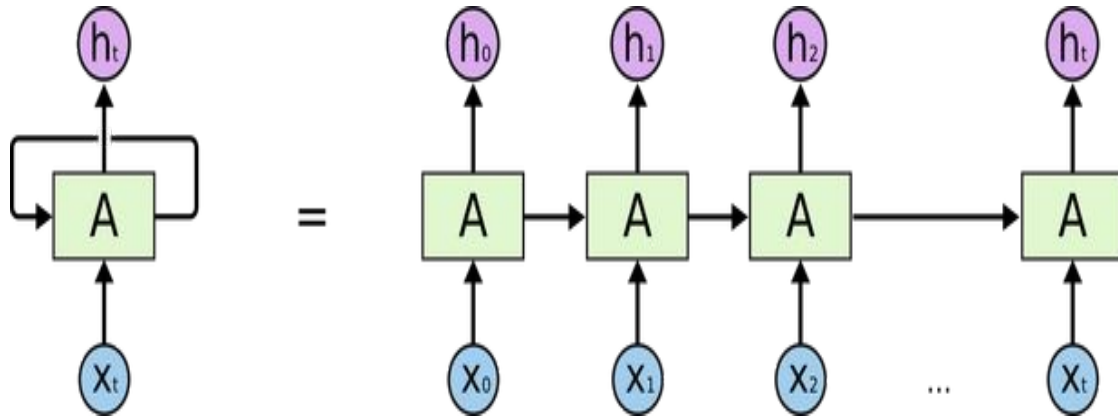# NLP Session

# Topics

- RNN
- LSTM
- Transformers

# RNN
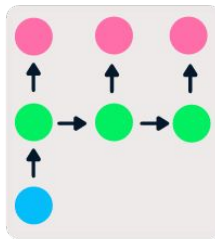
# Types of RNN

- One to One RNN
- One to Many RNN
- Many to One RNN
- Many to Many RNN
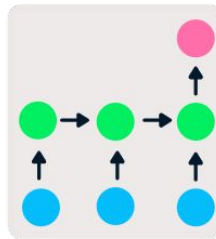
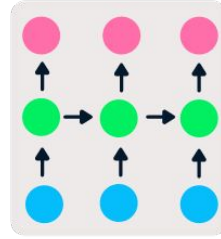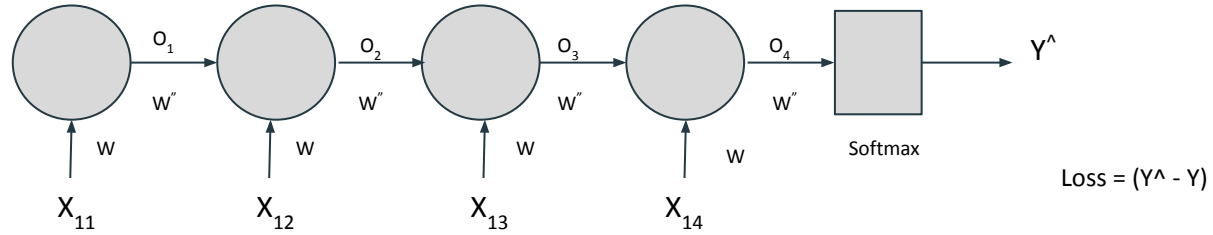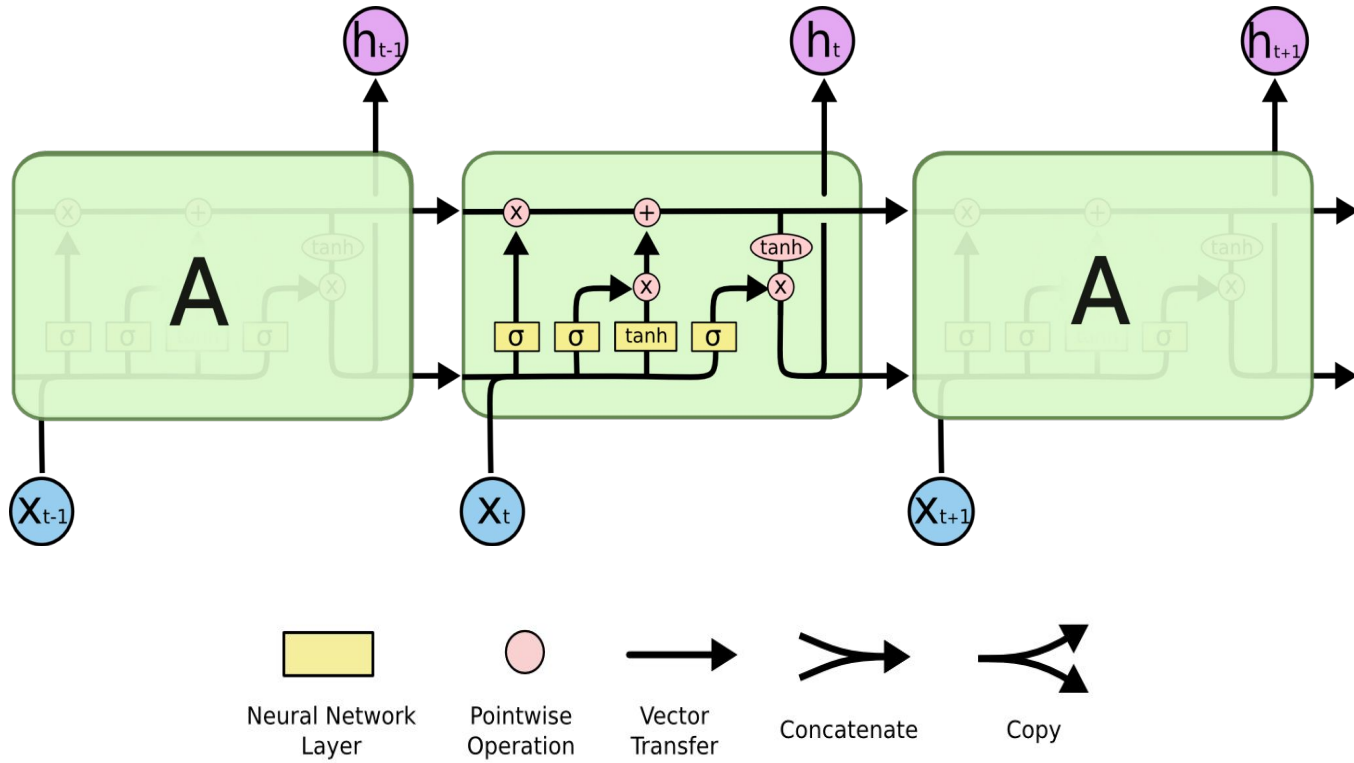# Forward and Backward Propagation
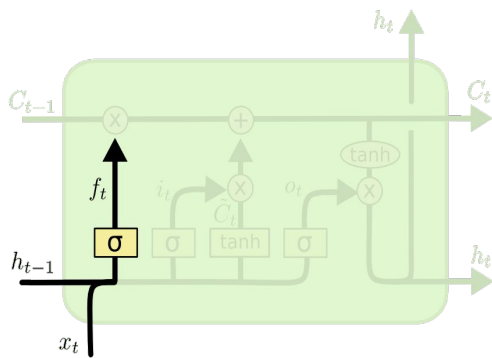
# LSTM

# Forget Gate Layer

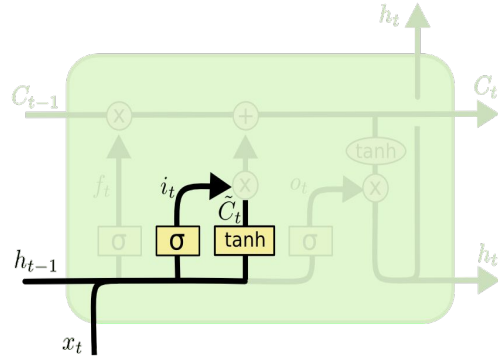Example 1 - Rohan likes pizza but he does not like burger.
Example 2 - Rohan likes pizza but his friend likes burger.



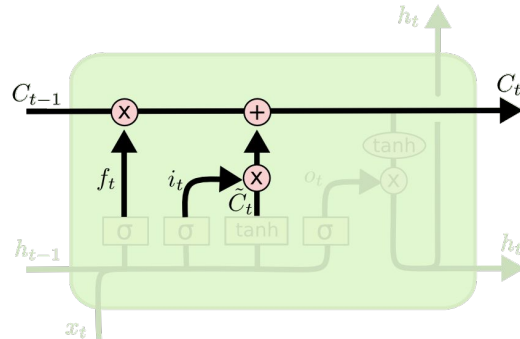$$f_t = \sigma \left( W_f \cdot [h_{t-1}, x_t] \ + \ b_f \right)$$

# Input Gate Layer



$$i_t = \sigma \left( W_i \cdot [h_{t-1}, x_t] + b_i \right)$$

$$\tilde{C}_t = \tanh(W_C \cdot [h_{t-1}, x_t] + b_C)$$
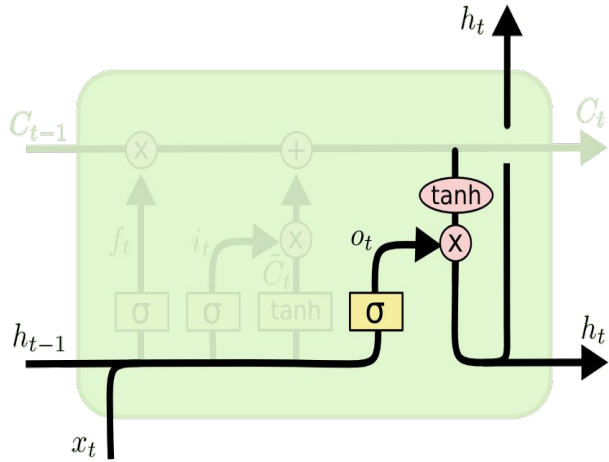
Combining
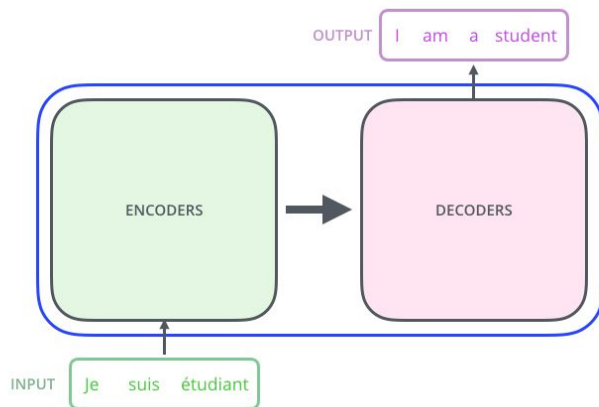


$$C_t = f_t * C_{t-1} + i_t * \tilde{C}_t$$
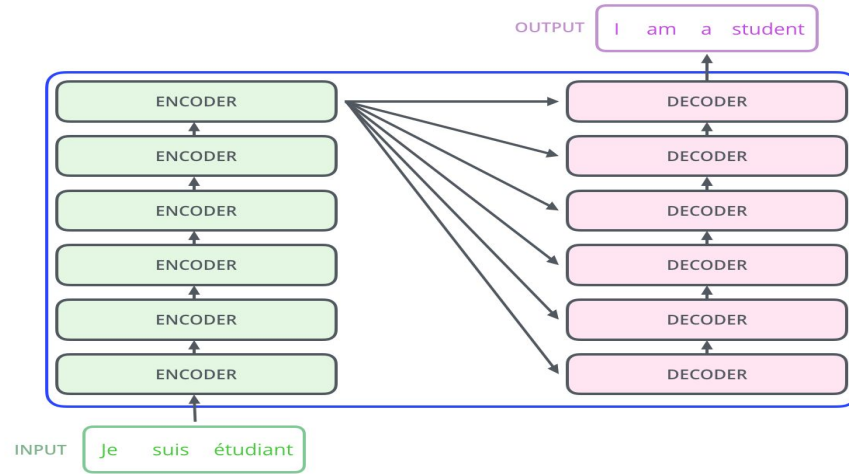
# Output Gate Layer



$$o_t = \sigma \left( W_o \left[ h_{t-1}, x_t \right] + b_o \right)$$

$$h_t = o_t * \tanh \left( C_t \right)$$

# Transformers

INPUT

Je    suis    étudiant

THE
TRANSFORMER

OUTPUT

I    am    a    student

OUTPUT    I    am    a    student

ENCODERS → DECODERS

INPUT    Je    suis    étudiant

OUTPUT: I am a student

ENCODER → DECODER (×6)

INPUT: Je suis étudiant

In Every Encoder -



ENCODER

Feed Forward Neural Network

Self-Attention
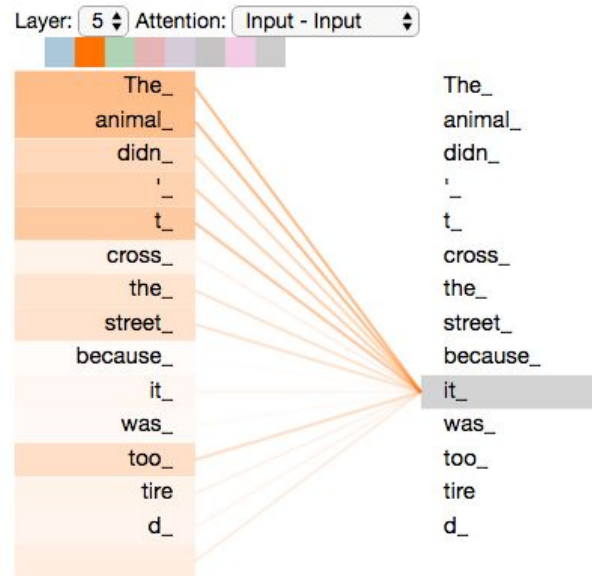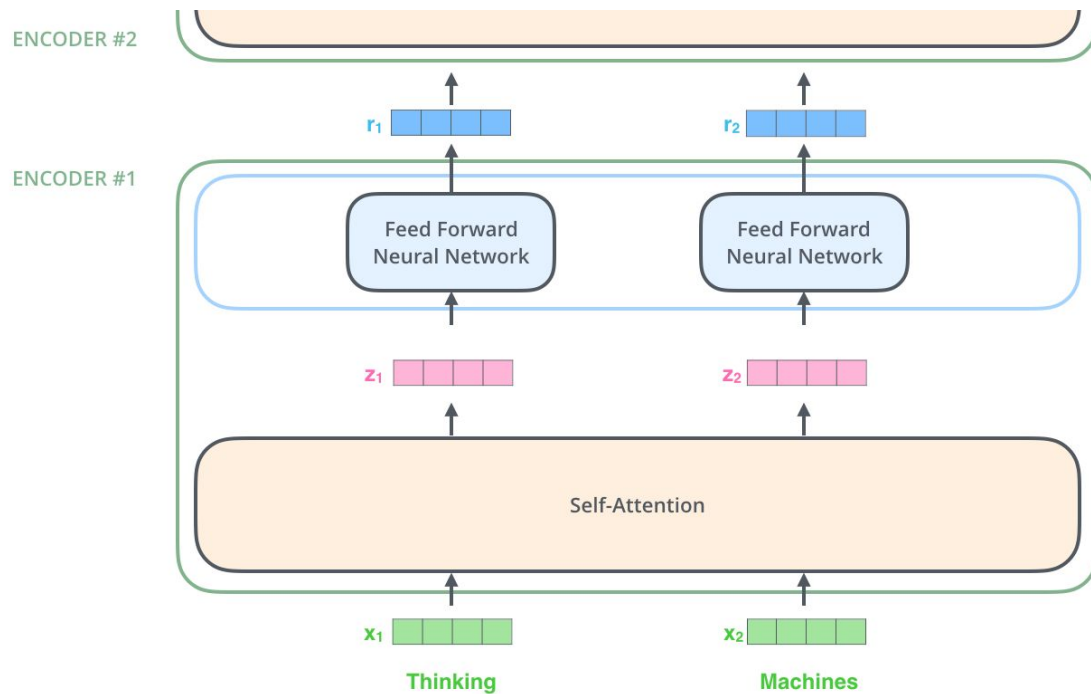
# Self - Attention

Why to use ?
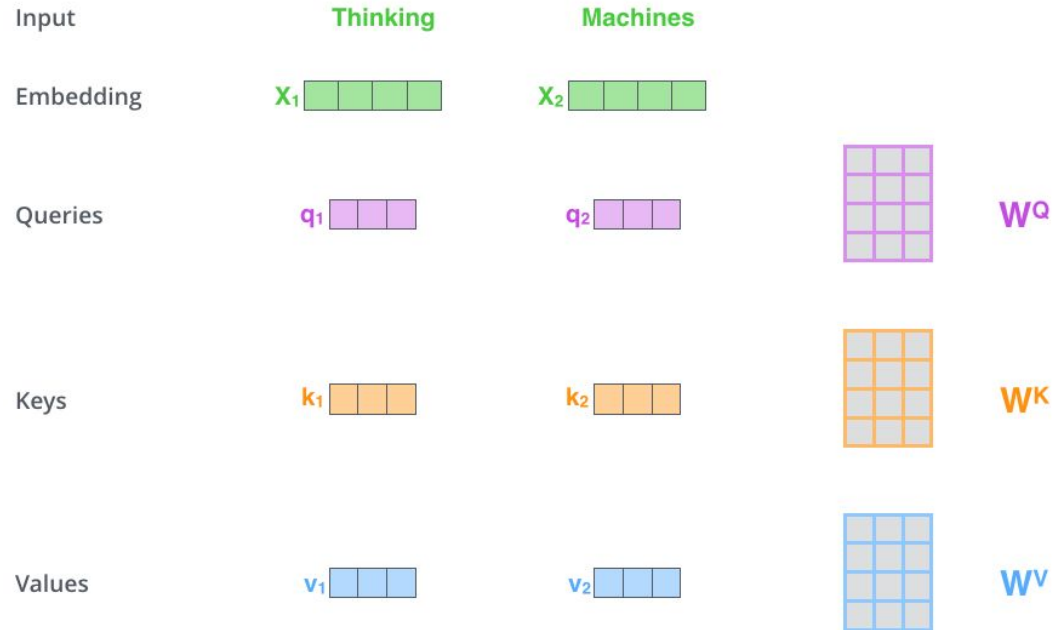
Example - The *animal* didn't cross the street because it was too tired

# Understanding self-attention

Step 1 - To create Queries, key and values vector

Step 2 - To calculate the score

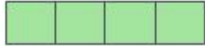# Step 3 - Divide the score by 8 and calculating the softmax

| | 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 |

# Step 4 - Multiply the value vector with 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 |
| Softmax X Value | $v_1$ | $v_2$ |
| Sum | $z_1$ | $z_2$ |

# Matrix Calculation of self-attention

# Multi-headed Attention

# Input to feed-forward neural network

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

=

# After Applying multi-headed attention

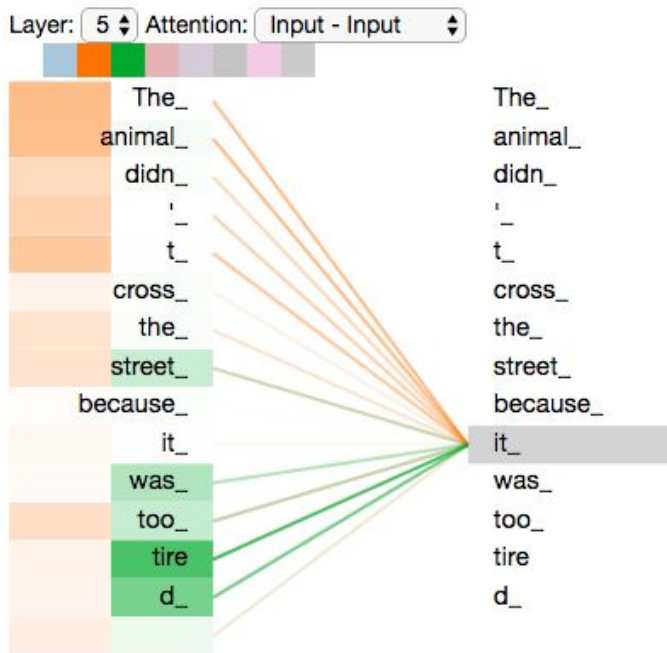Example - The <span style="color:orange">animal</span> didn't cross the street because <span style="color:red">it</span> was too <span style="color:green">tired</span>
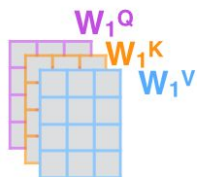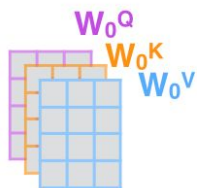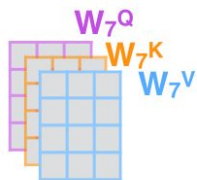
# Entire Process

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

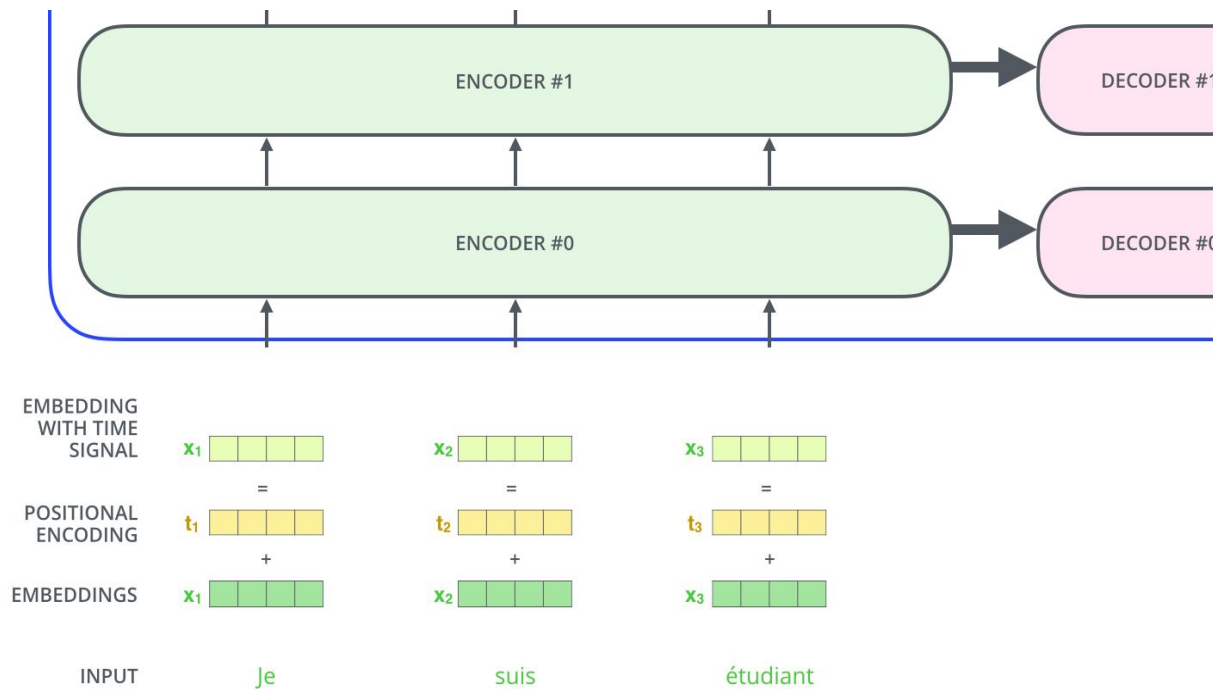5) Concatenate the resulting $Z$ matrices, then multiply with weight matrix $W^O$ to produce the output of the layer
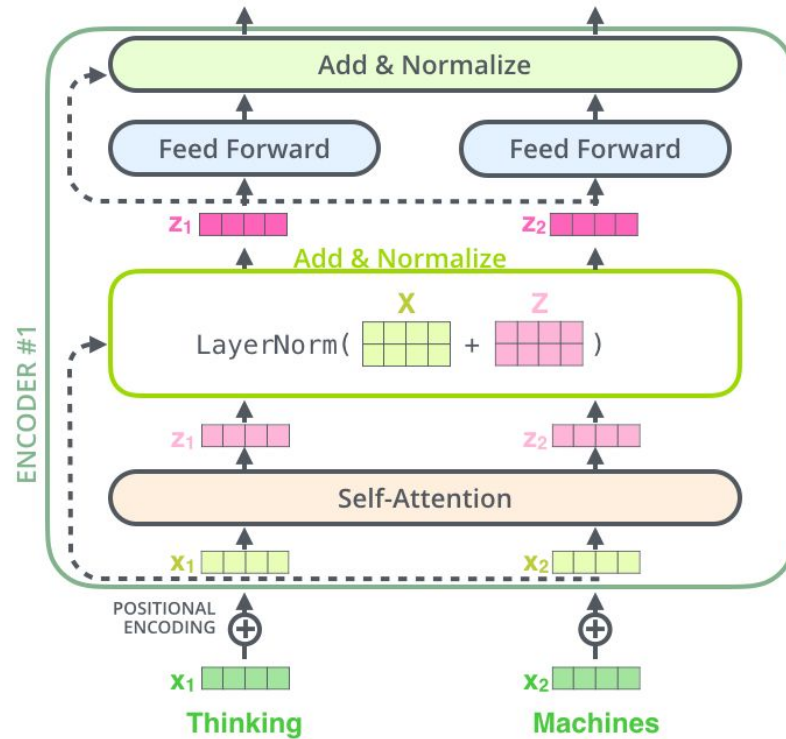


Thinking Machines

$X$

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

$Q_0$
$K_0$
$V_0$

$Z_0$

$W^O$

* In all encoders other than #0, we don't need embedding. We start directly with the output of the encoder right below this one
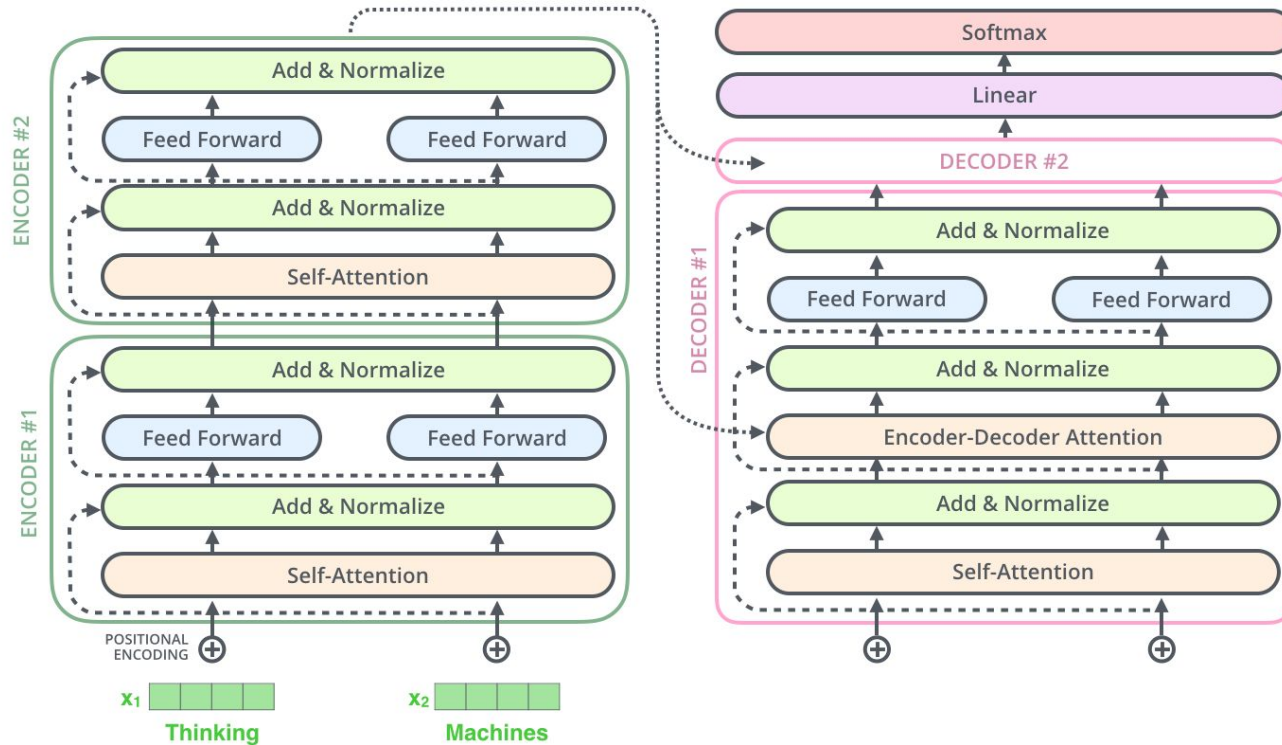
$W_1^Q$
$W_1^K$
$W_1^V$

$Q_1$
$K_1$
$V_1$

$Z_1$

$Z$

...

...

...

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

$Q_7$
$K_7$
$V_7$

$Z_7$

# Positional Encodings

# Layer Normalization

# Entire Process

# Output Visualization

Thank you