

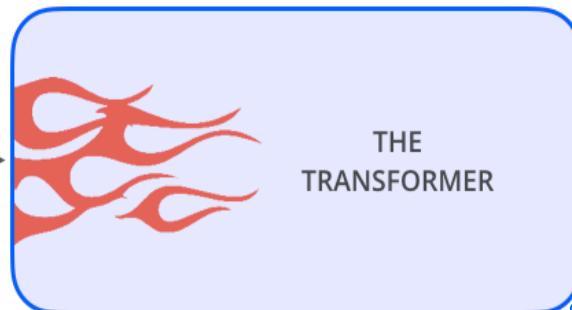
# This material is taken from Jay Alammar's blog "The Illustrated Transformer"

As aliens entered our planet and began to colonize earth a certain group of extraterrestrial



INPUT

Je suis étudiant

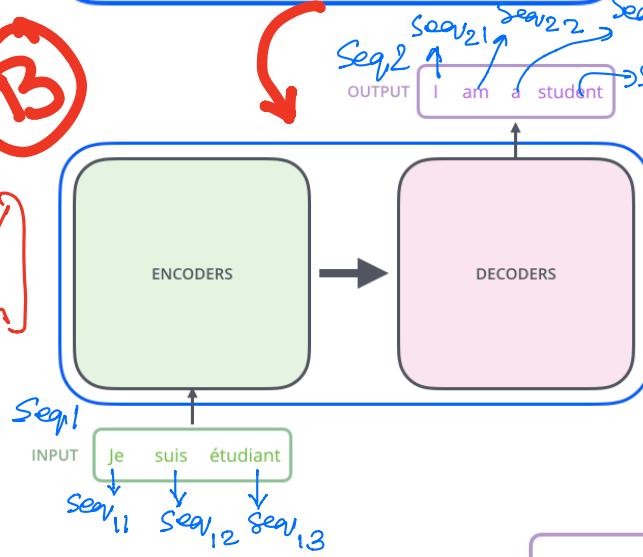


OUTPUT

I am a student



Encoder learns info for the entire Sequence [Some common/intermediate representation]  
Abstract & continuous representation



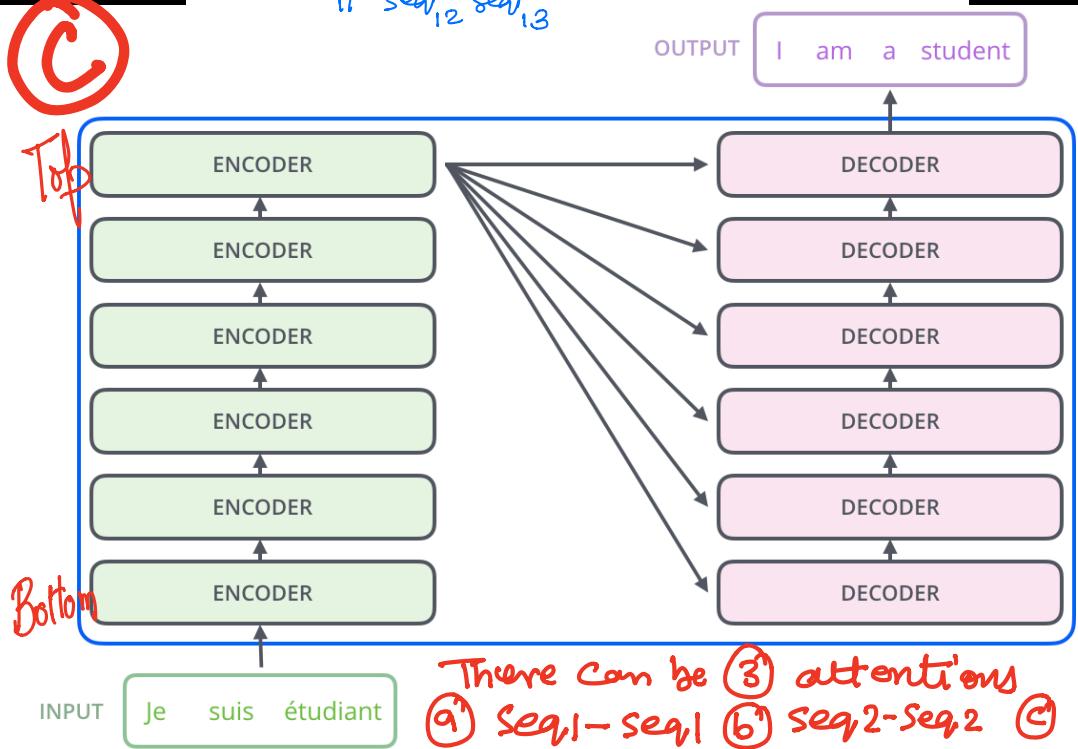
It takes that representation along with the previously generated / O/P word

↓  
Generate the O/P sequence Step by Step.

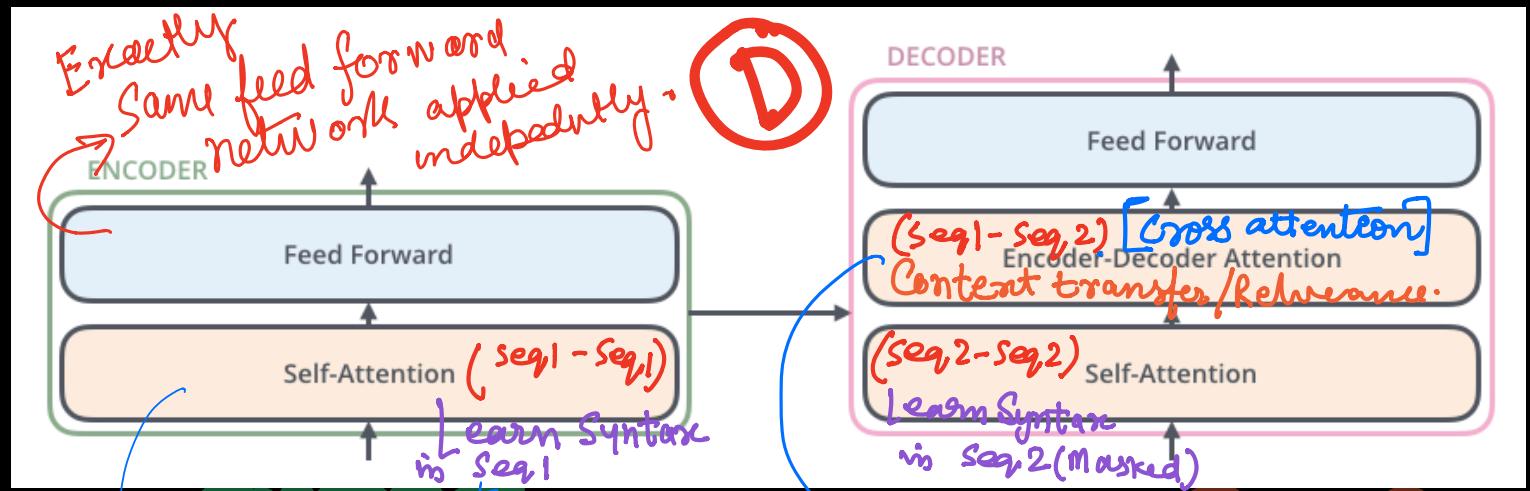


Top

All encoders are identical  
But do not share weights



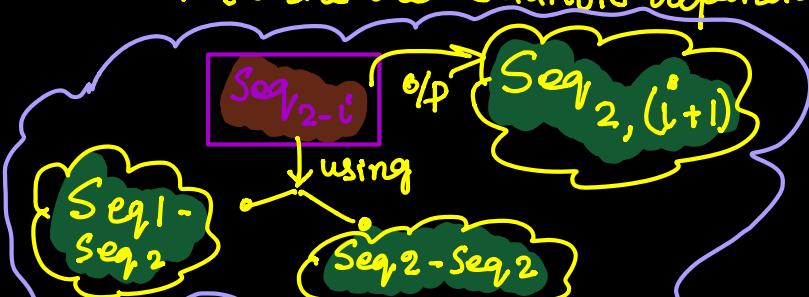
There can be (3) attentions  
① seq1-seq1 ② seq2-seq2 ③ seq1-seq2



Helps encoder to consider other input words as it encodes a specific word

Sequence  $\Rightarrow$  Syntax + Semantic

Encoder learns the semantic dependence  $\Leftrightarrow$  Attention



Additional layer in Decoder helps to focus on relevant input words as it decodes a specific word.

Word2Vec encoding

$x_1$  [green boxes]

Je

512-d vectors

$x_2$  [green boxes]

suis

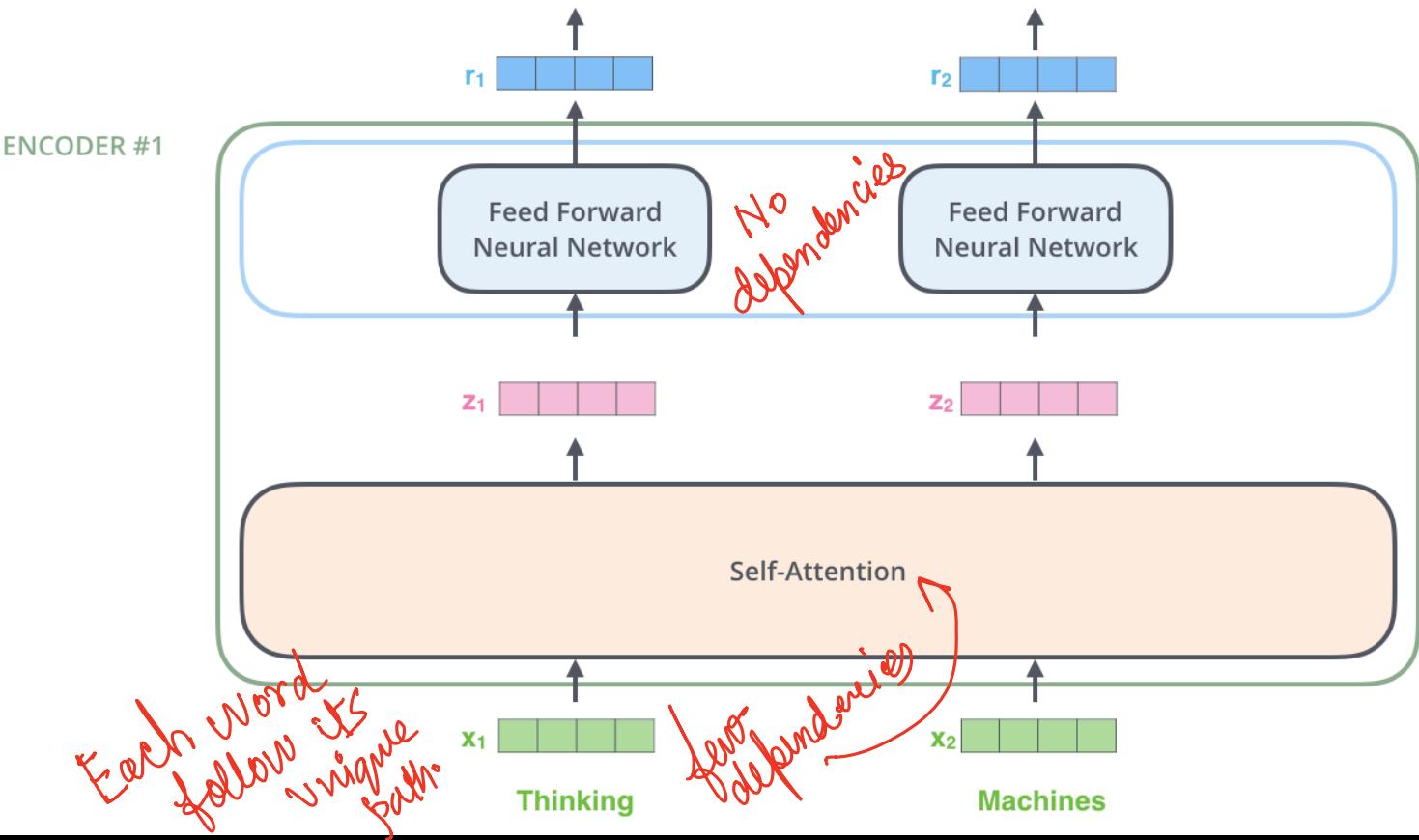
$x_3$  [green boxes]

étudiant

- ① Bottom encoder gets word embedding.
- ② Others get the O/p of previous encoders.
- ③ Encoder receives a list of embedding vectors of input sentences.

ENCODER #2

ENCODER #1



II

## Self-Attention (SA)

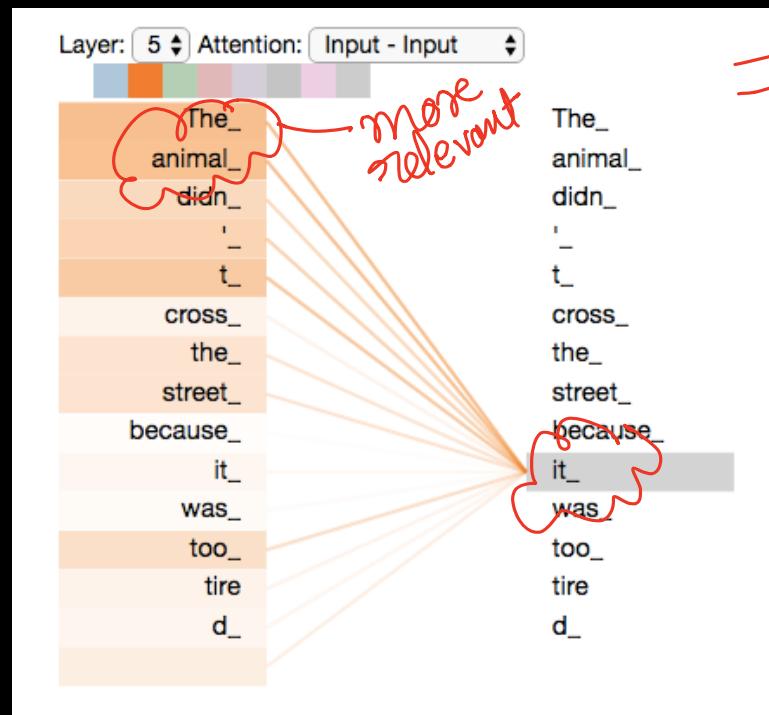
"The animal didn't cross the street because it was too tired"

Not Referencing

Referencing

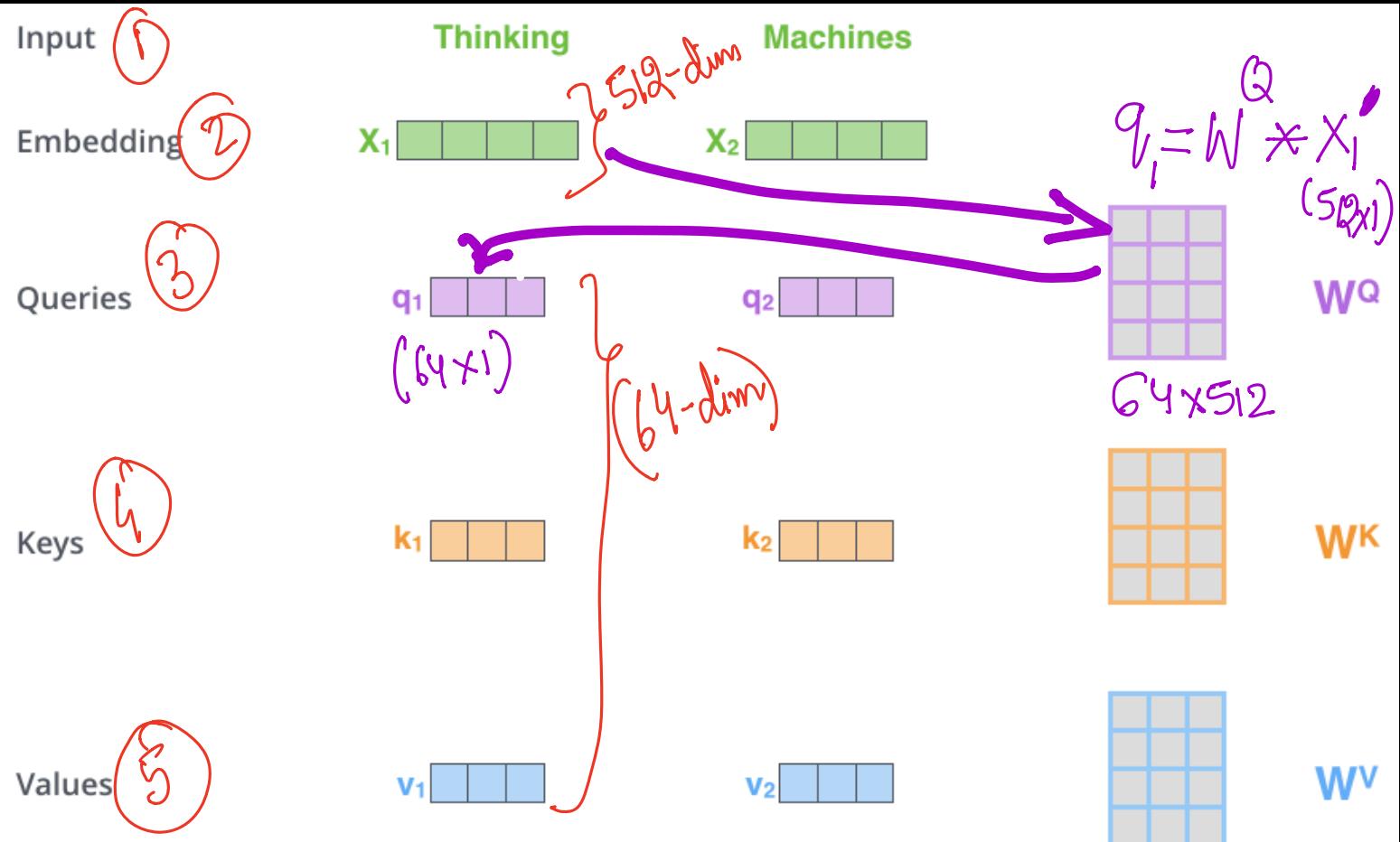
SA mostly discovers such references and uses them as clues to encode the word in terms of other related words

Never while processing "it", SA allows to associate with animal.



Showing the top encoders learned attention. SA gives more attention to "The" while encoding "It". Basically Matrices need to be learned.

III How SA is calculated :



Role of Query / Key / Value need to be understood.

① Input

② Embedding

③ Queries

④ Keys

⑤ Values

⑥ Score

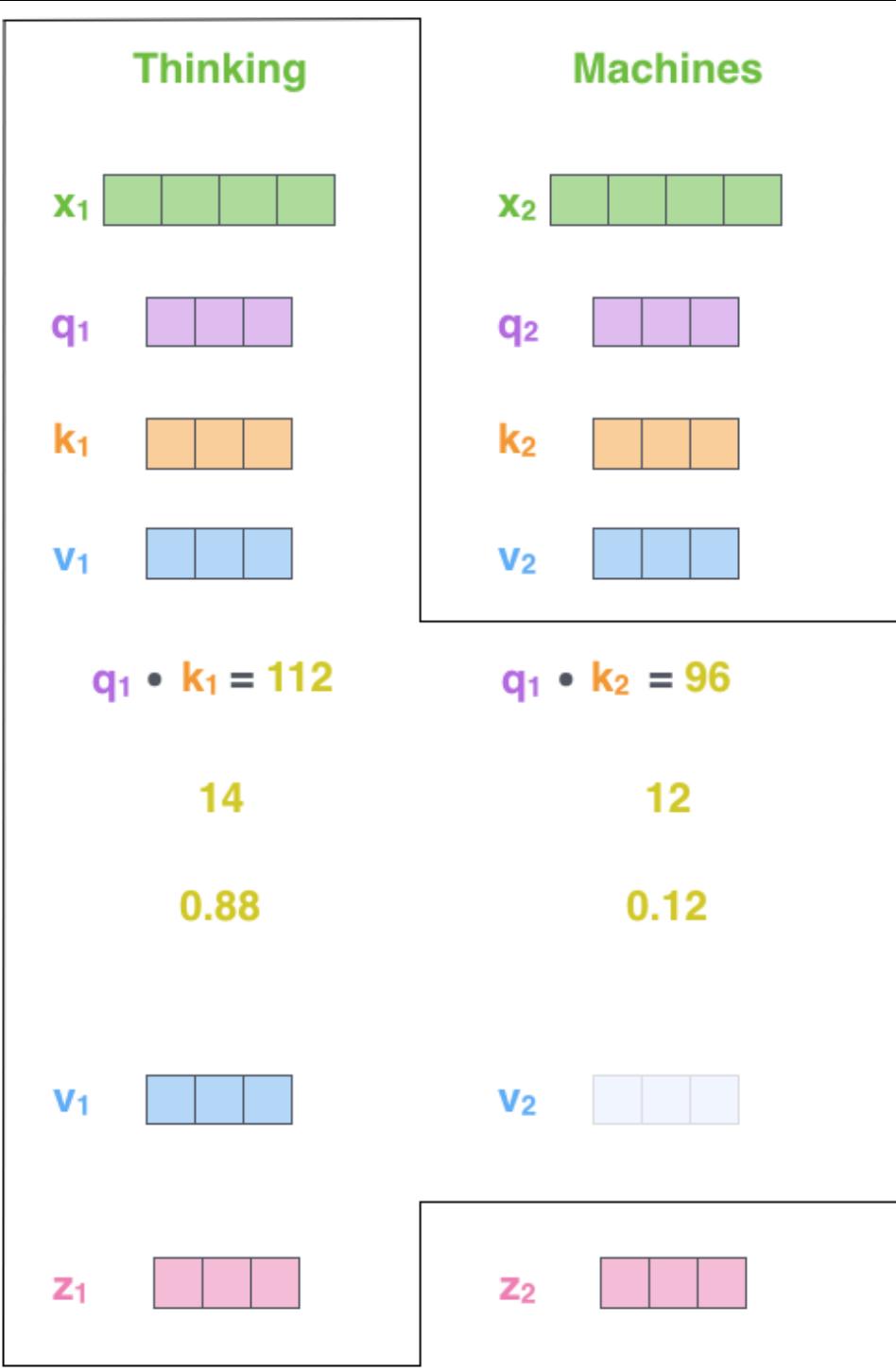
⑦ Divide by 8 ( $\sqrt{d_k}$ )

⑧ Softmax

⑨ Softmax

X  
Value

⑩ Sum



★  $z_i$ 's can be sent to feed forward MLP network.

★ These operations are done in matrix form for faster processing.

$$\begin{array}{c}
 \textbf{X} \\
 \xrightarrow{\substack{\rightarrow 512-d \\ \rightarrow \text{word-01} \\ \rightarrow \text{word-02}}} \\
 \boxed{\begin{matrix} & & \\ & & \\ & & \end{matrix}} \\
 2 \times 512
 \end{array}
 \times
 \begin{array}{c}
 \textbf{W}^Q \\
 \boxed{\begin{matrix} & & \\ & & \\ & & \end{matrix}} \\
 512 \times 64
 \end{array}
 =
 \begin{array}{c}
 \textbf{Q} \\
 \xrightarrow{\substack{\rightarrow Q-01 \\ \rightarrow Q-02}} \\
 \boxed{\begin{matrix} & & \\ & & \end{matrix}} \\
 2 \times 64
 \end{array}$$

$$\begin{array}{c}
 \textbf{X} \\
 \xrightarrow{\substack{\text{Input Seq.} \\ \text{Words} \rightarrow}} \\
 \boxed{\begin{matrix} & & \\ & & \\ & & \end{matrix}} \\
 2 \times 512
 \end{array}
 \times
 \begin{array}{c}
 \textbf{W}^K \\
 \boxed{\begin{matrix} & & \\ & & \\ & & \end{matrix}} \\
 512 \times 64
 \end{array}
 =
 \begin{array}{c}
 \textbf{K} \\
 \boxed{\begin{matrix} & & \\ & & \\ & & \end{matrix}} \\
 2 \times 64
 \end{array}$$

$$\begin{array}{c}
 \textbf{X} \\
 \boxed{\begin{matrix} & & & \\ & & & \\ & & & \end{matrix}} \\
 2 \times 512
 \end{array}
 \times
 \begin{array}{c}
 \textbf{W}^V \\
 \boxed{\begin{matrix} & & \\ & & \\ & & \end{matrix}} \\
 512 \times 64
 \end{array}
 =
 \begin{array}{c}
 \textbf{V} \\
 \boxed{\begin{matrix} & & \\ & & \\ & & \end{matrix}} \\
 2 \times 64
 \end{array}$$

$$\text{softmax} \left( \frac{\begin{array}{c} \textbf{Q} \\ \xrightarrow{\substack{Q_1 \\ Q_2}} \end{array} \times \begin{array}{c} \textbf{K}^T \\ \xrightarrow{\substack{64 \times 2 \\ \boxed{\begin{matrix} & \\ & \end{matrix}} = [2 \times 2]}} \end{array}}{\sqrt{d_k}} \right) * \begin{array}{c} \textbf{V} \\ \boxed{\begin{matrix} & \\ & \end{matrix}} = [2 \times 64] \end{array}$$

$\begin{array}{|c|c|} \hline Q_1 K_1 & Q_1 K_2 \\ \hline Q_2 K_1 & Q_2 K_2 \\ \hline \end{array}$

IV

## Multi-head attention $\stackrel{?}{=}$

- \* Attention allows model to focus on different positions. ( $Z_i$ : weighted sum of all encodings)

But dominated by its own embedding

- \* MHA enables attention layer to learn multiple representation Subspace.

What do they have?  
8 heads.

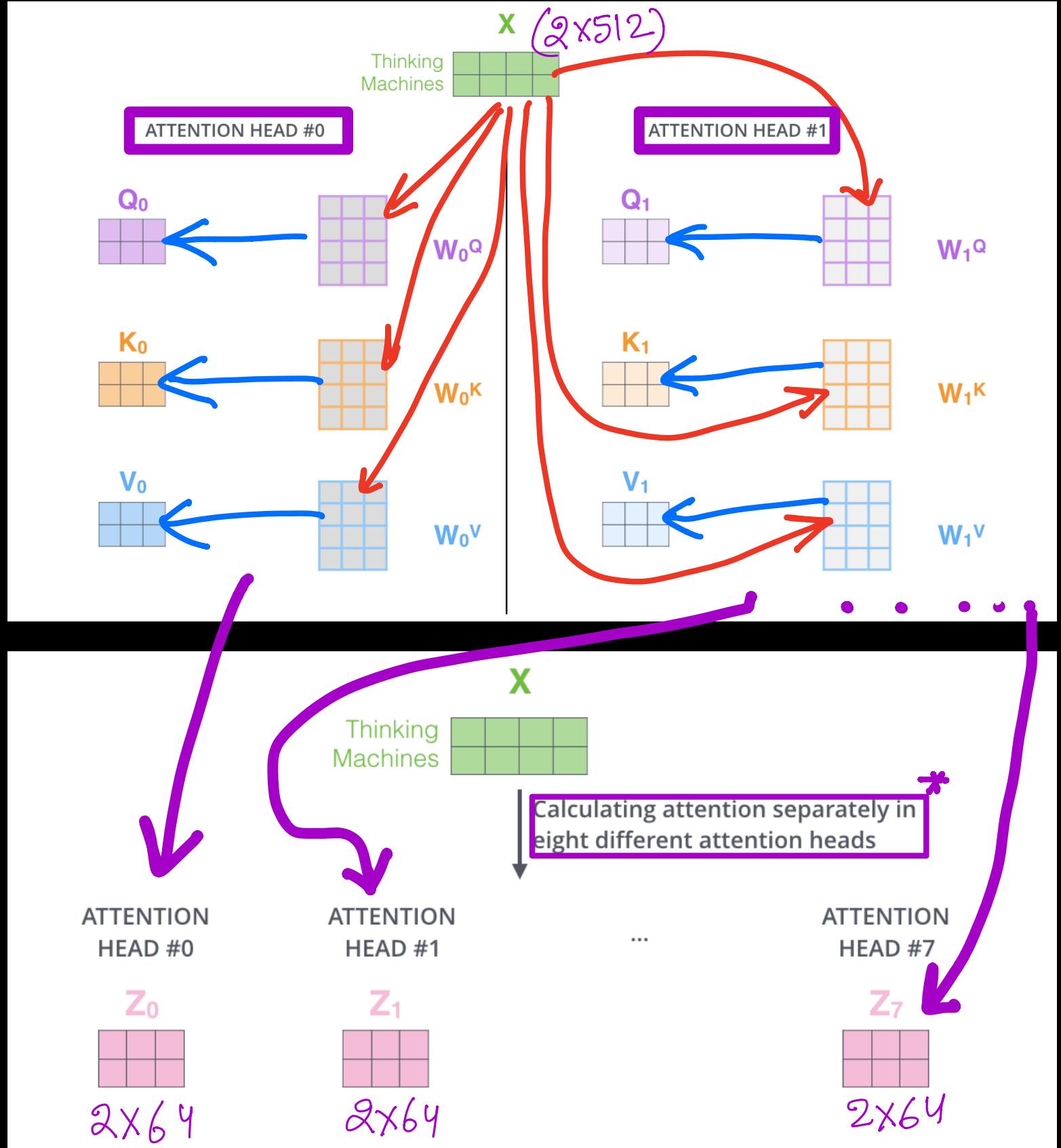
$H_1$

$H_2$

$$Q_1 \mid K_1 \mid V_1$$

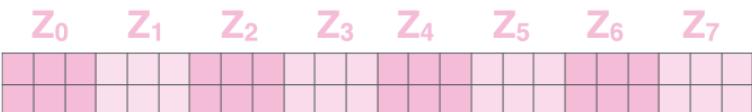
$$Q_2 \mid K_2 \mid V_2$$

Each input projects its into a different representation subspace.



Instead of one  $[Z]$  vector of size  $(2 \times 64)$   
we have got ⑧ such  $[2 \times 64]$  representations  
of  $(Z)$  vector.

1) Concatenate all the attention heads



$$Z_{Cat} = \text{Concat}_{i=0-7} [Z_i] [2 \times (64 \times 8)] \quad 2 \times 512$$

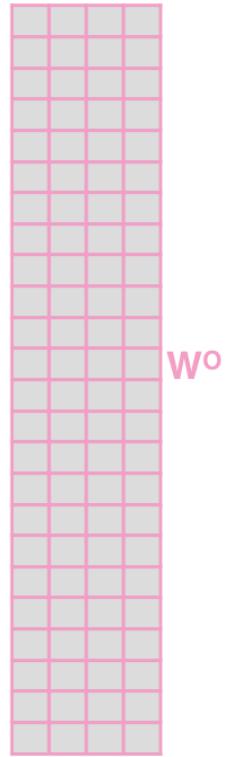
3) The result would be the  $Z$  matrix that captures information from all the attention heads. We can send this forward to the FFNN

$$= \begin{matrix} Z \\ \vdots \\ [2 \times 64] \end{matrix}$$

$$Z = Z_{Cat} * W_0 \quad [2 \times 64] \quad (2 \times 512) \quad (512 \times 64)$$

2) Multiply with a weight matrix  $W_0$  that was trained jointly with the model

$$x \quad W_0 \quad [512 \times 64]$$



# Recap till now.

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  $W_0$  to produce the output of the layer

Thinking Machines

$$X \quad \begin{matrix} \text{green} \\ \vdots \\ [2 \times 64] \end{matrix}$$

$$\begin{matrix} W_0^Q & W_0^K & W_0^V \\ \downarrow & \downarrow & \downarrow \\ \text{purple} & \text{orange} & \text{blue} \end{matrix}$$

$$\begin{matrix} Q_0 & K_0 & V_0 \\ \downarrow & \downarrow & \downarrow \\ \text{purple} & \text{orange} & \text{blue} \end{matrix}$$

$$Z_0 \quad \begin{matrix} \text{pink} \\ \vdots \\ [2 \times 64] \end{matrix}$$

$$W_0$$

$$Z \quad \begin{matrix} \text{pink} \\ \vdots \\ [2 \times 64] \end{matrix}$$

$$\begin{matrix} W_1^Q & W_1^K & W_1^V \\ \downarrow & \downarrow & \downarrow \\ \text{purple} & \text{orange} & \text{blue} \end{matrix}$$

$$\begin{matrix} Q_1 & K_1 & V_1 \\ \downarrow & \downarrow & \downarrow \\ \text{purple} & \text{orange} & \text{blue} \end{matrix}$$

$$Z_1 \quad \begin{matrix} \text{pink} \\ \vdots \\ [2 \times 64] \end{matrix}$$

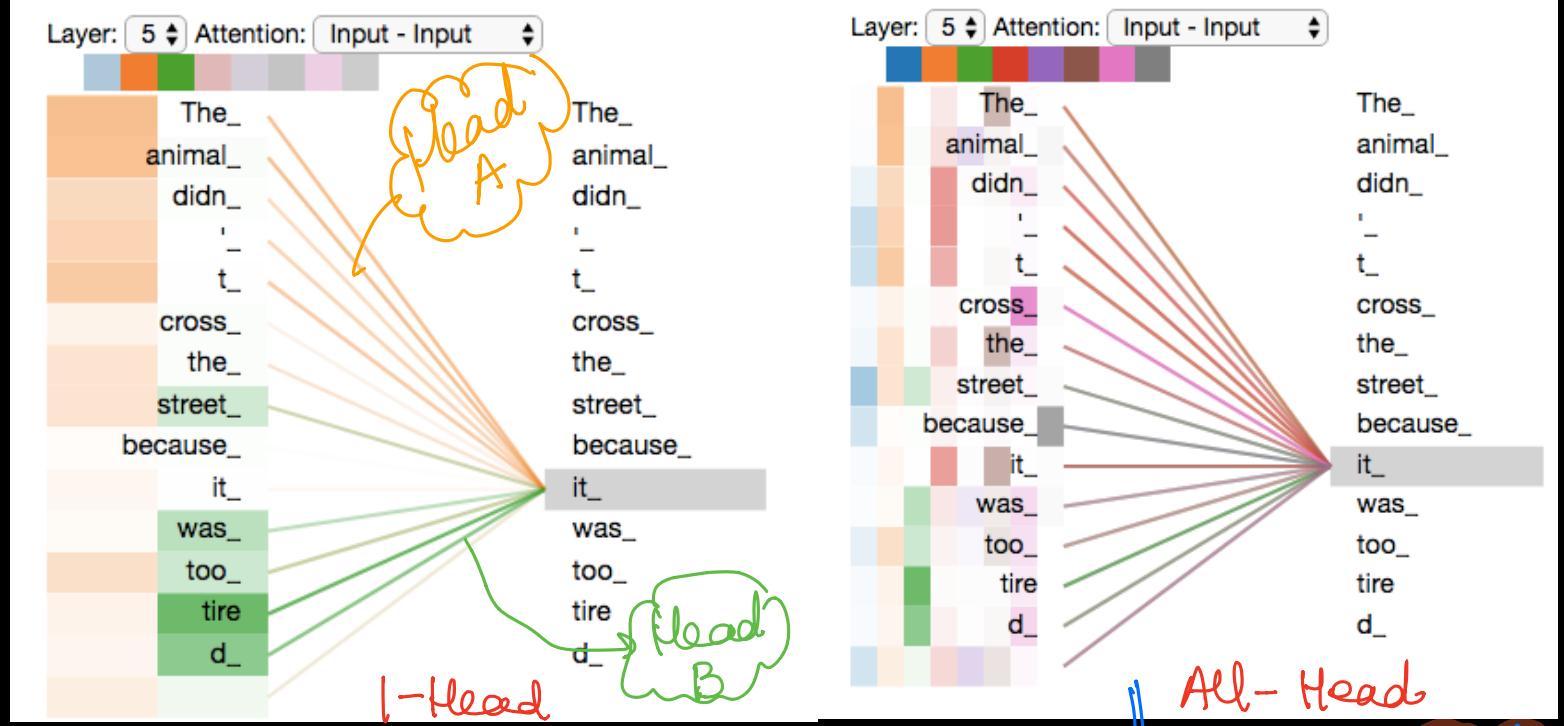
$$\dots \quad \begin{matrix} W_7^Q & W_7^K & W_7^V \\ \downarrow & \downarrow & \downarrow \\ \text{purple} & \text{orange} & \text{blue} \end{matrix}$$

$$\dots \quad \begin{matrix} Q_7 & K_7 & V_7 \\ \downarrow & \downarrow & \downarrow \\ \text{purple} & \text{orange} & \text{blue} \end{matrix}$$

$$Z_7 \quad \begin{matrix} \text{pink} \\ \vdots \\ [2 \times 64] \end{matrix}$$

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

$$R \quad \begin{matrix} \text{blue} \\ \vdots \\ [2 \times 64] \end{matrix}$$



While encoding "it"

One head is very focusing  
on "animal" or "The".  
(orange)

Other head may focus  
on "tired"

Accumulating  
the attention  
learned by all  
the heads.

together "it"  
is being getting  
Contributions from  
"Animal", "The", "Tired"

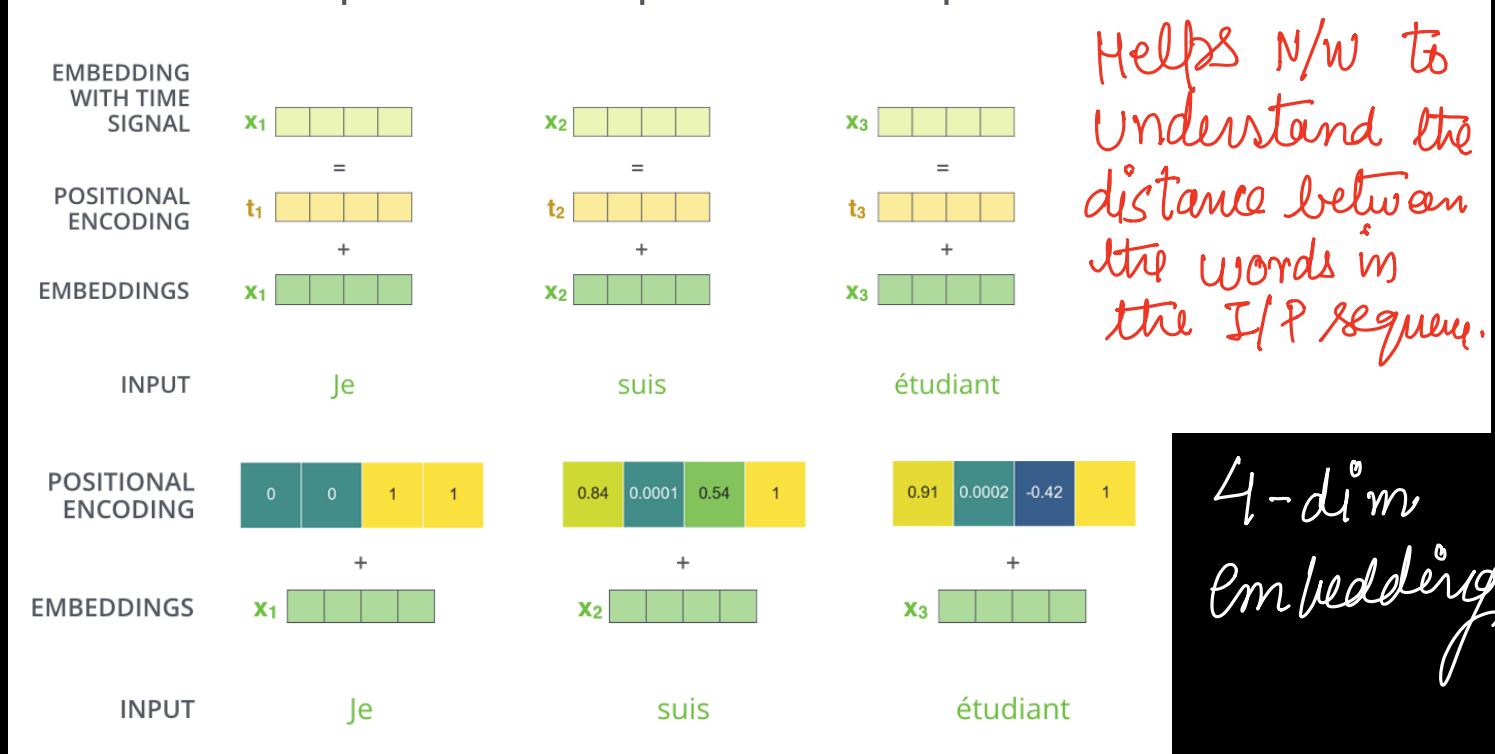
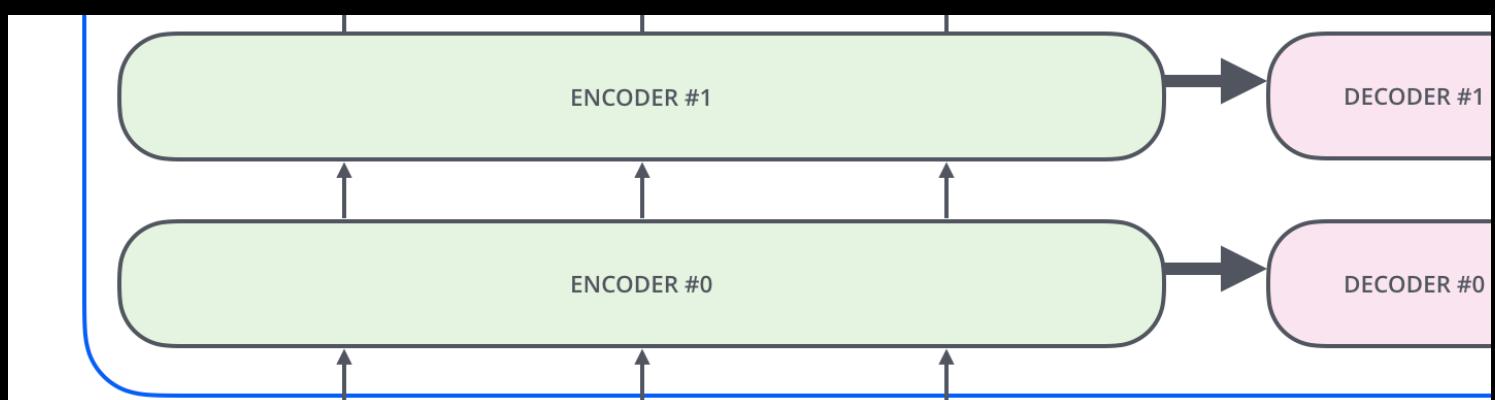
⌚ Adding more heads may be helpful

but later start to learn redundant  
attention or tend to overfit the data.

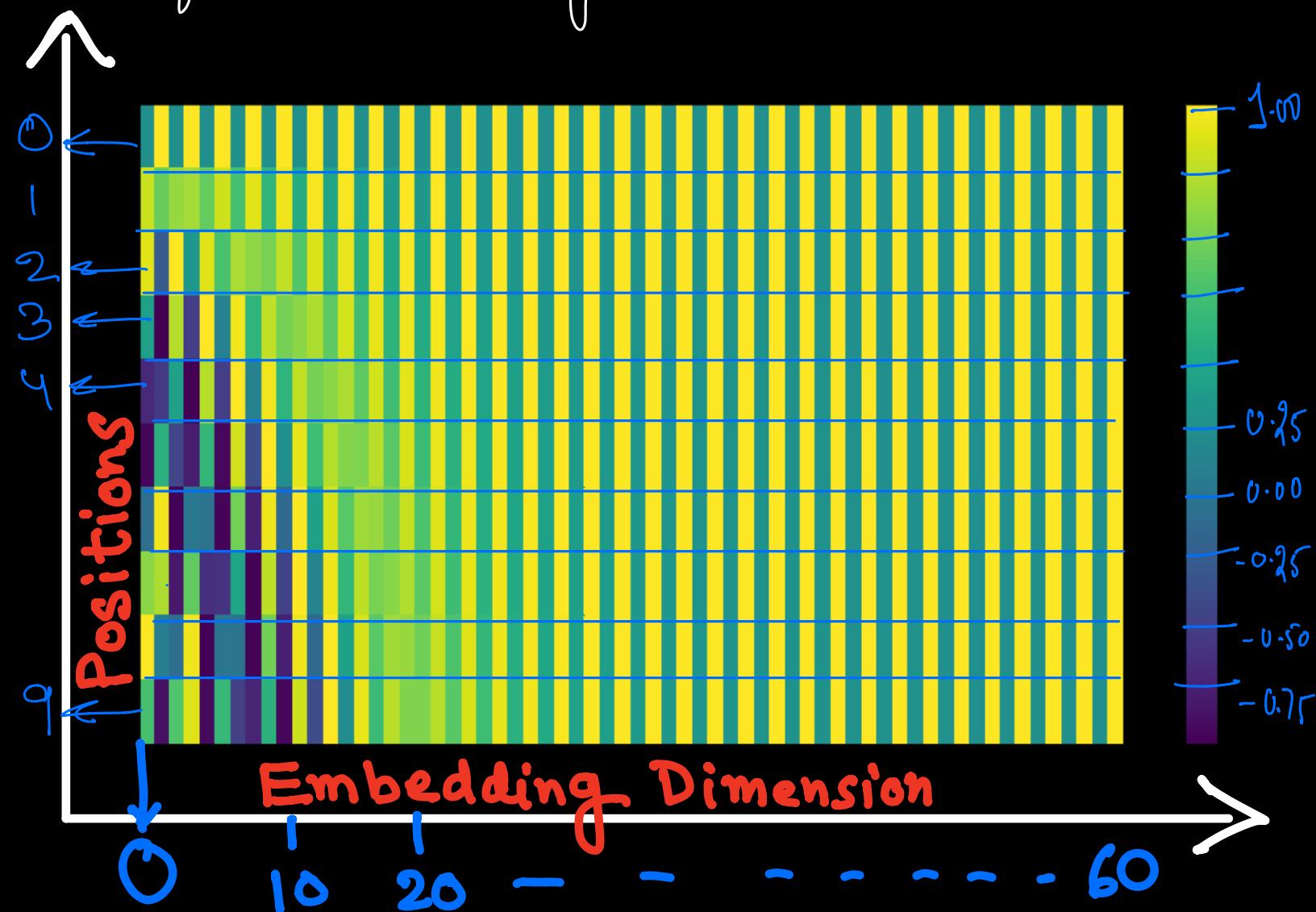
# Positional Embedding

- ⑥ Transformers are permutation equivariant.
- ⑦ Hence alongwith word embeddings we need to pass positional embeddings.

$$I/P = W \cdot E + P \cdot E$$

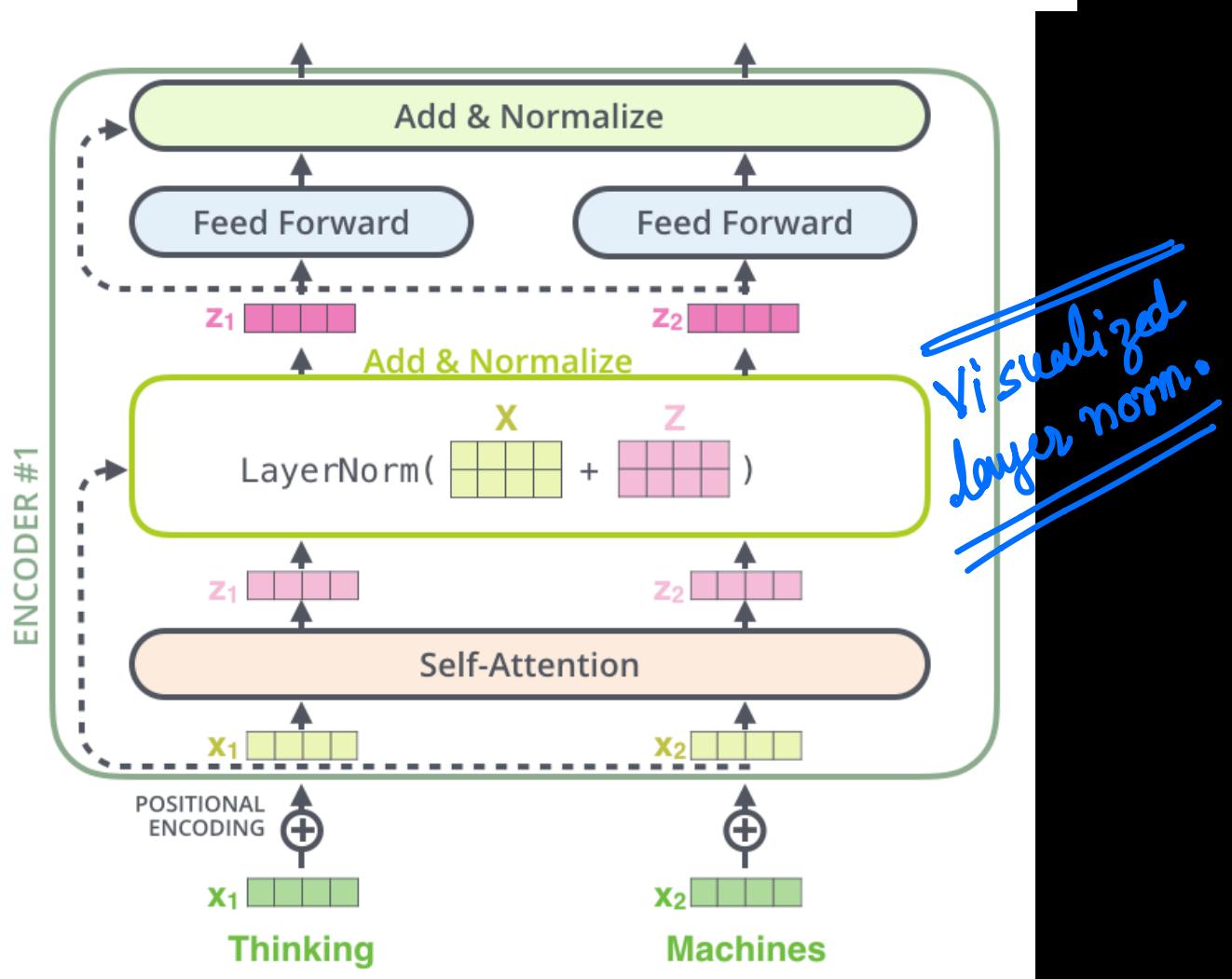
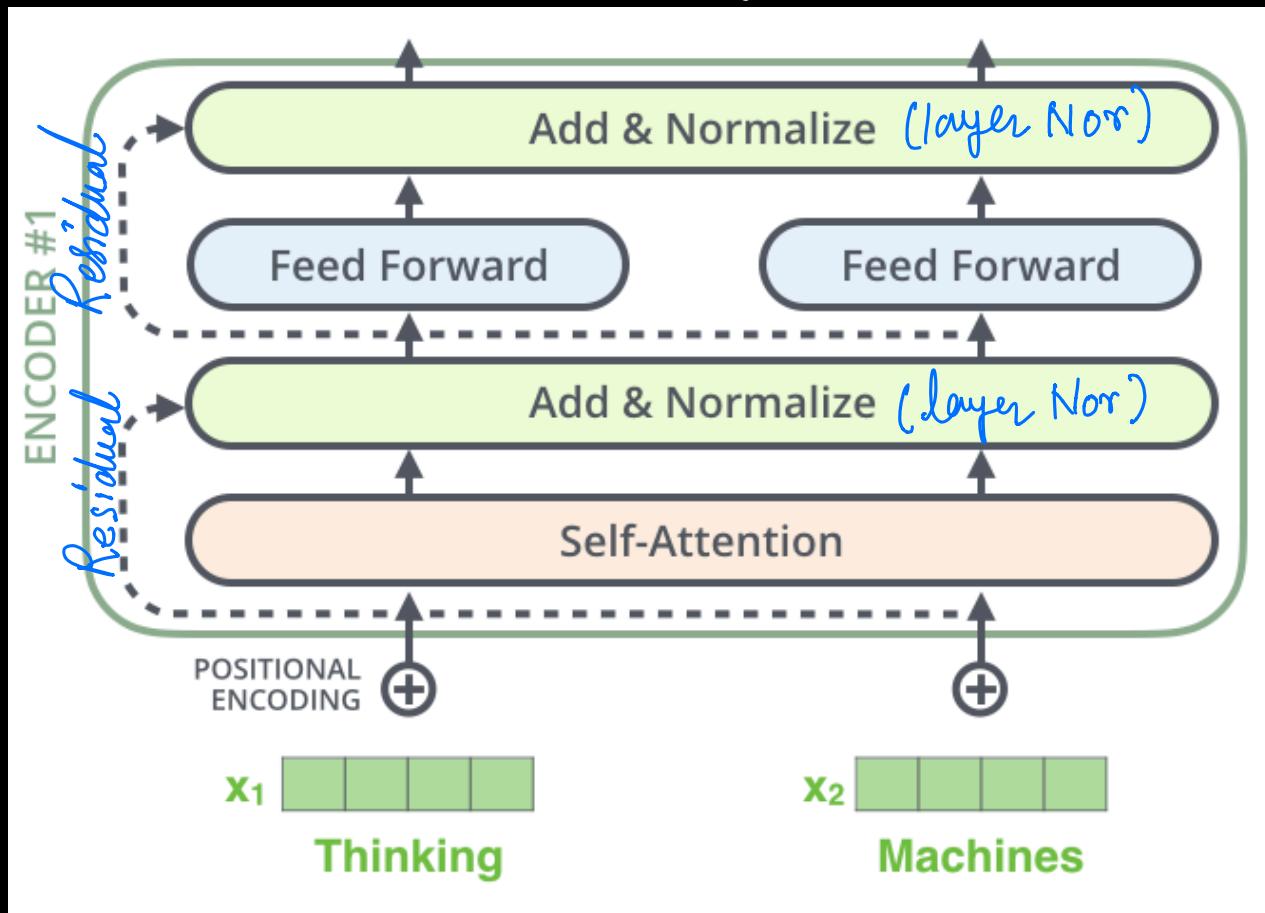


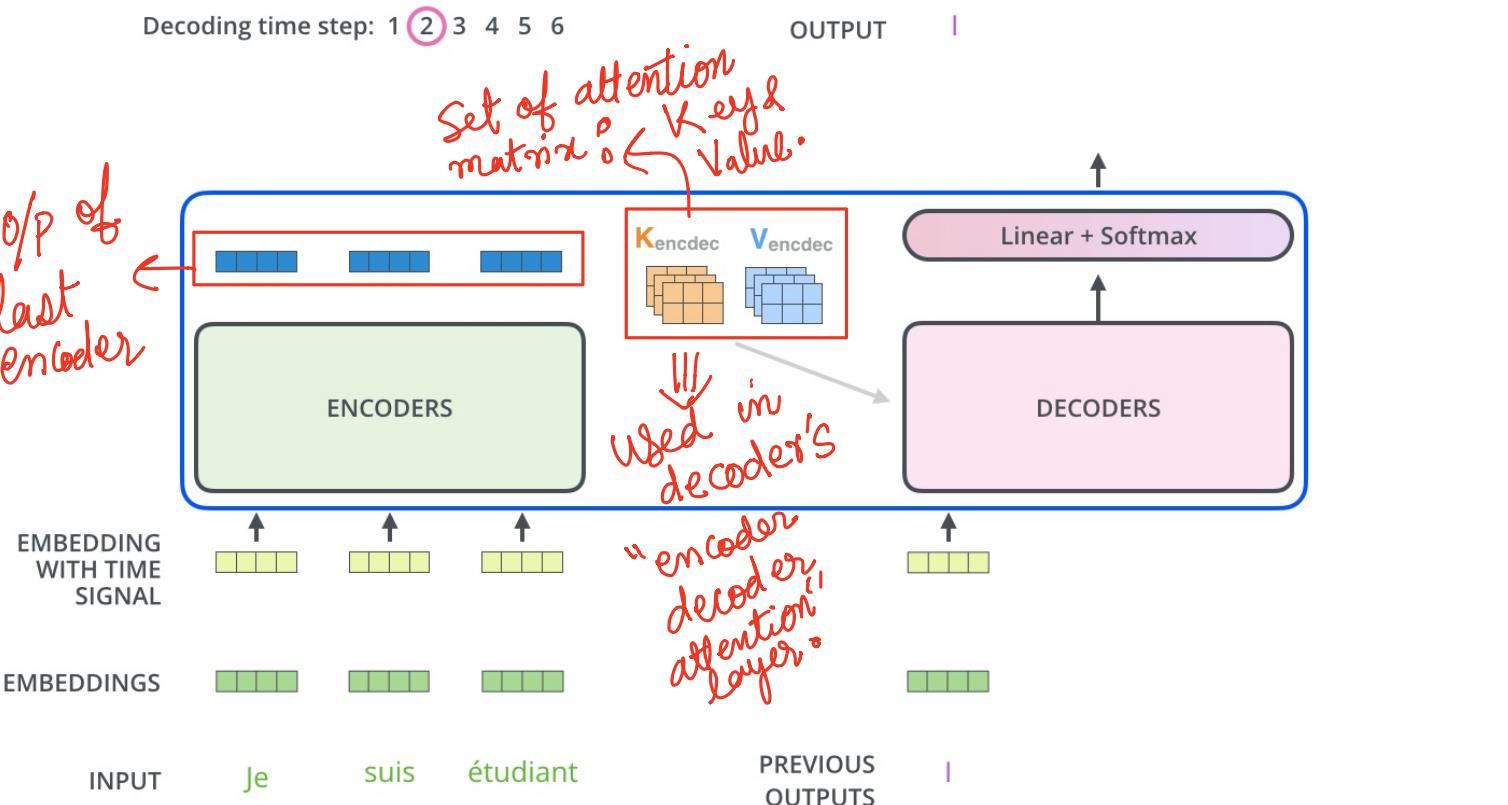
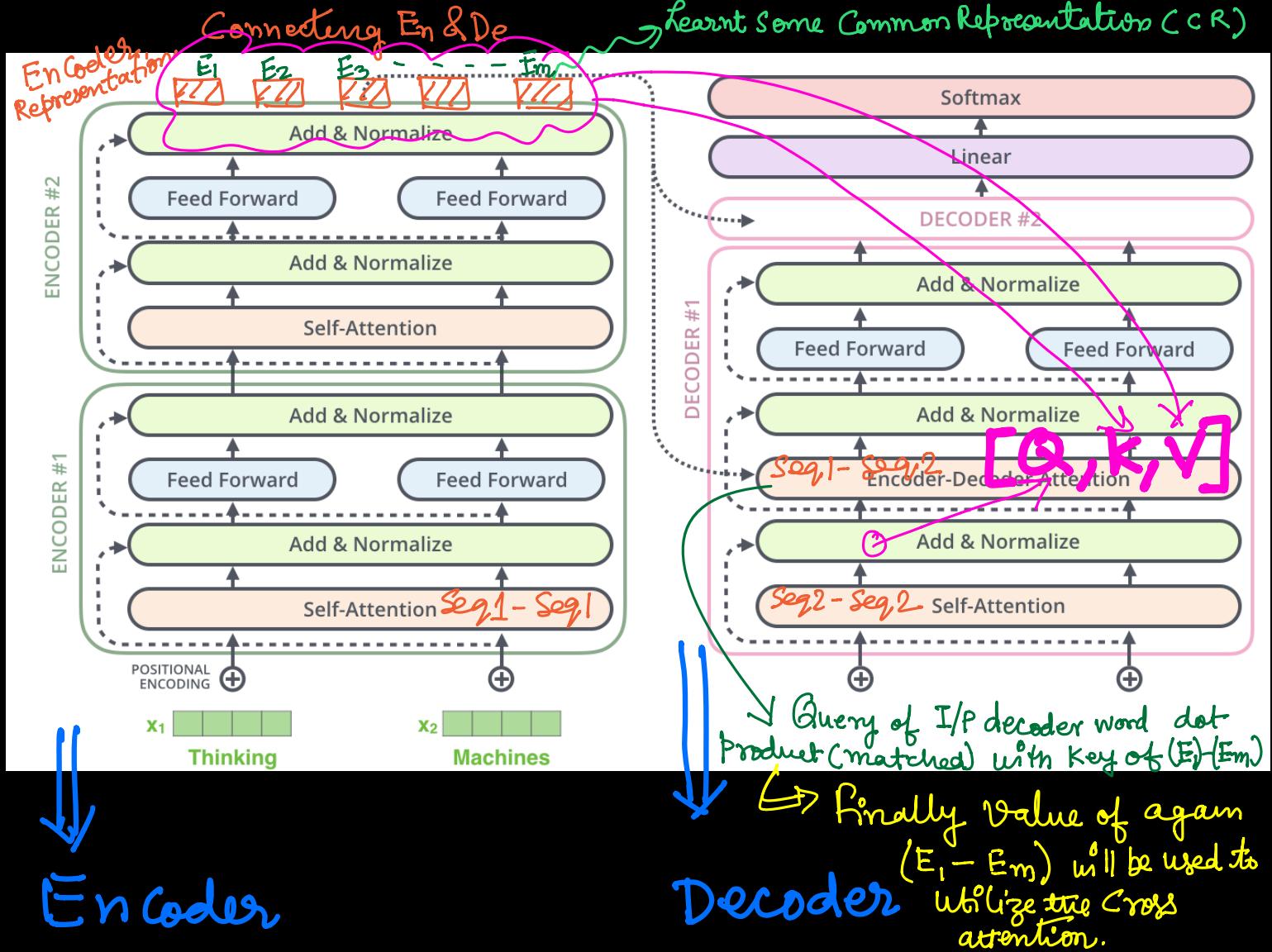
PE also helps to scale to unseen lengths of the testing sequences.



Both Sin & Cos functions are interleaved.

In encoder each sublayer has a residual connection, followed by layer normalization.





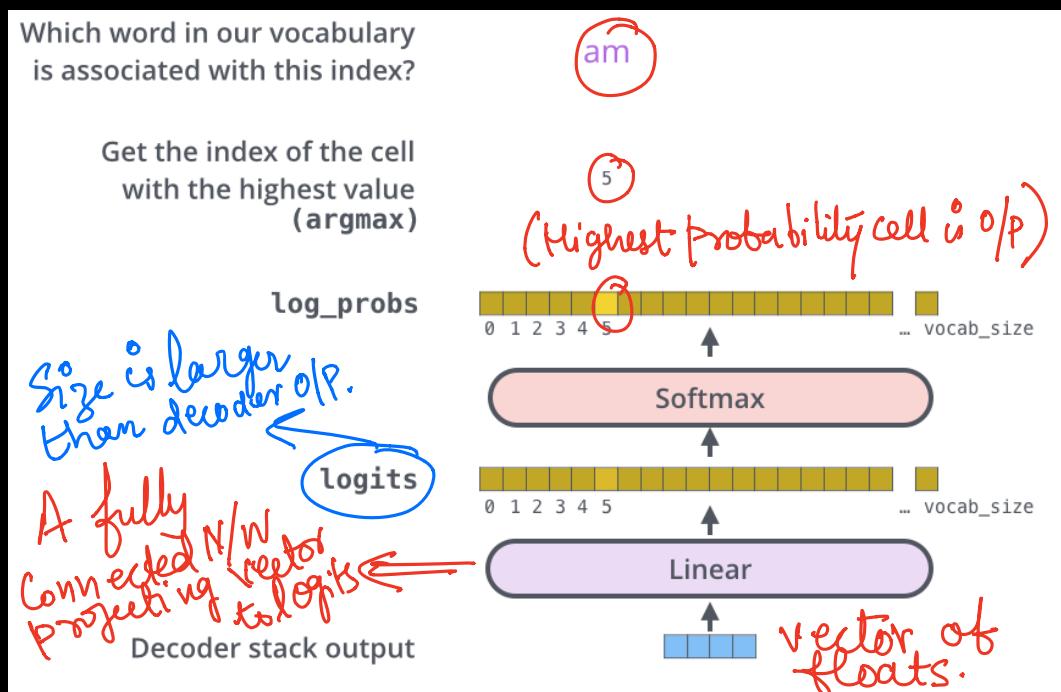
Decoder uses only previously shown input.

This is done by masking future positions by setting them to  $(-\infty)$  not  $(0)$  before Softmax step in self attention calculations.

\* Encoder Decoder Attention layer works similar to [MMSA] but only creates Query matrix & takes Key & Value (from previous decoder layer)  $[Q]$  (from encoder O/P layer)  $[K]$   $[V]$  matrix from O/P of encoder stack.

## Final linear & Softmax Layer

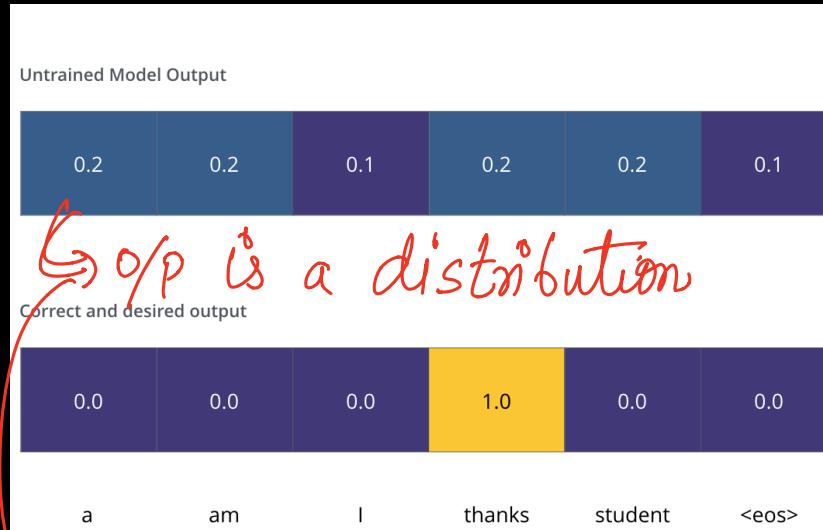
\* Decoder O/P : vector of floats.  $\Rightarrow$  Vector 2  
Word is required for getting O/P word.



\* Assuming that model has learned 1000 words as "O/P vocabulary".  
 $\Rightarrow$  logits's size is also 1K.

\* Finally Softmax Converts raw logits scores into probabilities. Highest probability Cell  $\Rightarrow$  O/P word.

## The loss fn:



Distributions can be compared using cross entropy/KL-Div.

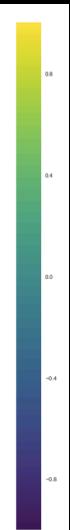
Assuming input

as: je suis

étudiant

expected O/P:

I am a student



Untrained model

initially produce random/in correct O/P.

assuming O/P vocab:  
a/am/I/thanks/  
student/<eos>

### Target Model Outputs

Output Vocabulary: a am I thanks student <eos>

position #1	0.0	0.0	1.0	0.0	0.0	0.0
-------------	-----	-----	-----	-----	-----	-----

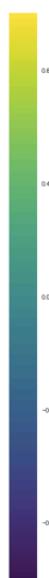
position #2	0.0	1.0	0.0	0.0	0.0	0.0
-------------	-----	-----	-----	-----	-----	-----

position #3	1.0	0.0	0.0	0.0	0.0	0.0
-------------	-----	-----	-----	-----	-----	-----

position #4	0.0	0.0	0.0	0.0	1.0	0.0
-------------	-----	-----	-----	-----	-----	-----

position #5	0.0	0.0	0.0	0.0	0.0	1.0
-------------	-----	-----	-----	-----	-----	-----

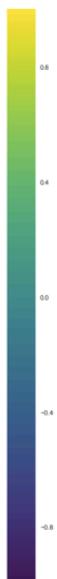
a am I thanks student <eos>



## Trained Model Outputs

Output Vocabulary: 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



After some initial training model may start to converge to the target distribution.

Model is selecting the "best match" highest probable word.

## Beam Search

Instead of choosing Top 1  $\Rightarrow$  Considering the beam-size = 2  
(say "I")

One can consider Top 2 words instead of just 1  
(say "I" & "am")

(For the first O/P position)  $\downarrow$   
we need to run model twice once assuming "I" as the next word and once assuming "am" as the next.

$\Rightarrow$  whichever version produces less error

At any time 2 partial hypotheses (unfinished translation) is kept in memory.

1 2 3 4 5 (O/P Vocab)

beam size  $\geq 2$   
max o/p length = 3

$\Rightarrow$  It will consider all of them and choose the one giving minimum error.

