parmar, Jakob Uszkoreit, Llion Jones, Aidan N. Gomez, Łukasz Kaiser, and Illia Polosukhin. 2017. **Attention is all you need**. In Proceedings of *the 31st International Conference on Neural Information Processing Systems* (NIPS'17),