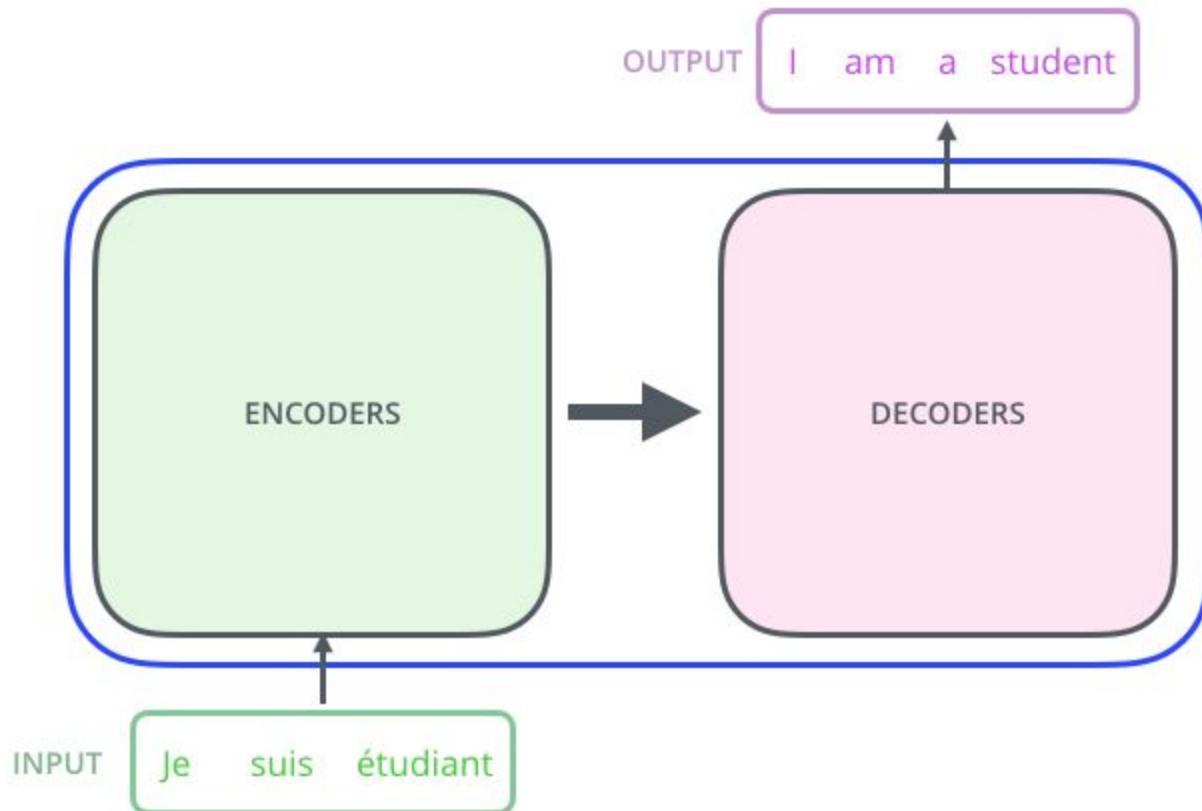
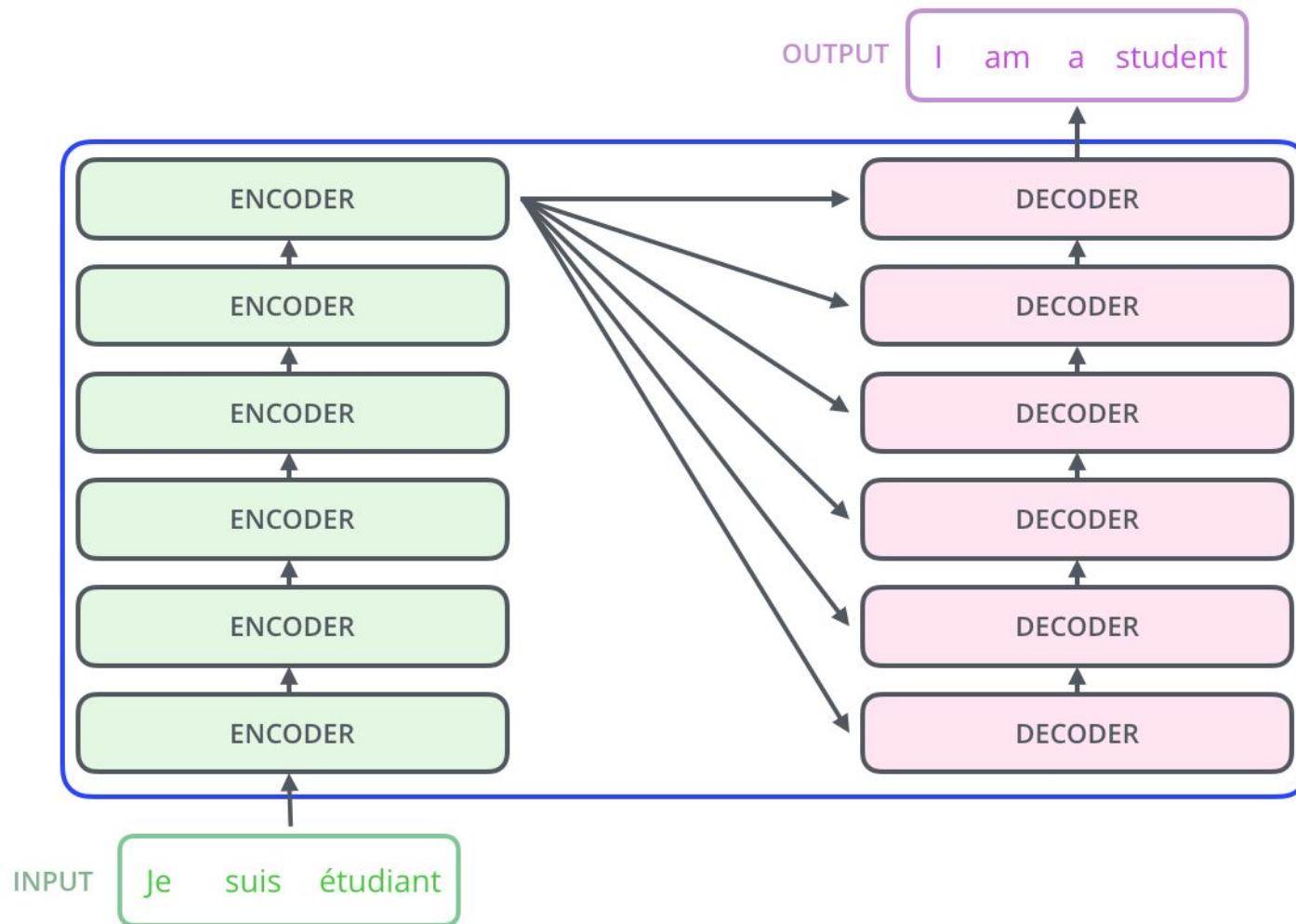
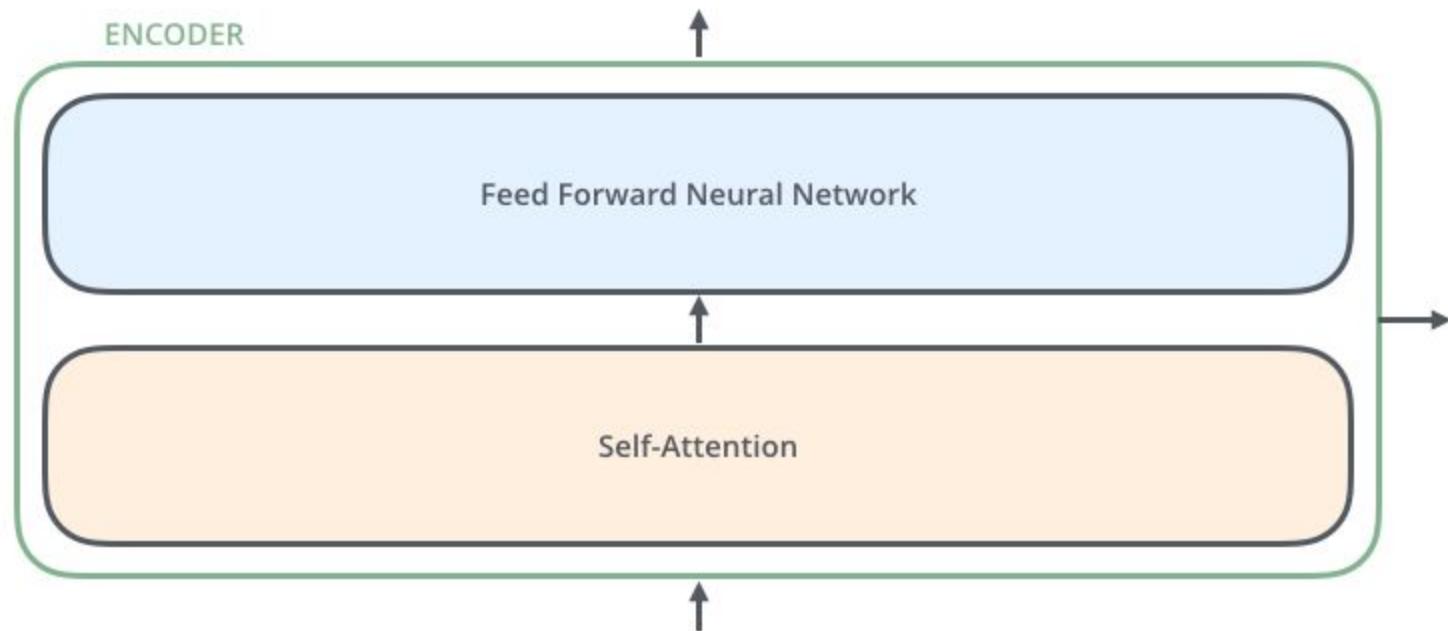
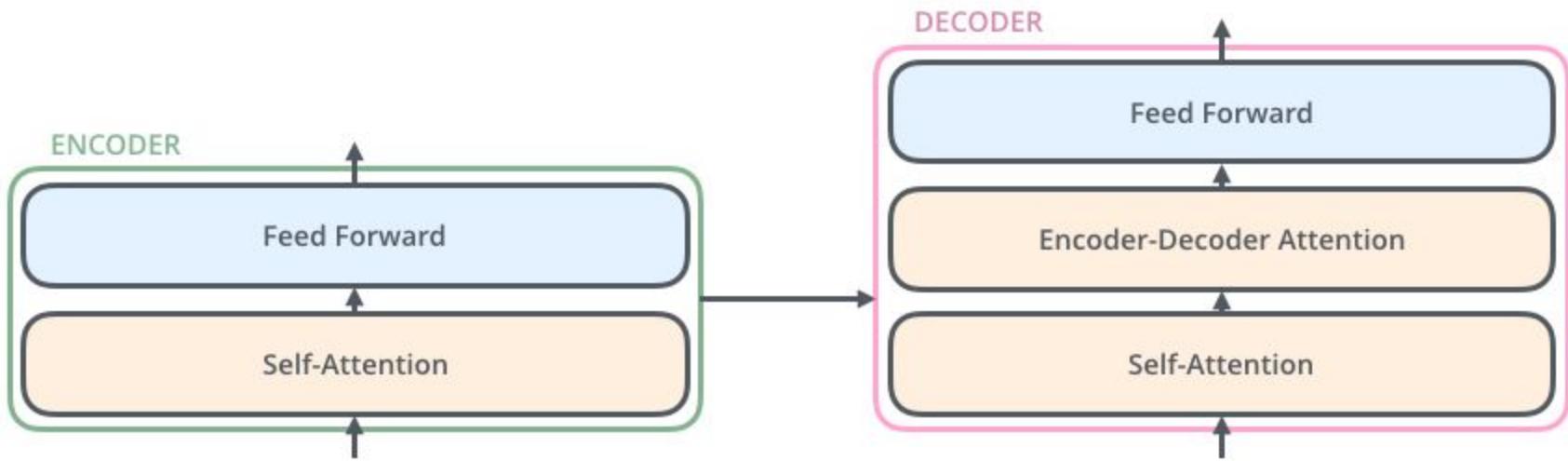

Transformer

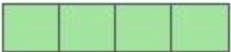




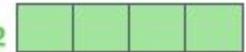




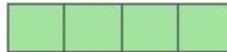


x_1 

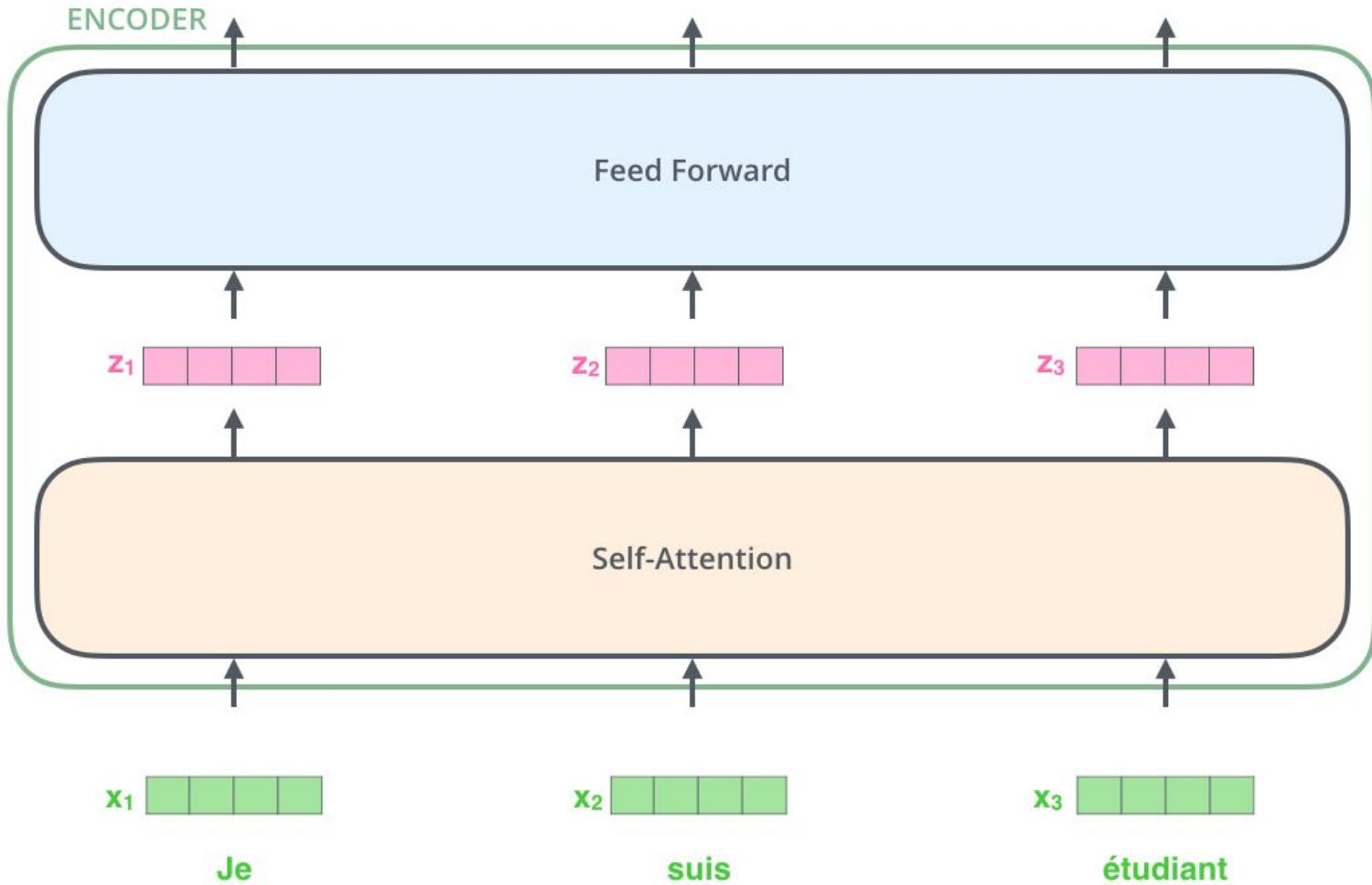
Je

x_2 

suis

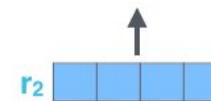
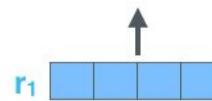
x_3 

étudiant



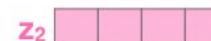
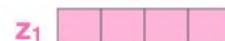
ENCODER #2

ENCODER #1

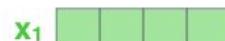


Feed Forward
Neural Network

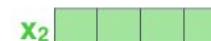
Feed Forward
Neural Network



Self-Attention



Thinking



Machines

Self Attention

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

Layer: 5 Attention: Input - Input



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

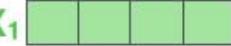
The_
animal_
didn_
'
t_
cross_
the_
street_
because_
it_
was_
too_
tire
d_

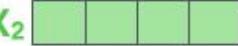
Input

Thinking

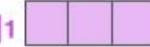
Machines

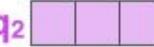
Embedding

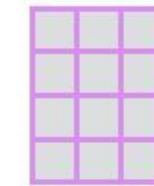
X_1 

X_2 

Queries

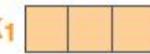
q_1 

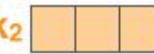
q_2 



W^Q

Keys

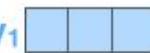
k_1 

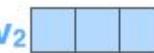
k_2 



W^K

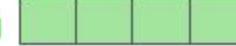
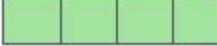
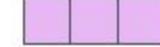
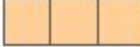
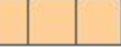
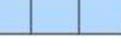
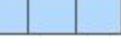
Values

v_1 

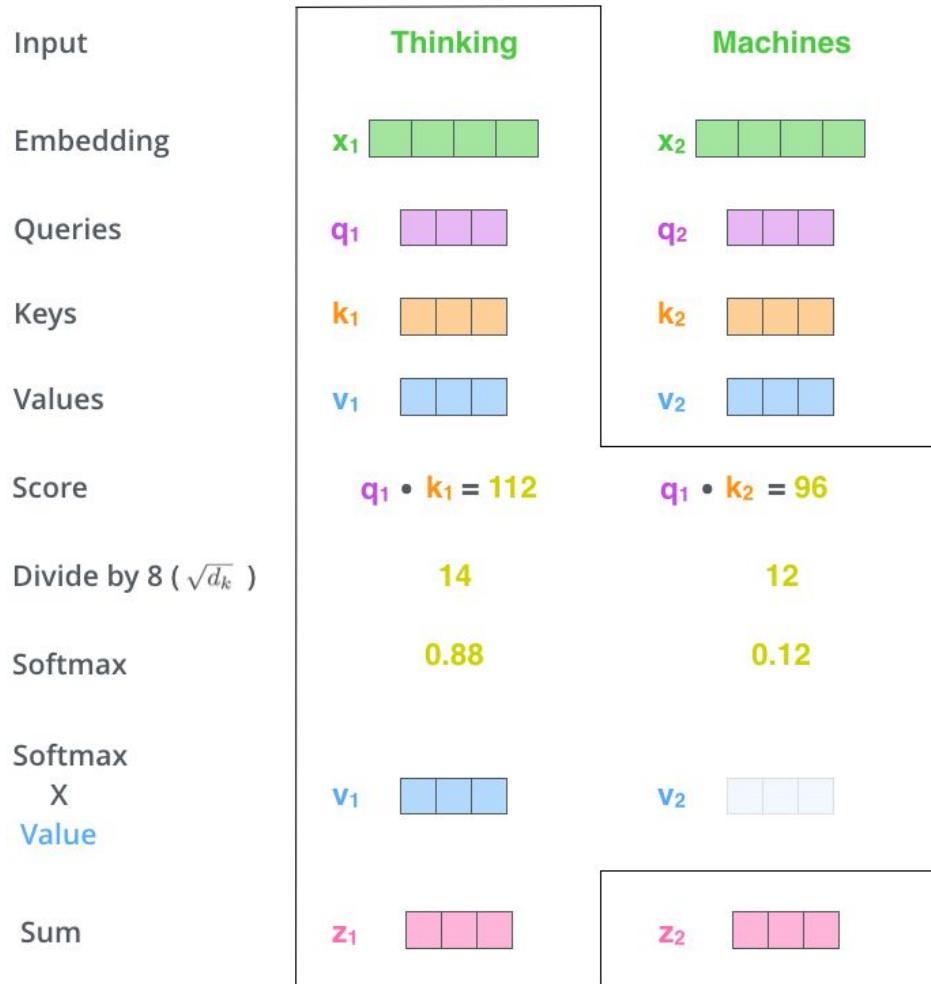
v_2 



W^V

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		
Embedding	x_1	
Queries	q_1	
Keys	k_1	
Values	v_1	
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



$$\begin{matrix} \textcolor{green}{X} \\ \\ \begin{array}{|c|c|c|c|} \hline & & & \\ \hline & & & \\ \hline & & & \\ \hline \end{array} \end{matrix} \times \begin{matrix} \textcolor{violet}{W^Q} \\ \\ \begin{array}{|c|c|c|c|c|c|} \hline & & & & & \\ \hline \end{array} \end{matrix} = \begin{matrix} \textcolor{violet}{Q} \\ \\ \begin{array}{|c|c|c|} \hline & & \\ \hline & & \\ \hline & & \\ \hline \end{array} \end{matrix}$$

$$\begin{matrix} \textcolor{green}{X} \\ \\ \begin{array}{|c|c|c|c|} \hline & & & \\ \hline & & & \\ \hline & & & \\ \hline \end{array} \end{matrix} \times \begin{matrix} \textcolor{orange}{W^K} \\ \\ \begin{array}{|c|c|c|c|c|c|} \hline & & & & & \\ \hline \end{array} \end{matrix} = \begin{matrix} \textcolor{orange}{K} \\ \\ \begin{array}{|c|c|c|} \hline & & \\ \hline & & \\ \hline & & \\ \hline \end{array} \end{matrix}$$

$$\begin{matrix} \textcolor{green}{X} \\ \\ \begin{array}{|c|c|c|c|} \hline & & & \\ \hline & & & \\ \hline & & & \\ \hline \end{array} \end{matrix} \times \begin{matrix} \textcolor{blue}{W^V} \\ \\ \begin{array}{|c|c|c|c|c|c|} \hline & & & & & \\ \hline \end{array} \end{matrix} = \begin{matrix} \textcolor{blue}{V} \\ \\ \begin{array}{|c|c|c|} \hline & & \\ \hline & & \\ \hline & & \\ \hline \end{array} \end{matrix}$$

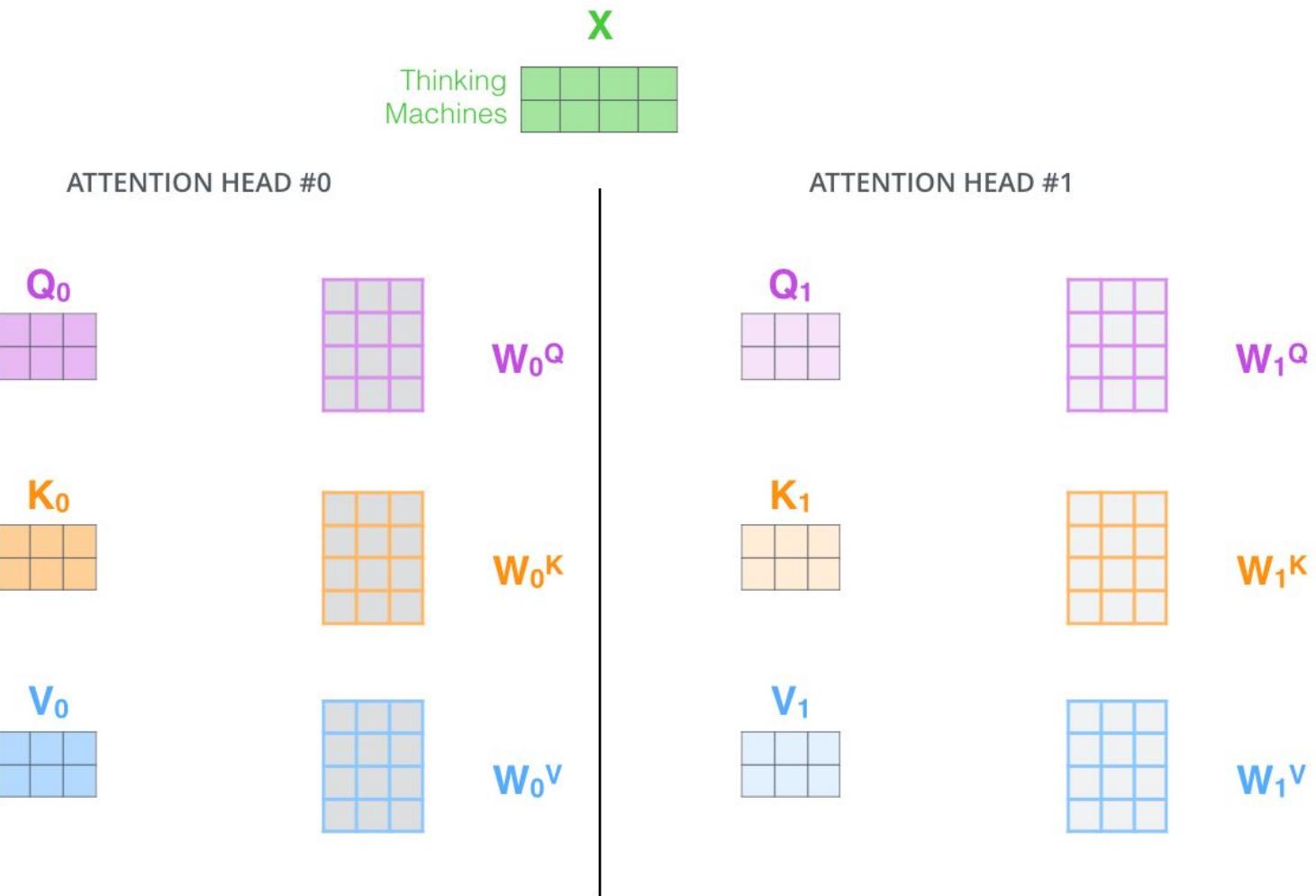
$$\text{softmax} \left(\frac{\mathbf{Q} \times \mathbf{K}^T}{\sqrt{d_k}} \right) \mathbf{V}$$

The diagram illustrates the computation of the query matrix \mathbf{Q} multiplied by the transpose of the key matrix \mathbf{K}^T , scaled by the square root of the dimension d_k . The result is then multiplied by the value matrix \mathbf{V} .

$$= \mathbf{Z}$$

The resulting matrix \mathbf{Z} is shown as a 2x4 grid of pink squares.

Multi-headed Attention



X

Thinking
Machines

Calculating attention separately in
eight different attention heads

ATTENTION
HEAD #0

Z_0

ATTENTION
HEAD #1

Z_1

...

ATTENTION
HEAD #7

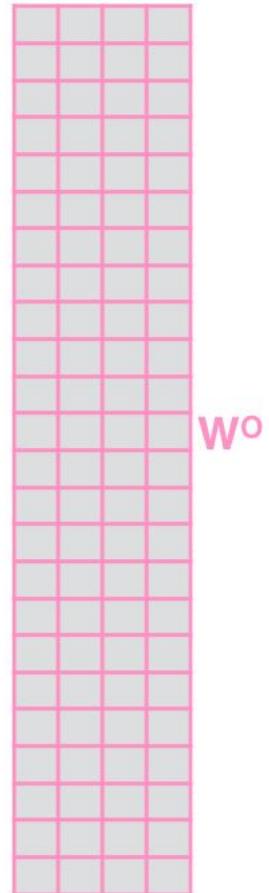
Z_7

1) Concatenate all the attention heads



2) Multiply with a weight matrix W^o that was trained jointly with the model

X



3) The result would be the Z matrix that captures information from all the attention heads. We can send this forward to the FFNN

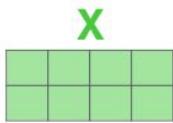
$$= \begin{matrix} Z \\ \hline \end{matrix}$$

Multi Headed Self Attention

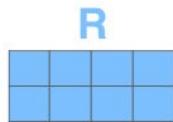
1) This is our input sentence*

2) We embed each word*

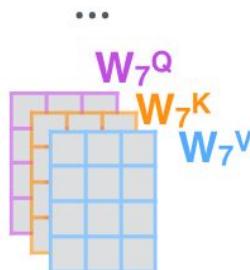
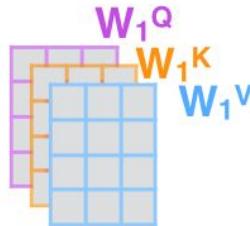
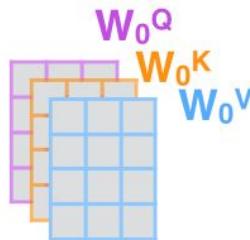
Thinking Machines



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



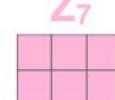
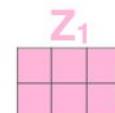
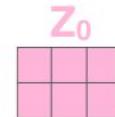
3) Split into 8 heads.
We multiply X or R with weight matrices



4) Calculate attention using the resulting $Q/K/V$ matrices



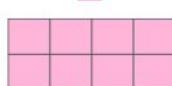
5) Concatenate the resulting Z matrices, then multiply with weight matrix W^o to produce the output of the layer



W^o



Z



Layer: 5

Attention: Input - Input



The_
animal_
didn_
'
t_
cross_
the_
street_
because_

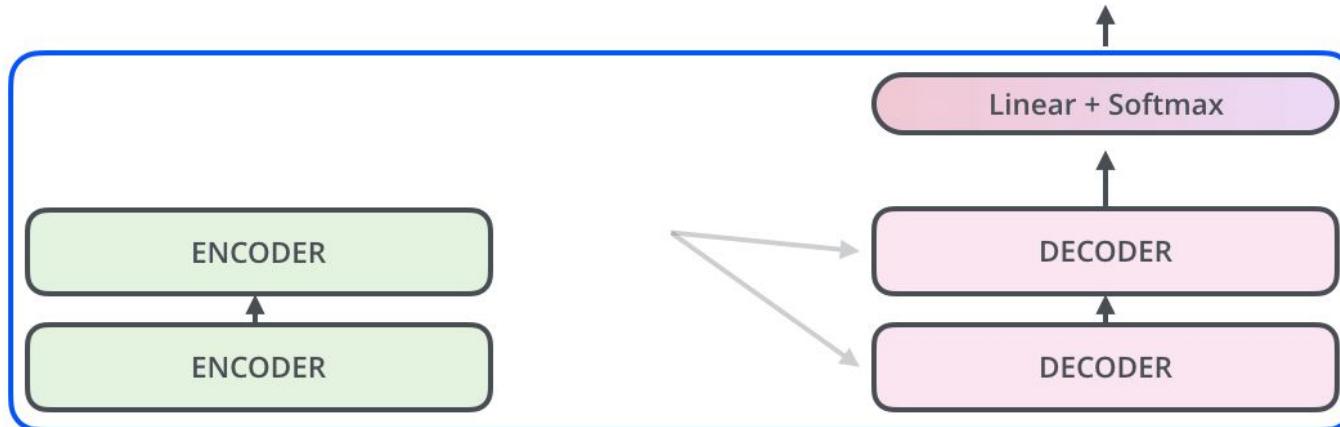
it_
was_
too_
tire
d_

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



Decoding time step: 1 2 3 4 5 6

OUTPUT



EMBEDDING
WITH TIME
SIGNAL



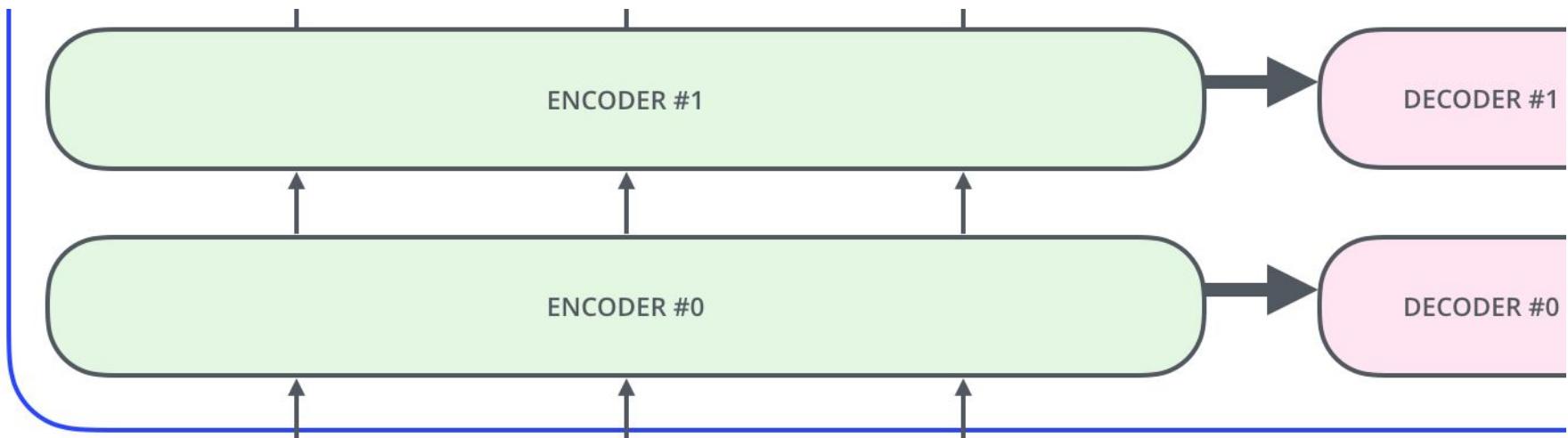
EMBEDDINGS



INPUT

Je suis étudiant

Positional Encoding



EMBEDDING
WITH TIME
SIGNAL

$$\mathbf{x}_1 \quad \boxed{\text{light green}} \quad \boxed{\text{light green}} \quad \boxed{\text{light green}}$$

$$\mathbf{x}_2 \quad \boxed{\text{light green}} \quad \boxed{\text{light green}} \quad \boxed{\text{light green}}$$

$$\mathbf{x}_3 \quad \boxed{\text{light green}} \quad \boxed{\text{light green}} \quad \boxed{\text{light green}}$$

=

POSITIONAL
ENCODING

$$\mathbf{t}_1 \quad \boxed{\text{yellow}} \quad \boxed{\text{yellow}} \quad \boxed{\text{yellow}}$$

$$\mathbf{t}_2 \quad \boxed{\text{yellow}} \quad \boxed{\text{yellow}} \quad \boxed{\text{yellow}}$$

$$\mathbf{t}_3 \quad \boxed{\text{yellow}} \quad \boxed{\text{yellow}} \quad \boxed{\text{yellow}}$$

+

EMBEDDINGS

$$\mathbf{x}_1 \quad \boxed{\text{green}} \quad \boxed{\text{green}} \quad \boxed{\text{green}}$$

$$\mathbf{x}_2 \quad \boxed{\text{green}} \quad \boxed{\text{green}} \quad \boxed{\text{green}}$$

$$\mathbf{x}_3 \quad \boxed{\text{green}} \quad \boxed{\text{green}} \quad \boxed{\text{green}}$$

INPUT

je

suis

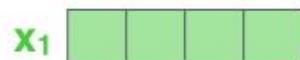
étudiant

POSITIONAL ENCODING

0	0	1	1
---	---	---	---

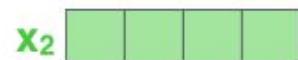
+

EMBEDDINGS



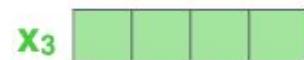
0.84	0.0001	0.54	1
------	--------	------	---

+



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

+



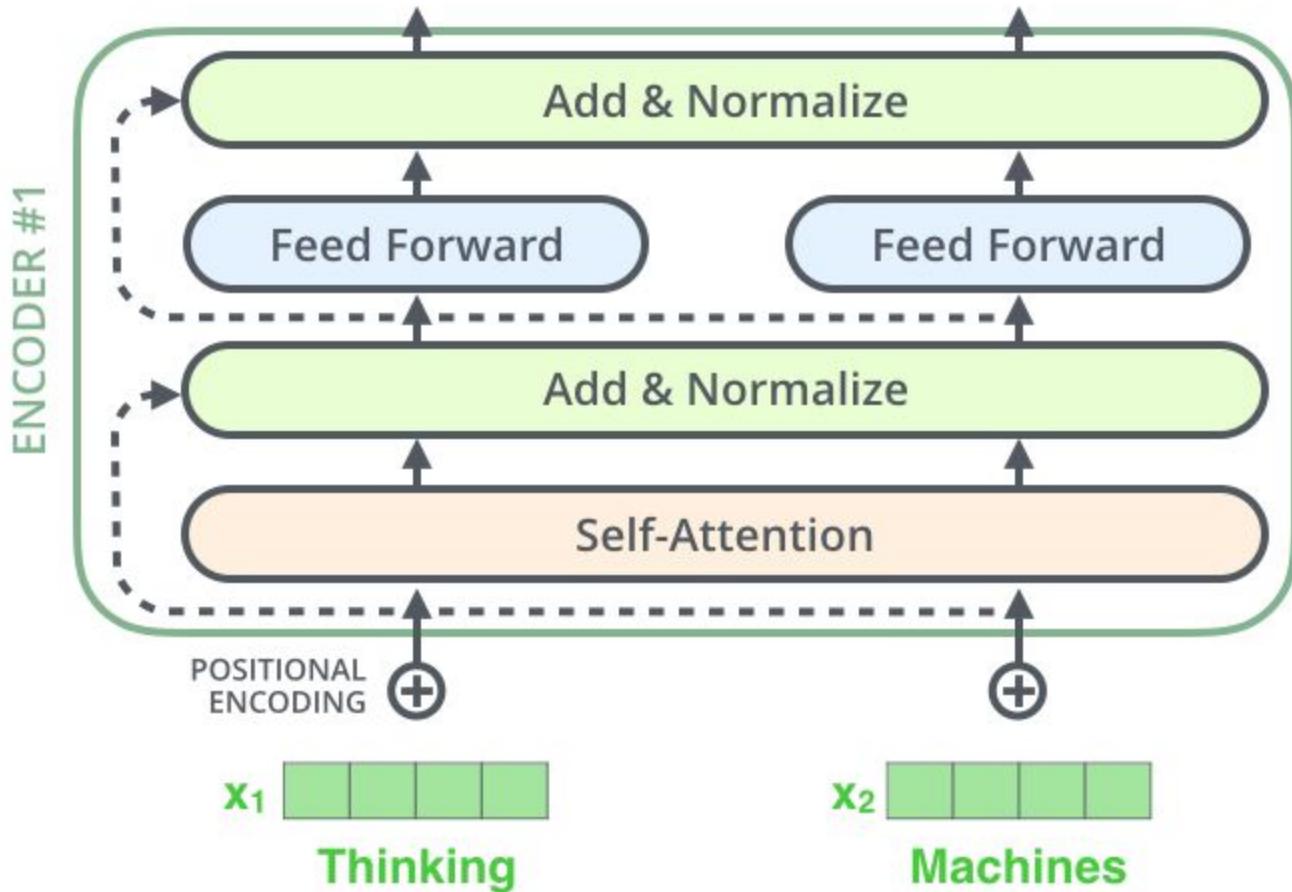
INPUT

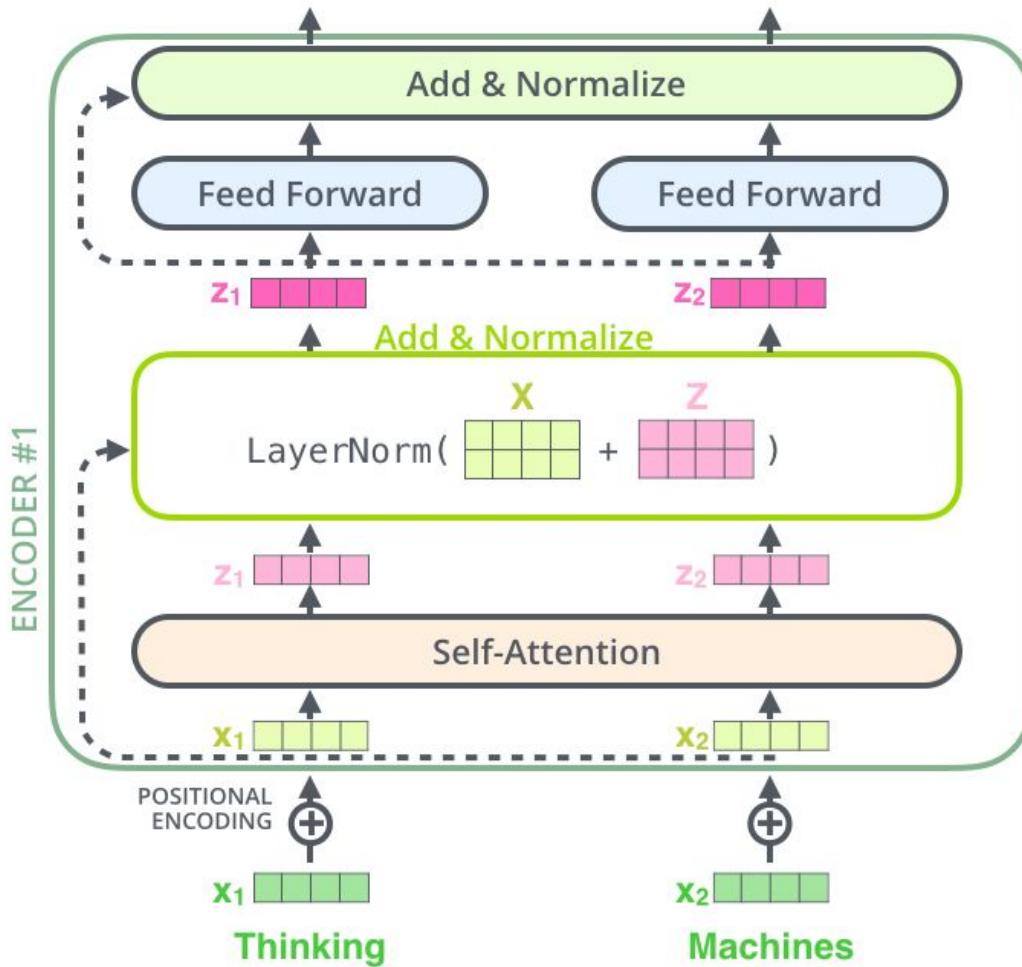
Je

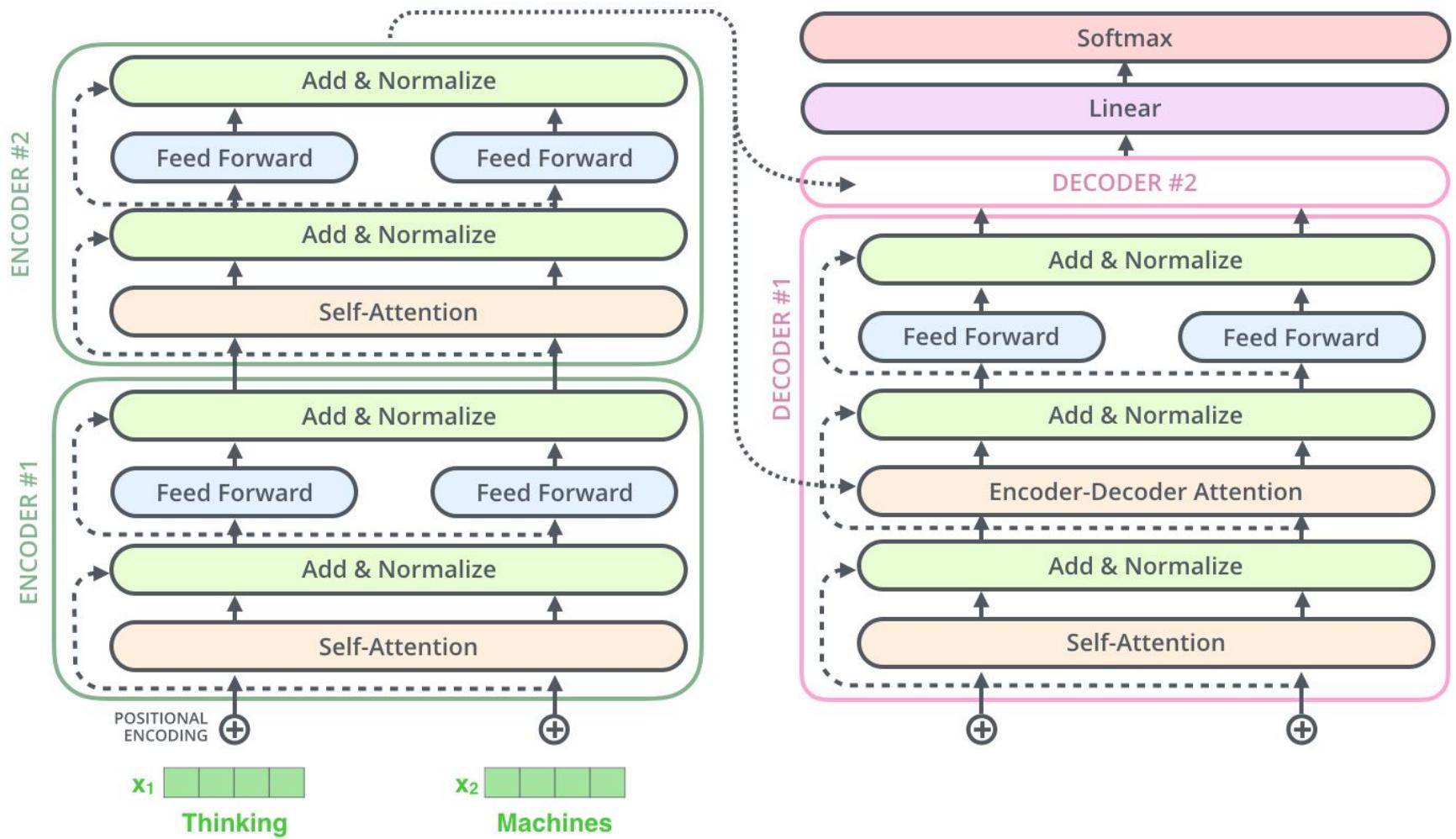
suis

étudiant

Residuals

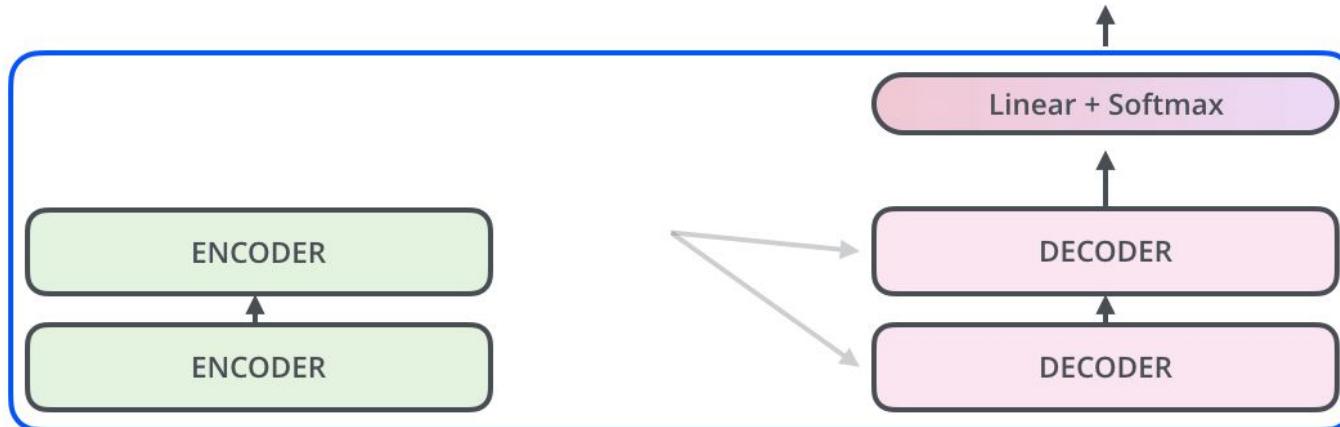






Decoding time step: 1 2 3 4 5 6

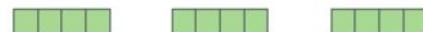
OUTPUT



EMBEDDING
WITH TIME
SIGNAL



EMBEDDINGS



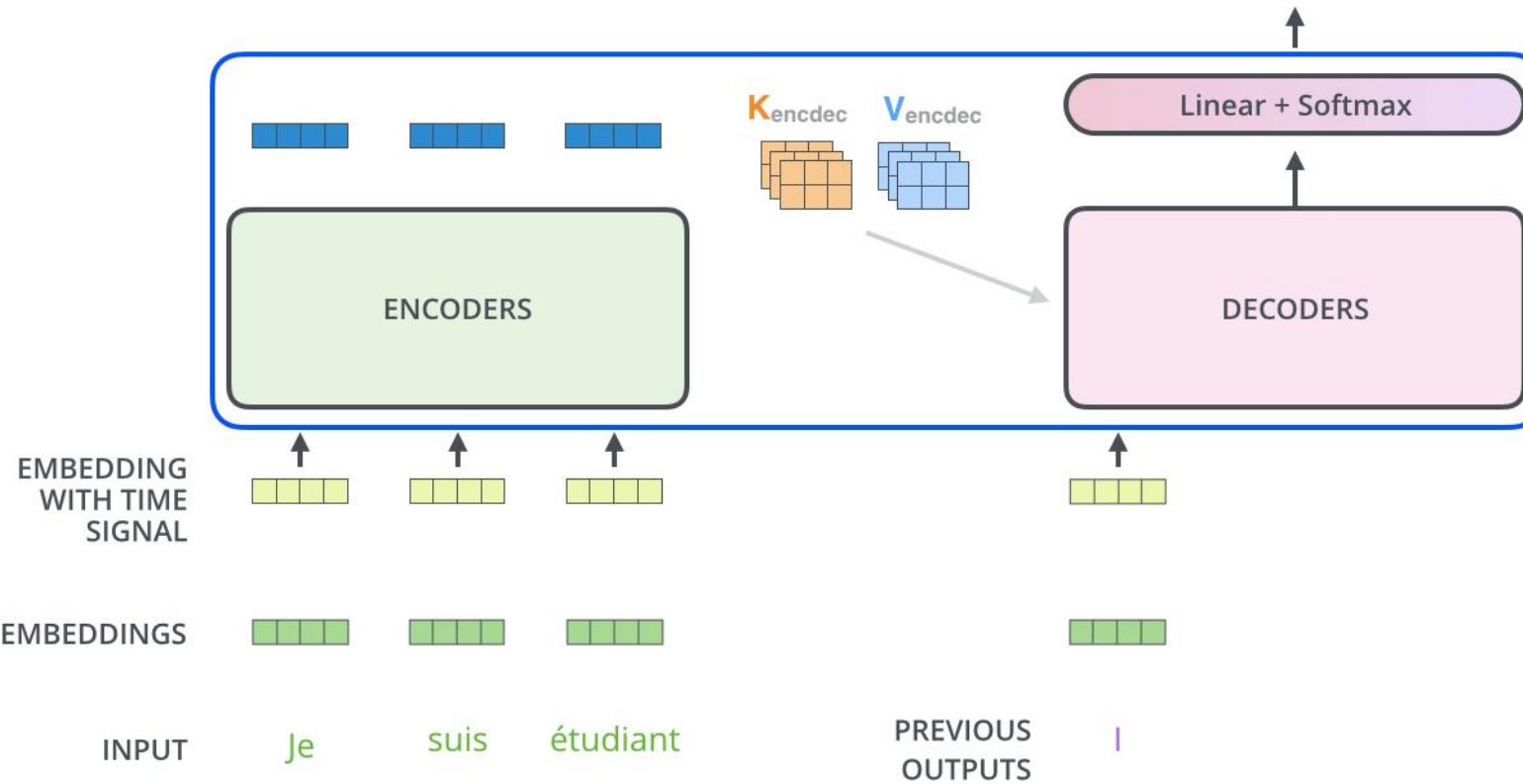
INPUT

Je suis étudiant

Decoding time step: 1 2 3 4 5 6

OUTPUT

|



Which word in our vocabulary
is associated with this index?

am

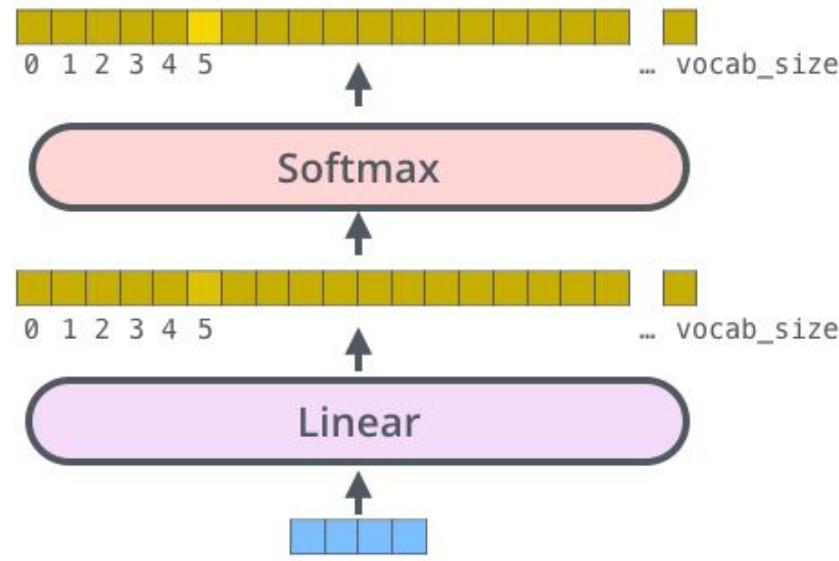
Get the index of the cell
with the highest value
(argmax)

5

log_probs

logits

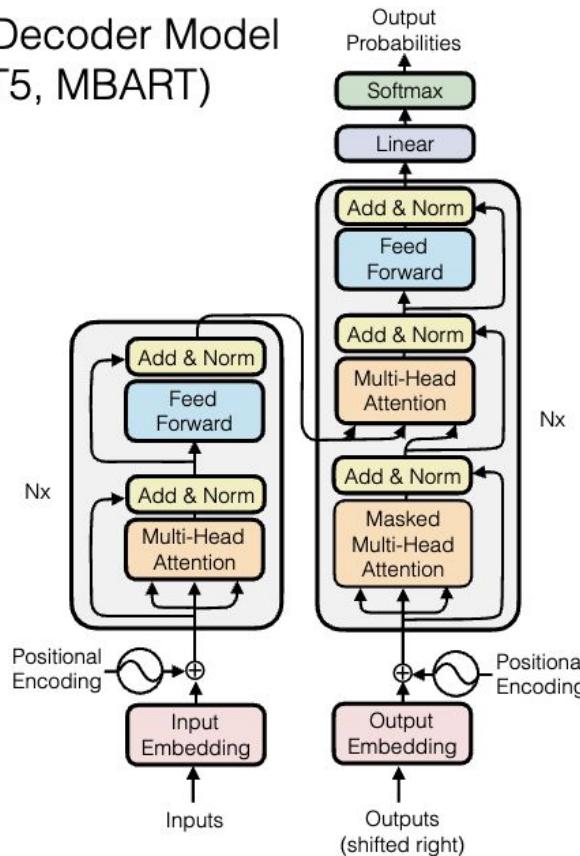
Decoder stack output



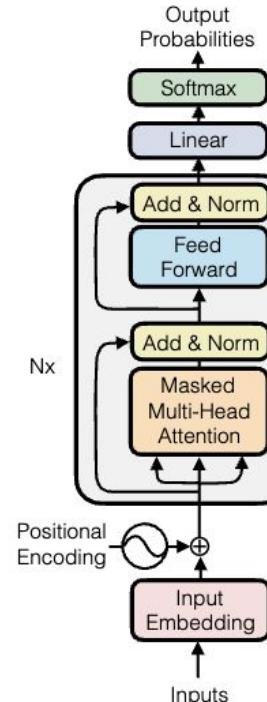
LLMs

Two Types of Transformers

Encoder-Decoder Model
(e.g. T5, MBART)



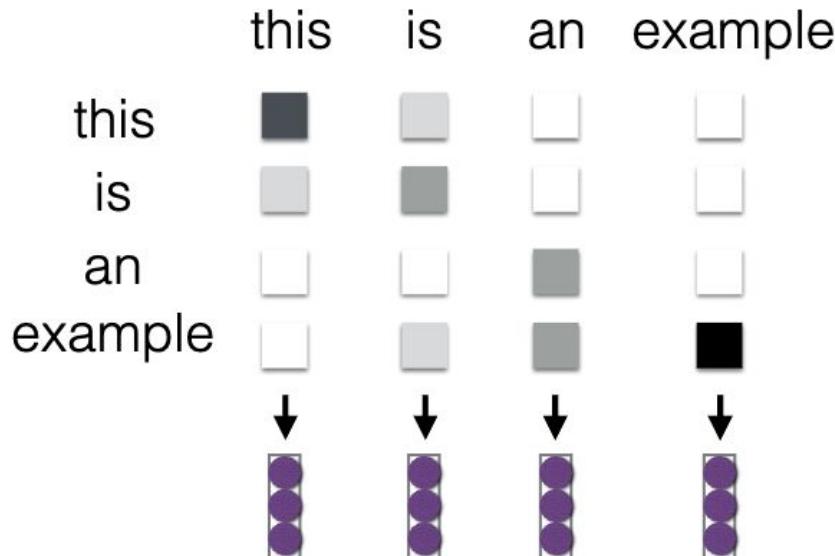
Decoder Only Model
(e.g. GPT, LLaMa)



Self Attention

(Cheng et al. 2016, Vaswani et al. 2017)

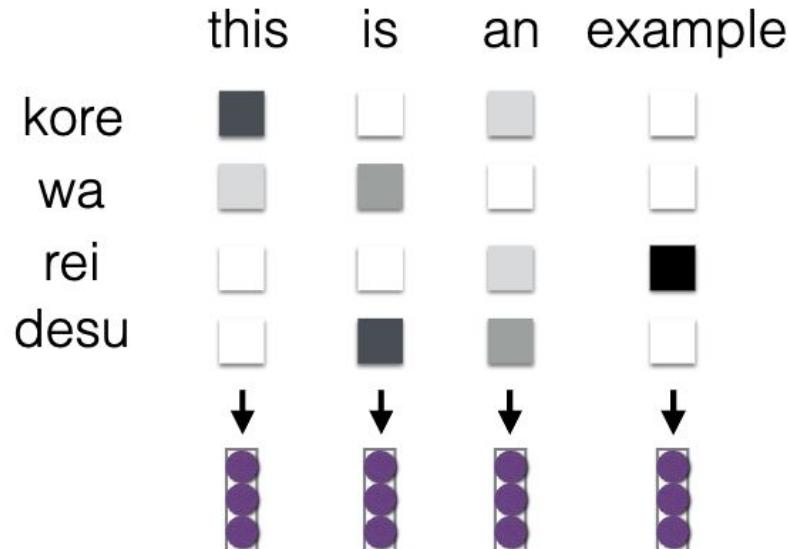
- Each element in the sequence attends to elements of that sequence



Cross Attention

(Bahdanau et al. 2015)

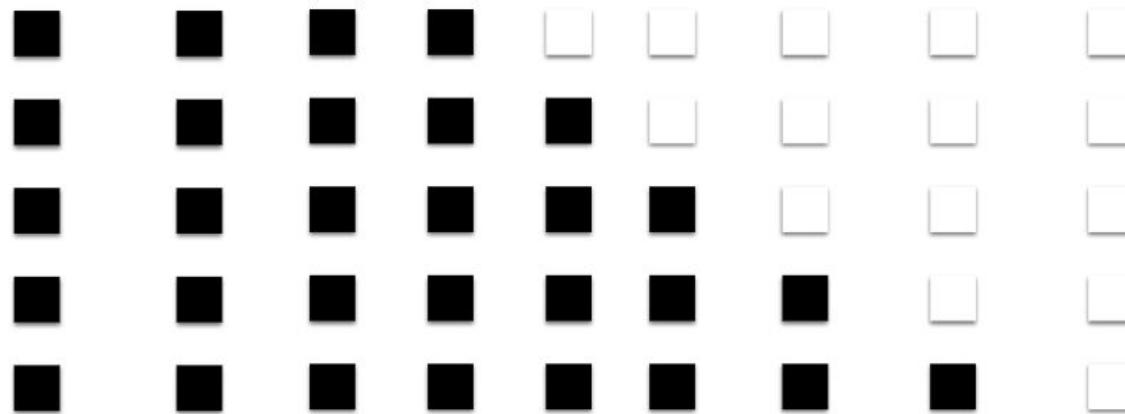
- Each element in a sequence attends to elements of another sequence



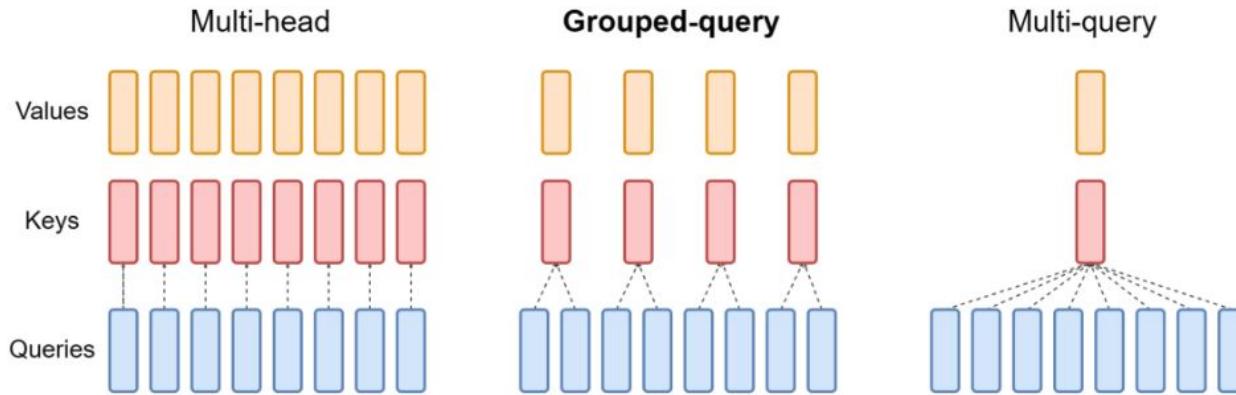
Masking for Language Model Training

- Mask the attention from future timesteps
 - Prevents the model from cheating when predicting the next token

kono eiga ga kirai I hate this movie </s>



Grouped-query attention



- Shares key and value heads for each *group* of query heads
- Saves on memory, which leads to faster inference

Original Transformer vs. Llama

	Vaswani et al.	LLama	Llama 2
Norm Position	Post	Pre	Pre
Norm Type	LayerNorm	RMSNorm	RMSNorm
Non-linearity	ReLU	SwiGLU	SwiGLU
Positional Encoding	Sinusoidal	RoPE	RoPE
Attention	Multi-head	Multi-head	Grouped-query