

## → Natural Language Processing

going towards → transformers (different type of neural network architecture)  
 OR transformer architecture (majorly used for NLP & time series type of data)

contains → encoder & decoder unit

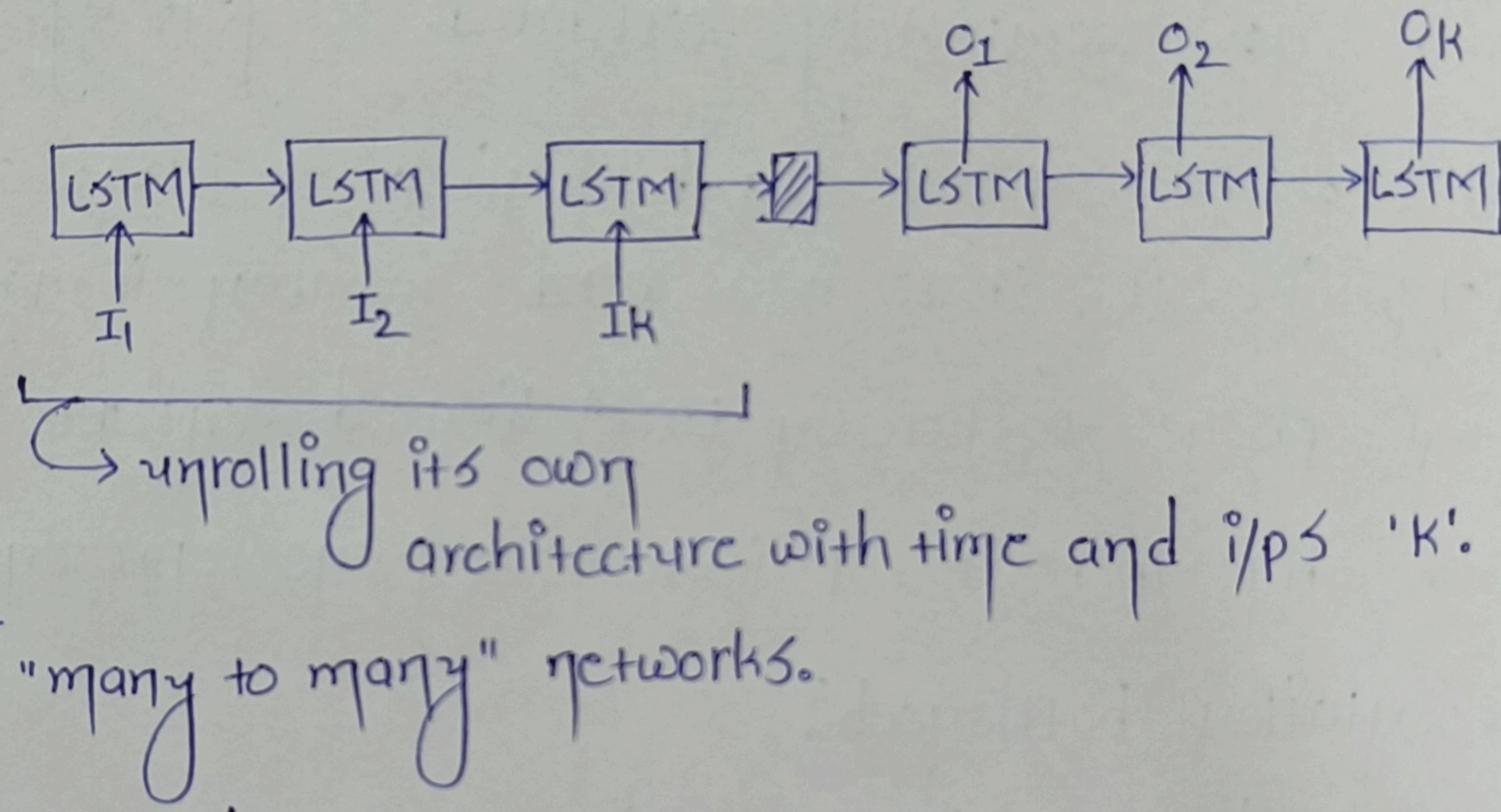
\* → what is latent representation? (see Page-16)

\* → what is <sos>? (see Page-15)

→ Recurrent Neural Network and Long Short Term Network

working model  
of LSTM

Here inputs are many and outputs are also many so these type of networks are also known as "many to many" networks.



→ This specific problem / process known as "sequence to sequence" modelling. It is basically used for language translation. e.g. Hindi to

English, German to Gujarati, and etc.

→ Issues OR problems with LSTM and RNN:

long term dependencies are difficult to manage.

(1) long sentences (not able to give much more efficiency on it)

(2) large amount of memory requirement

(3) only one representation (Representation Bottleneck)

(4) attention or other model utilization

(5) recurrence (due to this it is difficult to solve sentences)

(6) regularization is difficult. (Prone to overfitting)

(7) exploding & vanishing problem, (8) sequential's i/p & o/p

LSTM & GRU are used to overcome this issue.

\* Transformer can be subject.

Page-2

→ Issues (1) and (7) are related to each other.

\* Sir has said; we can solve exploding by using clipping? But how?  
what is it? (see Page-22)

→ LSTM can extract long term dependencies better than RNN. But at certain extent LSTM is also not working. To overcome this issue and want to extract more long term dependencies then more memory is required.

→ LSTM are now completely outdated. (1 year ago)

→ In LSTM stored context can be limited in terms of amount of memory. like if there is a word "it" then it may represent cat, dog, and microphone. So here in LSTM limited context can be stored.

That is why there is an issue "only one representation" OR "representation bottleneck". solution of this issue

→ We also need to apply attention to LSTM due to its limited representation. Using LSTM along with attention gives good representation in terms of context.

→ Major issue with the LSTM and RNN is recurrence / they are sequential in nature.

→ In the LSTM for regularization, Batch Normalization can not be used instead of it Layer Normalization is used. This is also issue with the transformers.

# → Batch Normalization vs. Layer Normalization

Page-3

↳ covariates in data  
OR in batches / covariate shift

→ Difference b/w data  
of various batches

Batch Normalization is not applied over the weights of the layer of neural network. For the regularization of weights of the layer(s) of neural network "L<sub>1</sub> & L<sub>2</sub> regularization" can be applied.

Handling the covariate shift is all about the batch normalization.

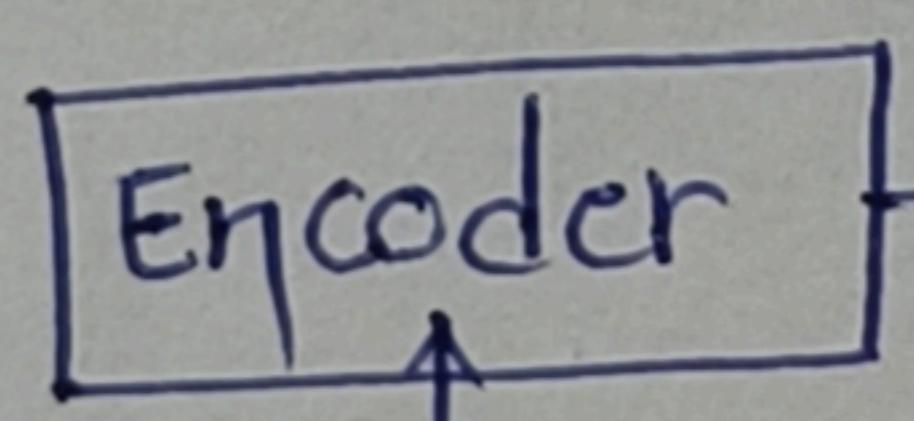
\* → what is the truncated backpropagation in time? (See Page-20)

→ Natural language Processing ≈ Transformer (Nowadays Reality)

→ Landmark paper: "Attention is all you need"

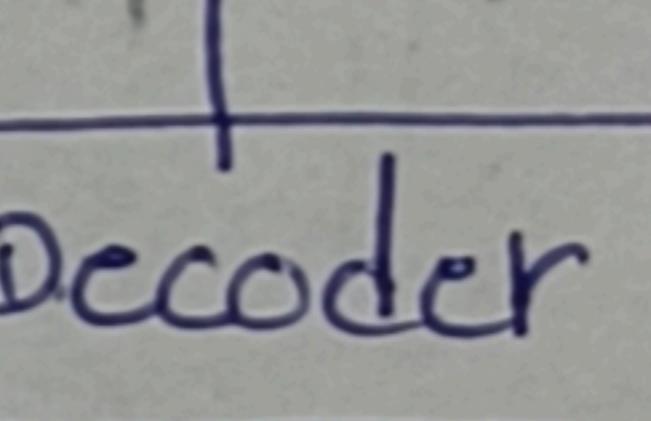
Transformer Abstract View:

Encoder is learning  
based on self-attention.



Representation

Je vais bien (French)



Decoder is learning based  
on self-attention and  
also based on cross-atten-  
tion.

GPT  
Generalized Pre-trained  
Transformer

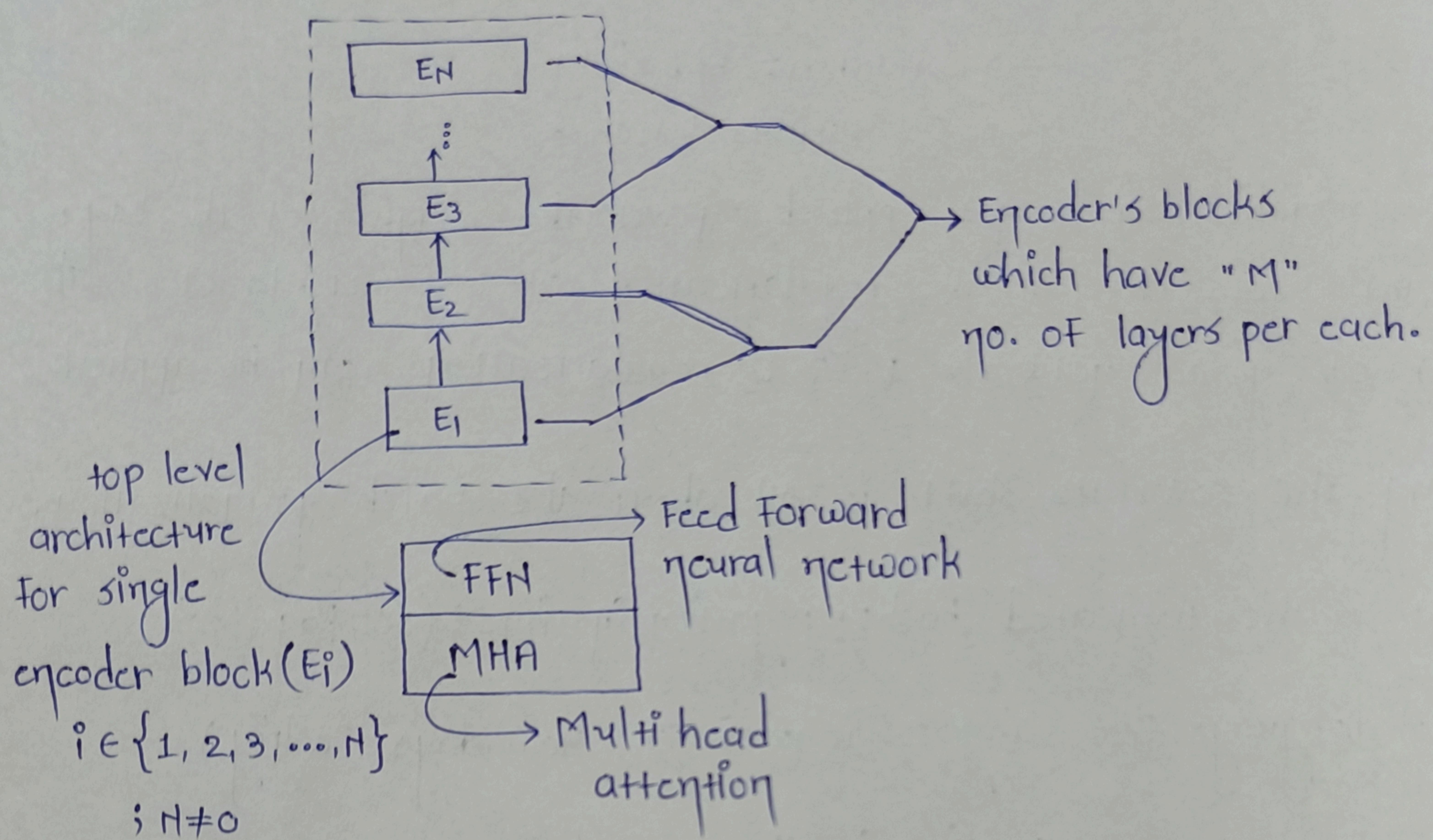
{ Bidirectional Encoder  
Representation from  
Transformer }

\* → what is difference b/w Encoder architecture in transformer and  
auto encoder architecture?

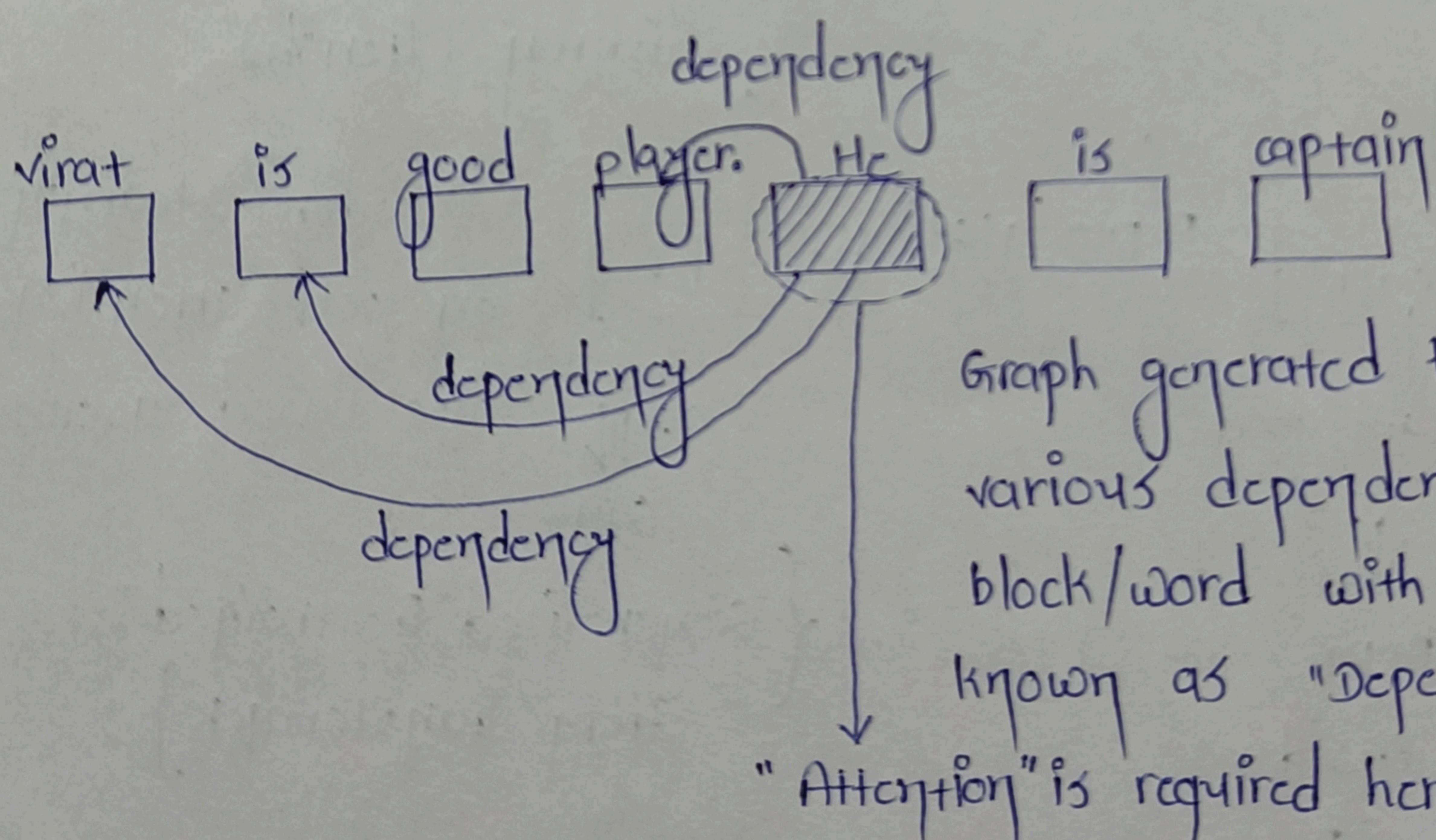
→ Encoder explanation:

Page-4

Use "top level architecture" instead of "abstract view" word.



consider a scenario,



i.c.<sub>o</sub>(2)

A dog ate the food because it was hungry.

A dog ate the food because it was hungry.

given attention to  
each & every word  
within same  
sentence known  
as "self-attention".

→ An attention given to each & every word of OTHER sentences known as "cross - attention."

c.g. (1)      Representation of "I" in i/p matrix will be  $X_I$

I	$x_{11}$ 1.76	$x_{12}$ 2.22	$x_{13}$ 3.79	...	$x_{1,512}$ 6.66
am	$x_{21}$ 7.77	$x_{22}$ 0.631	$x_{23}$ 9.86	...	$x_{2,512}$ 5.35
good	$x_{31}$ 11.44	$x_{32}$ 10.10	$x_{33}$ 19.04	...	$x_{3,512}$ 3.33

Input matrix also known as "Embedding Matrix" represented using " $X$ ".

\* → How to generate vector for given specific word?

and the dimension of that vector is  $(1, 512)$ ?

\* → Is Word2Vec is embedding methods? (YES) "Attention is all you need" (see paper & glove)

Each and every is being seen using 3 rules:

consider single word vector of I,

$$I \Rightarrow [x_{11} \ x_{12} \ x_{13} \ \dots \ x_{1,512}]$$

→ Query; if you apply Query() Function on "I" i.e. Query(I) then it will return " $X_I^Q$ ".

→ Key; if you apply Key() Function on "I" i.e. Key(I) then it

value will return " $X_I^K$ ".

; represented using

" $Val_I$ "

→ Additional:

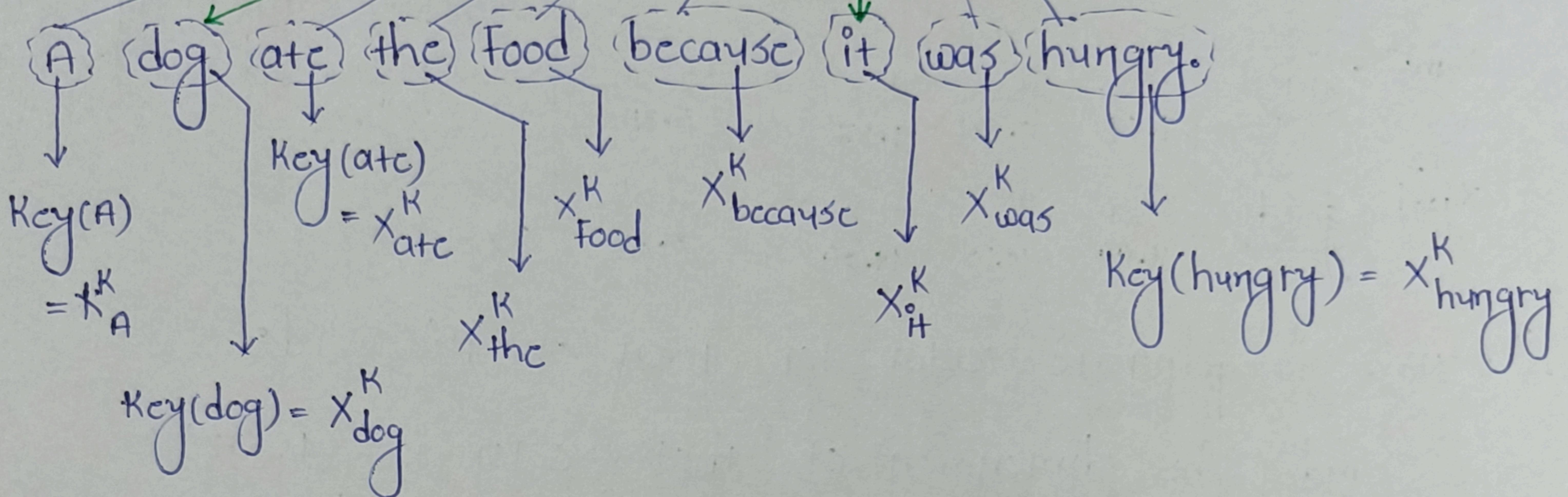
An embedding is a vector representation of word or phrase. It is a way of representing text data in a way that is understandable by machines. Typically a fixed-length vector, such as 300 or 1024 dimensions.

cog. (2)

A dog ate the food because it was hungry.

Query( $it$ ) will be

applied for the main word.



Here in given example, it is trying to find attention value to each & every word within same sentence with respect to that main word, here it is ["it"]. For the main Query value will be found and Key value will be found for opposite compared words.

\* → what is word embedding? & what is "without context" concept?  
(see Page 20)

And to calculate the attention score b/w main word "it" and other word "because" we will do dot product of query representation of "it" i.e.  $x_{it}^Q$  and key representation of "because" i.e.  $x_{because}^K$ .

\* → Here question is why only b/w query & key values?  
and why only dot product based attention score?  
(see paper "Attention is all you need")

→ e.g. (3)

consider a sentence "I am good"

representation will be,

$$x_1 \Leftarrow x_I = [x_{11} \ x_{12} \ x_{13} \ \dots \ x_{1,512}]_{1 \times 512}$$

$$x_2 \Leftarrow x_{\text{am}} = [x_{21} \ x_{22} \ x_{23} \ \dots \ x_{2,512}]_{1 \times 512}$$

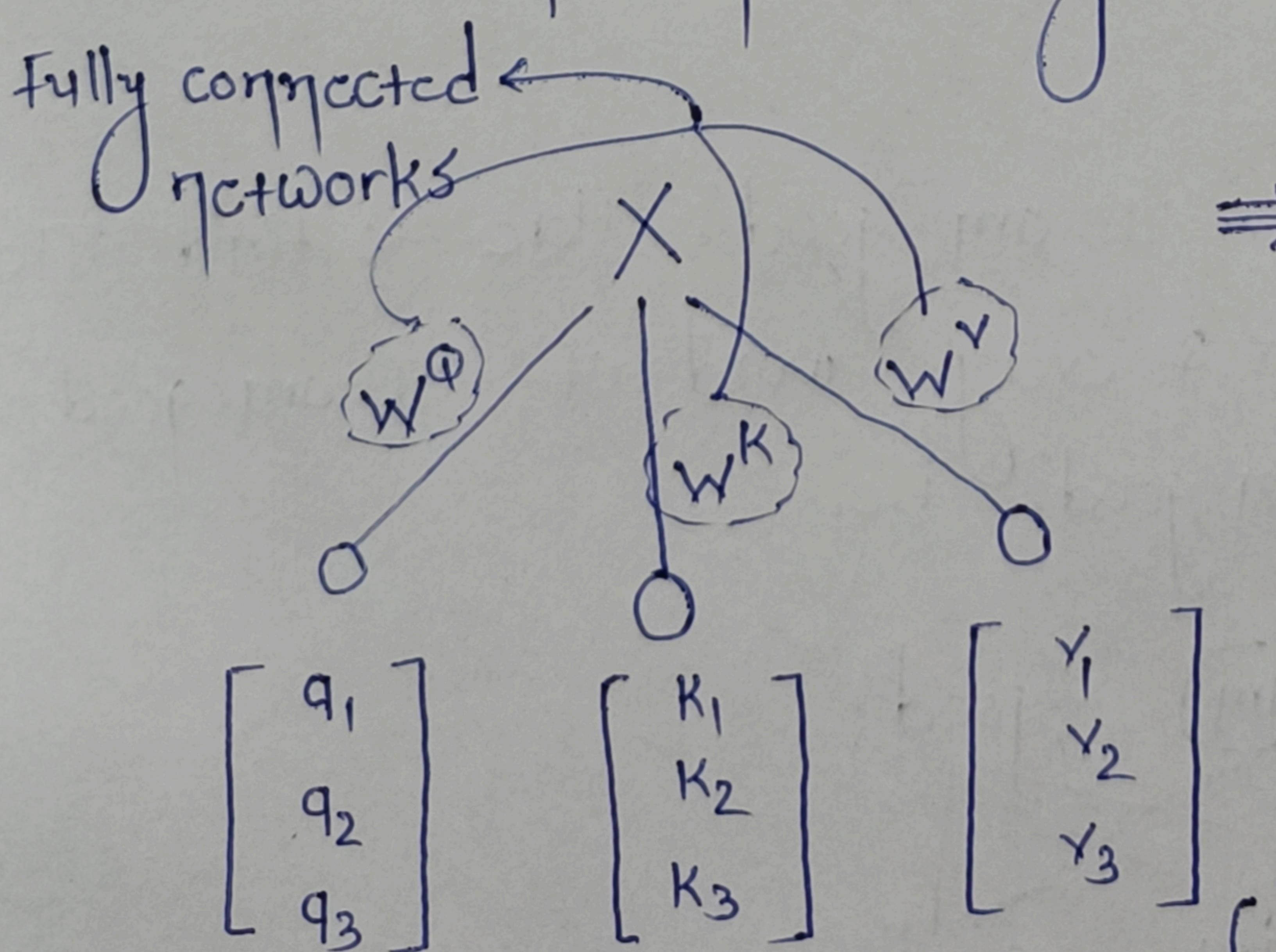
$$x_3 \Leftarrow x_{\text{good}} = [x_{31} \ x_{32} \ x_{33} \ \dots \ x_{3,512}]_{1 \times 512}$$

$[x_1]_{512 \times 1}$ ,  $[x_2]_{512 \times 1}$ , and  $[x_3]_{512 \times 1}$

column matrix of all will be input matrix,

input matrix/vector  $x = \begin{bmatrix} x_I \\ x_{\text{am}} \\ x_{\text{good}} \end{bmatrix} \sim \begin{bmatrix} x_1 \\ x_2 \\ x_3 \end{bmatrix}$

For this input matrix three different neural networks (Fully connected networks) will be trained. One for query, one for key, and one for value. Those neural networks are represented using  $w^Q$ ,  $w^K$ , and  $w^V$  accordingly.



$\Rightarrow [x_1]$  from  $X$  (input matrix) is given

$\Rightarrow$  to the  $w^Q$ ,  $w^K$ , and  $w^V$  to get  $q_1$ ,  $k_1$ , and  $v_1$  accordingly. So if

dimension of  $x_1$  is  $(512 \times 1)$  and  
 " " "  $q_1$  is  $(64 \times 1)$  then  
 " " "  $w^Q$  will be  $(64 \times 512)$

This is similar for remaining i/p vectors ( $x_2$  and  $x_3$ ) and  $w^K$  and  $w^V$ , & their outputs ( $k_i$  &  $v_i$ );  $i \in \{1, 2, 3\}$

Because there is multiplication b/w  $w^Q$  and input vector  $x_1$ .

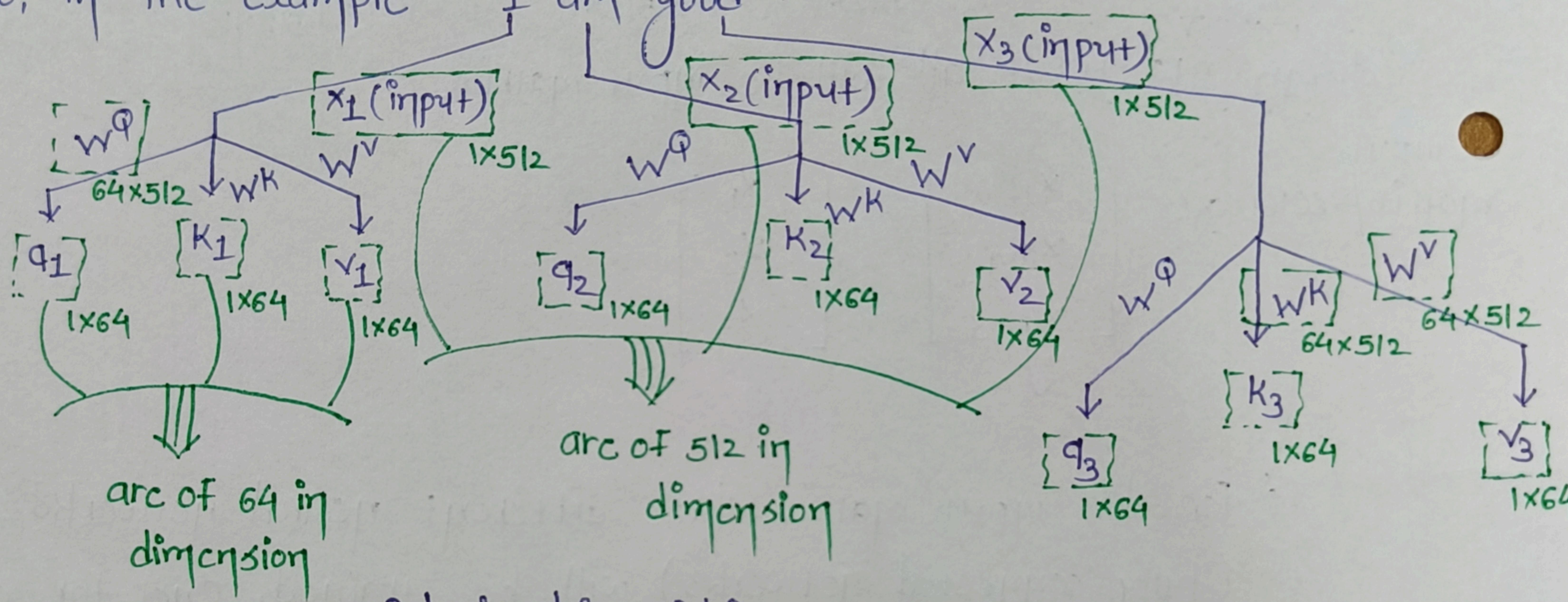
$$\therefore w^Q * x_1 = q_1$$

→  $\begin{bmatrix} q_1 \\ q_2 \\ q_3 \end{bmatrix}$ ,  $\begin{bmatrix} K_1 \\ K_2 \\ K_3 \end{bmatrix}$ , and  $\begin{bmatrix} v_1 \\ v_2 \\ v_3 \end{bmatrix}$  are represented by capital letters

$Q$ ,  $K$ , and  $V$ .

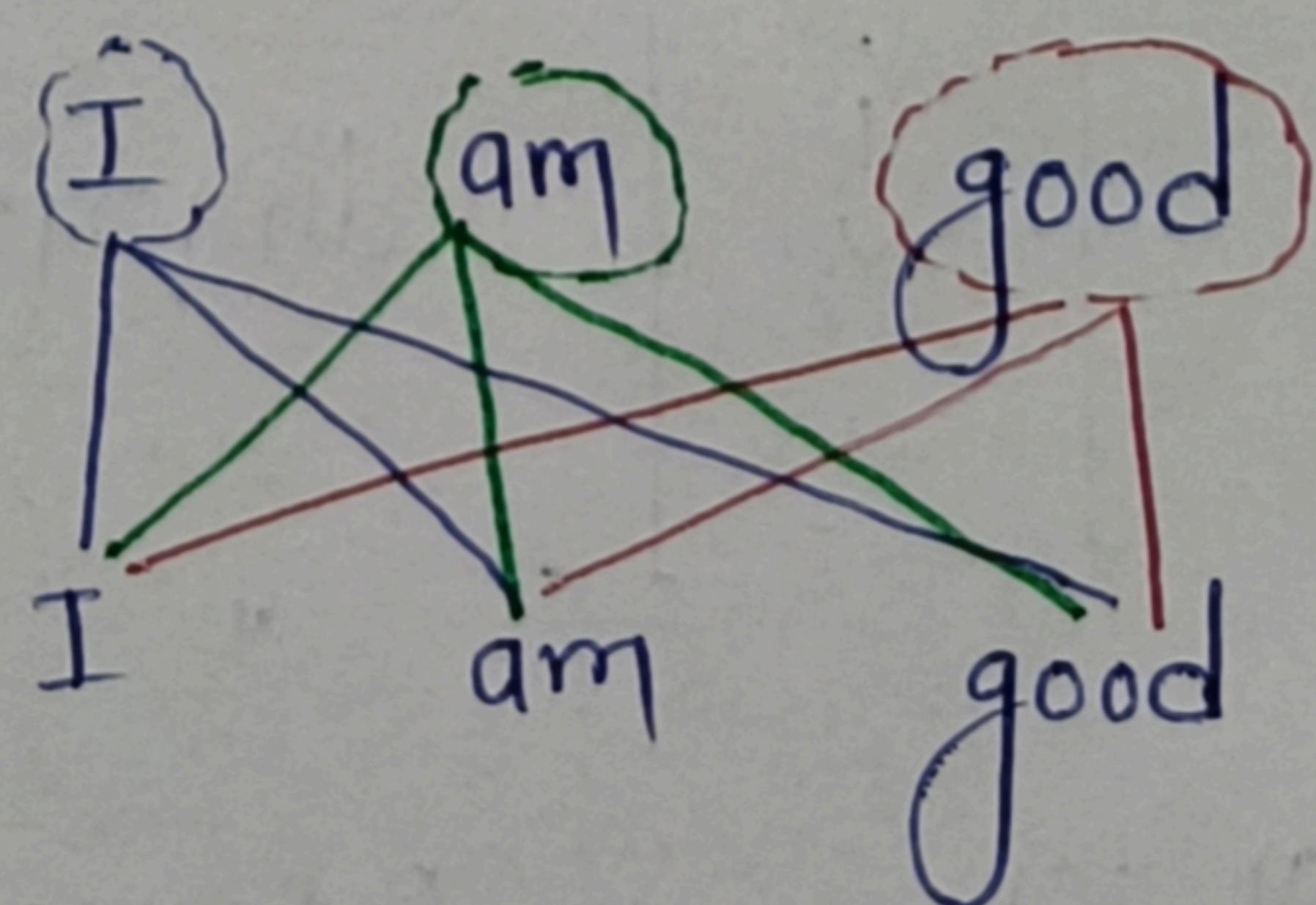
pronunciation of all is  $Q$ -matrix,  $K$ -matrix, and  $V$ -matrix respectively.

So, in the example "I am good".



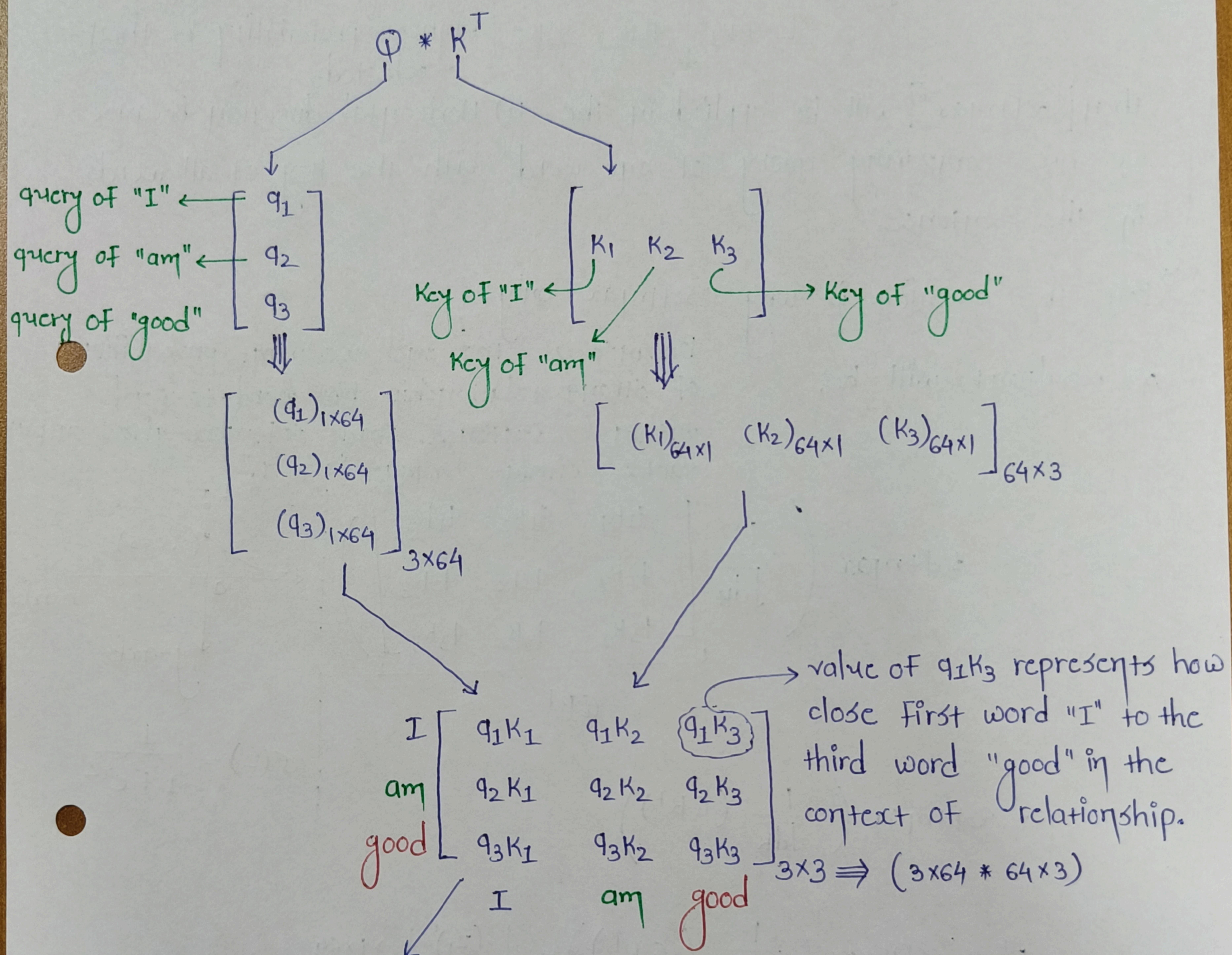
→ Now, self-attention scores are calculated based on  $Q$ -matrix and  $K$ -matrix. That is why here  $Q$ -matrix and  $K$ -matrix are calculated.

now consider same example, "I am good". Here to find attention score / relation value of each & every word of "I am good" with the same sentence "I am good", i.e.



so, those attention scores will be found based on matrix generated by matrix multiplication of Q-matrix and transpose of K-matrix, i.e.

Page-9

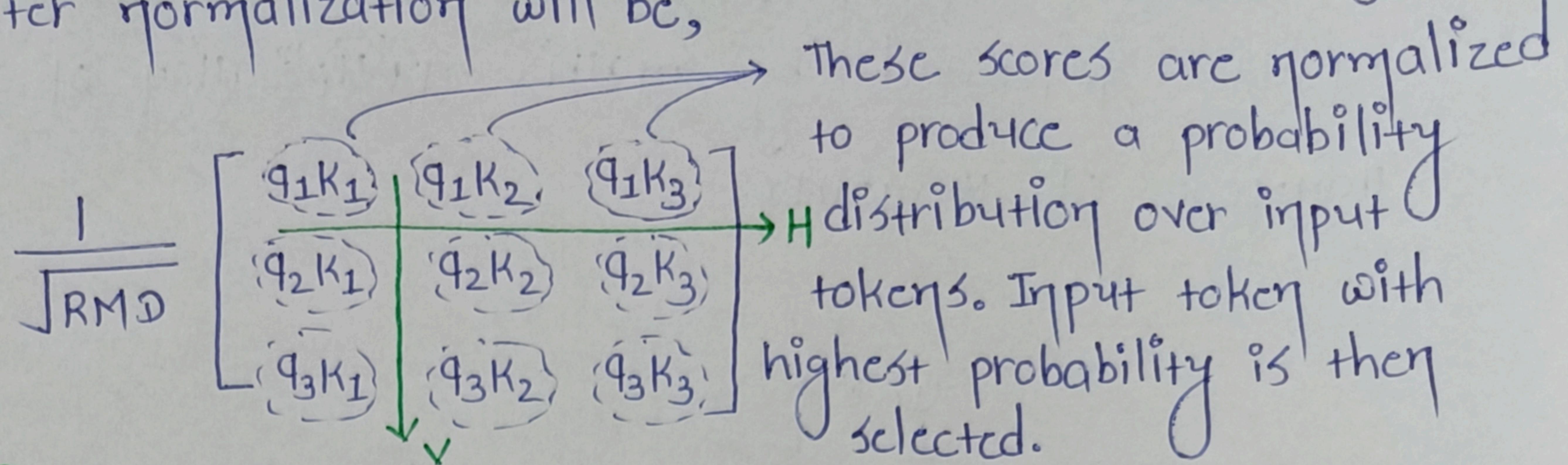


Based on this matrix, we attention given to specific word by specific word can be known. size of the matrix (resultant of  $QKT$ ) is based on no. of words particular sentence contains.

→ Now, it is important to normalize the matrix values. So for that resultant matrix of  $QKT$  is multiplied by  $\frac{1}{\sqrt{RMD}}$ ; RMD = Resultant Matrix Dimension

of particular query or key vector. i.e. 64

so matrix after normalization will be,



then **"softmax"** will be applied in the (H) Horizontal direction because we are comparing query of one word with the key of all words in the sentence.

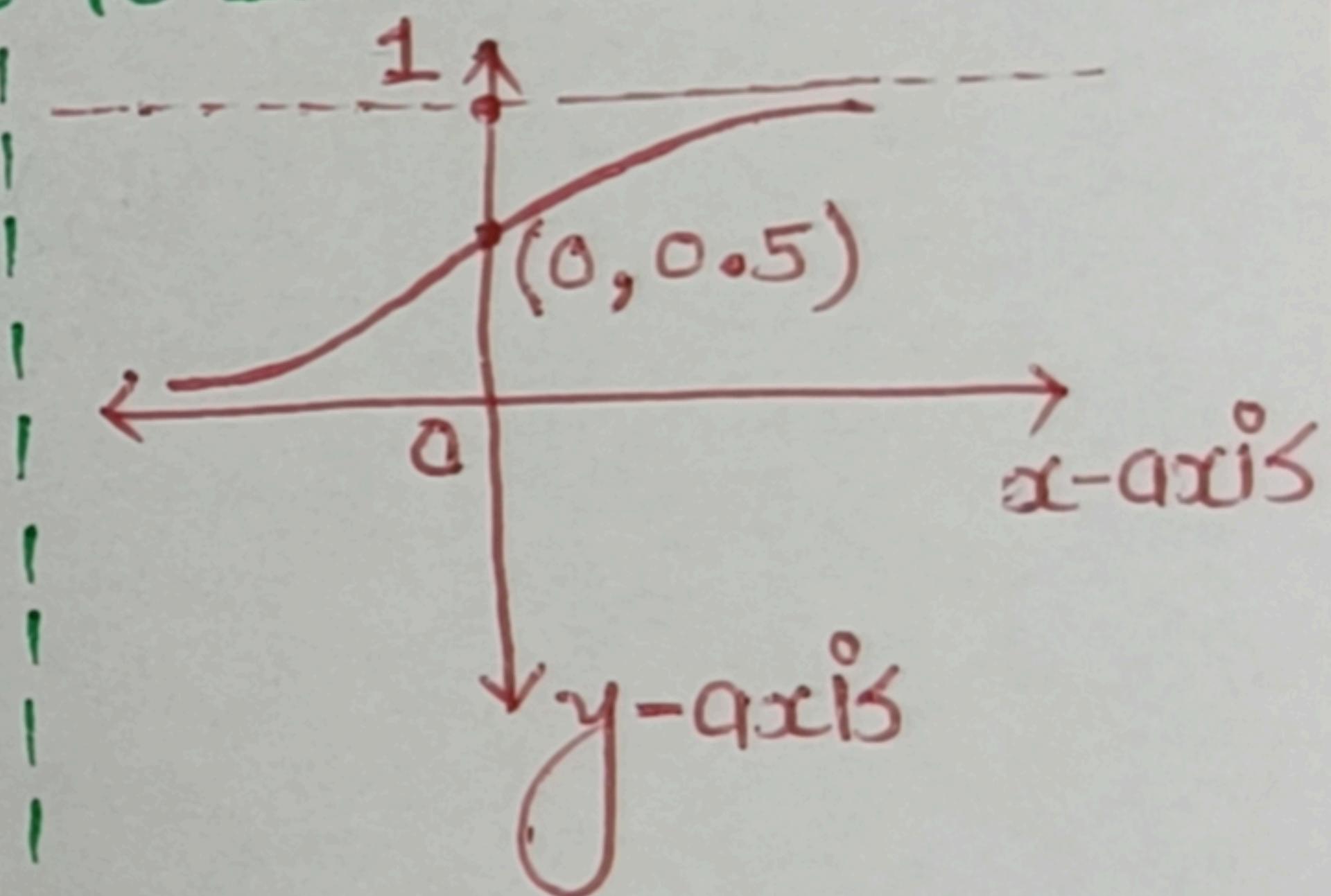
But the question is why "softmax" only?

so, resultant will be,

Because at the end we want probability of strong relationship b/w word-to-word within sentence. And softmax gives output values within range "0 TO 1."

$$\text{softmax} \left( \frac{1}{\sqrt{d_K}} \begin{bmatrix} q_1 k_1 & q_1 k_2 & q_1 k_3 \\ q_2 k_1 & q_2 k_2 & q_2 k_3 \\ q_3 k_1 & q_3 k_2 & q_3 k_3 \end{bmatrix} \right)$$

$\Phi K^T$



$$\sigma(z) = \frac{1}{1 + e^{-z}}$$

$$\sim \text{softmax} \left( \frac{1}{\sqrt{d_K}} (\Phi K^T) \right)$$

$$z' = \text{softmax} \left( \frac{1}{\sqrt{d_K}} (\Phi K^T) \right) \# \begin{bmatrix} (z_1)' \\ (z_2)' \\ (z_3)' \end{bmatrix} \rightarrow 3 \times 64$$

$$z = z' * \checkmark \begin{bmatrix} v_1 \rightarrow I \\ v_2 \rightarrow am \\ v_3 \rightarrow good \end{bmatrix} = \begin{bmatrix} (z_1)' \\ (z_2)' \\ (z_3)' \end{bmatrix} \rightarrow 1 \times 64$$

output matrix is of " $\eta \times 64$ " in dimension, here  $\eta$  is total no. of words in sentence. Here  $\checkmark$  represents value-matrix, typically size of value-matrix equals to the key-matrix.

→ Additional:

The input and output of both functions (sigmoid & softmax) are slightly different since sigmoid receives just one input and only outputs a single number that represents the probability of belonging to class 1 (remember that we only have 2 classes so the probability of belonging to class 2 =  $1 - P(\text{class 1})$ ). While on the other hand softmax is vectorized, meaning that takes a vector with the same number of entries as classes we have and outputs another vector where each component represents the probability of belonging to that class.

Sigmoid

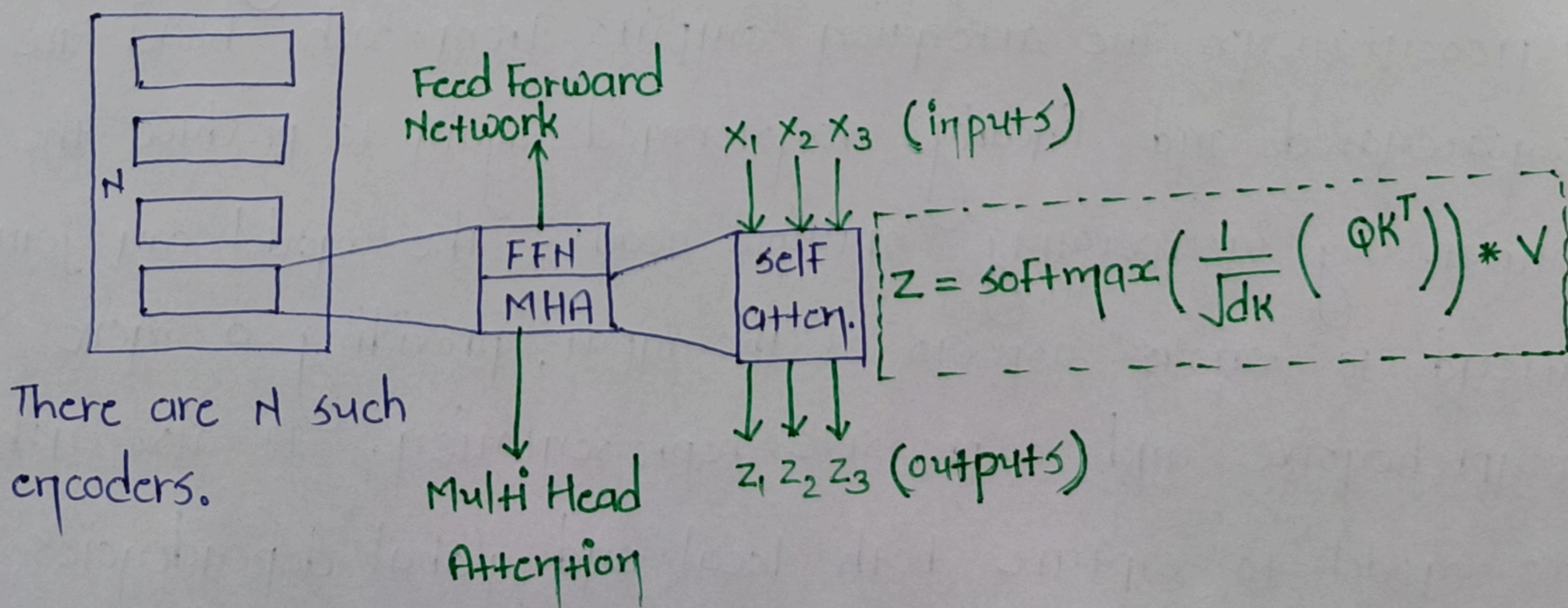
$$\text{out} = P(Y = \text{class 1} | x)$$

$$; x = \text{input} \in (-\infty, +\infty)$$

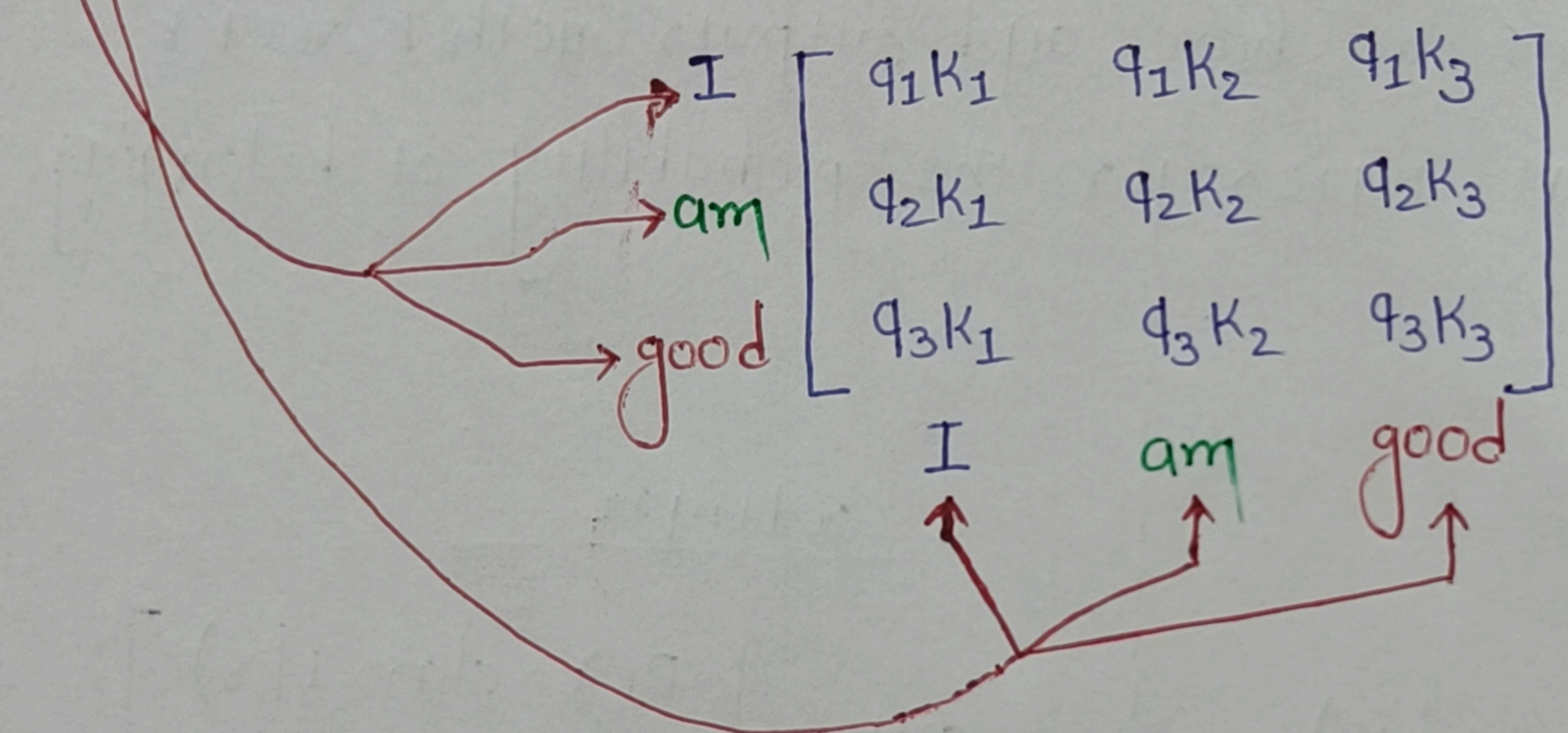
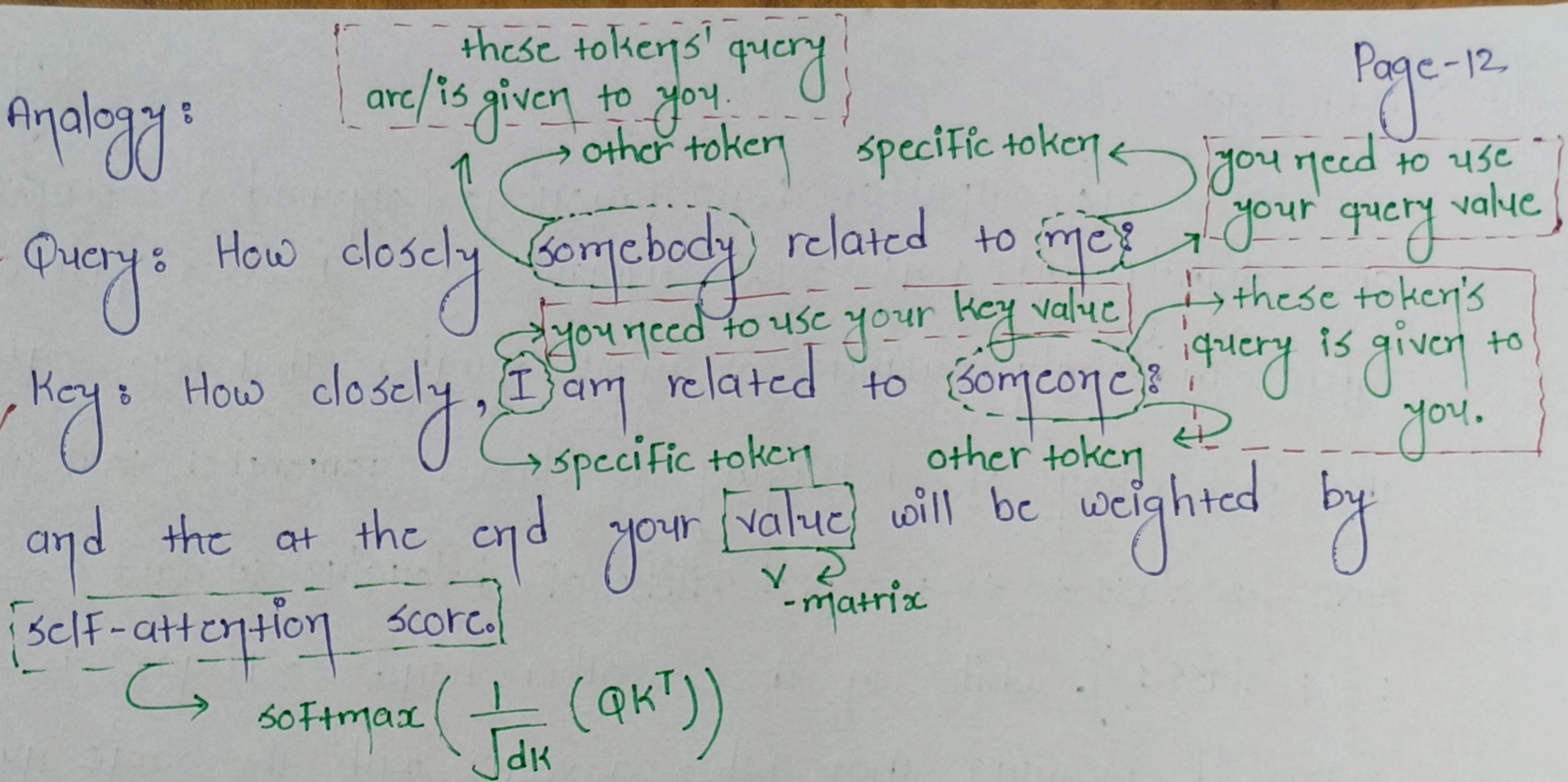
Softmax

$$\text{out} = \begin{bmatrix} P(Y = \text{class 1} | x) \\ P(Y = \text{class 2} | x) \\ \vdots \\ P(Y = \text{class K} | x) \end{bmatrix}$$

→ Now, summarizing encoder part with an example "I am good".



→ Analogy:



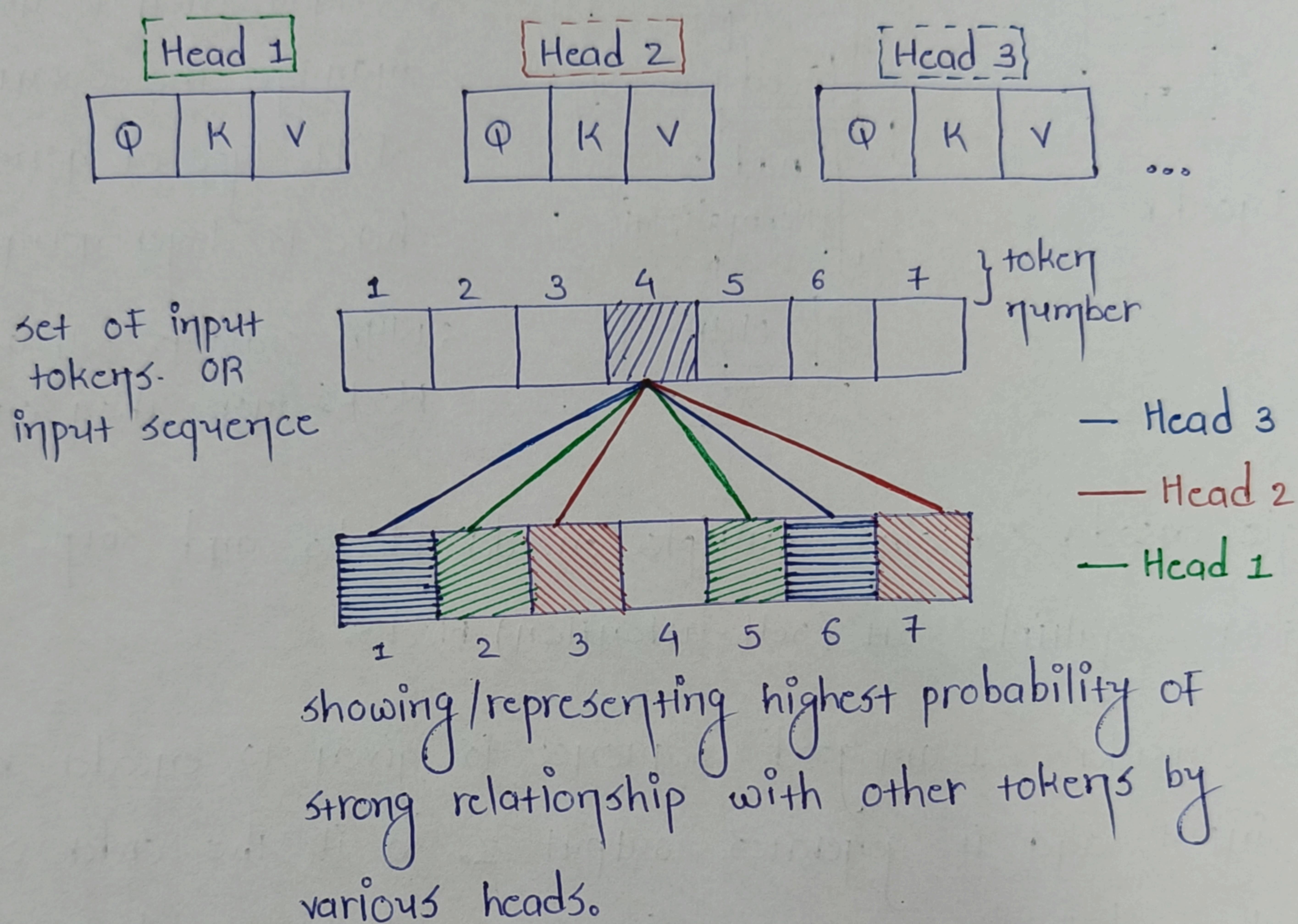
→ Multi-head mechanism:

Typically, multi-head is used to include various aspects or contexts in the training from the input sequences. In the multi-head mechanism, each head learns its own set of query, key, and value representations. The attention outputs from each head are then concatenated and linearly transformed again to produce the final output. By incorporating multiple heads, the model can jointly attend to various aspects of the input, providing a more comprehensive and expressive representation. It also enables the model to capture both local and global dependencies, as different heads can attend to different ranges of the input sequence.

All the attention heads in the multi-head mechanism of transformer are working simultaneously to capture multiple perspectives and contexts from the input sequence.

Page-13

graphical representation:

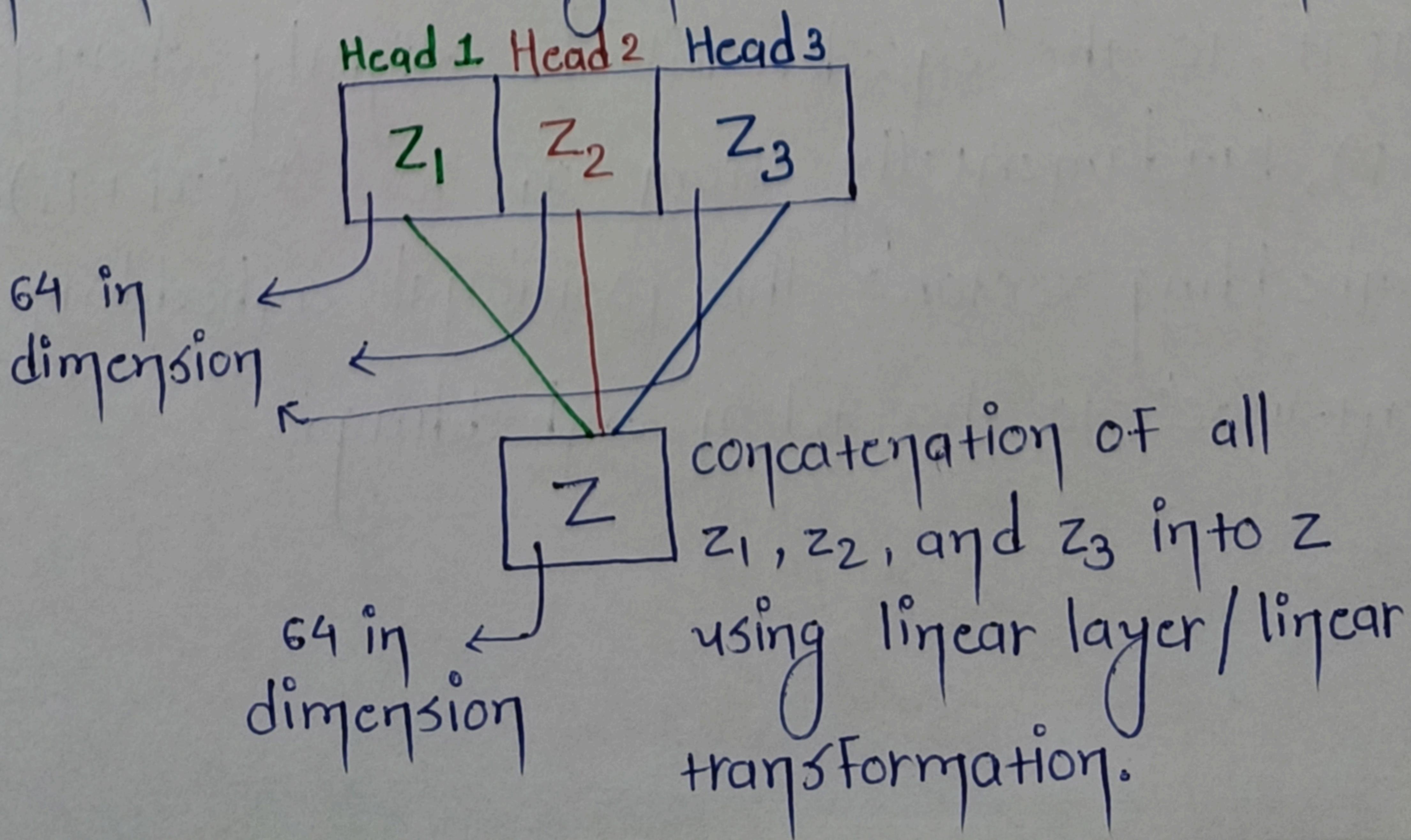


For token 4, Head 1 gives high attention to token 1 & 6.

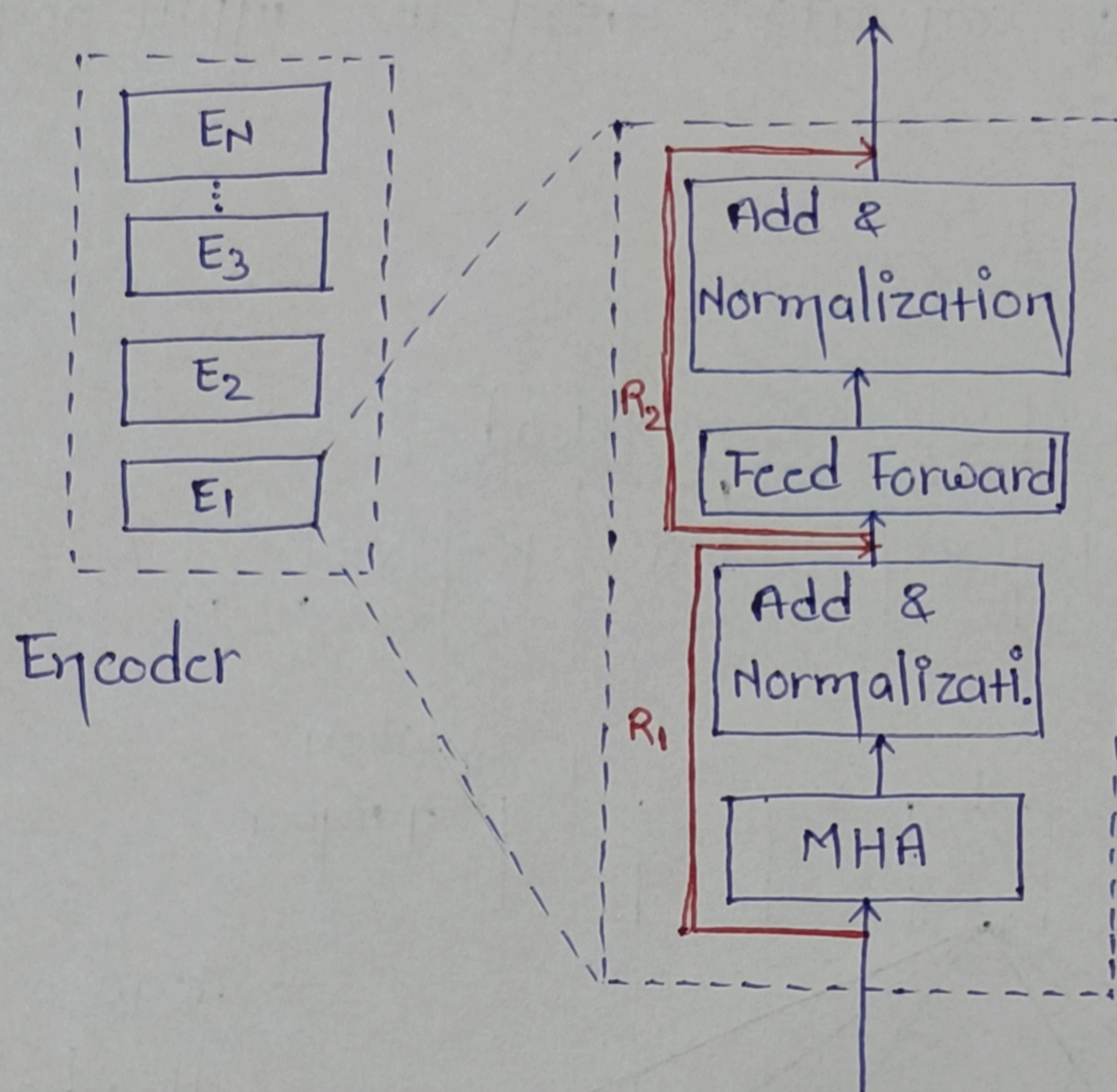
Head 2 " " " " token 3 & 7.

Head 3 " " " " token 1 & 6.

All attention heads will generate z-matrices independently.



→ At last, drawing top-level diagram of single encoder block.



- ;  $R_1$  and  $R_2$  are residual connections.
- ; Normalization is used to manage the covariate shift. Type of normalization here is layer normalization
- ; MHA contains self-attention blocks (more than one.).

single encoder contains multiple encoder blocks and single MHA contains multiple SA (self-attention) blocks:

→ Now, consider "I am good" sequence is given to encoder block as an input and it generates output  $z$ . So if the order of tokens is changed in the sequence (sequence will be "good am I") and again give it to encoder then also output will be same i.e.  $z$ . This happens due to "Positional Embedding". {output of encoder is "position invariant" in nature.}

Now, the input to the encoder will be slightly changed from  $x$  to  $(x + p)$ , fundamentally from  $x_i$  to  $(x_i + p_i)$ . Here  $p_i$  is positional embedding vector. This positional embedding vector is added element-wise to the token embeddings.

Recurrent Neural Network have an inbuilt mechanism that deals with the order of sequences. The transformer model, does not use recurrence or convolution and treats each data point as independent of the other. Hence, positional information is added to the model explicitly to retain the information regarding the order of words in a sentence. ["Positional Encoding"] is the scheme through which the knowledge of the order of objects in a sequence is maintained.

- \* → It can be learnable/trainable parameter. But this is not a case for token embedding also learnable parameter? all time.

Both [token embeddings] and [positional embeddings] contribute to the overall input representation in the transformer model. These embeddings allow the model to capture relationships between words and the order of tokens, enabling effective modeling of context and semantic information.

- \* → what is SOS and EOS in decoder side?

Start of Sequence      End of Sequence

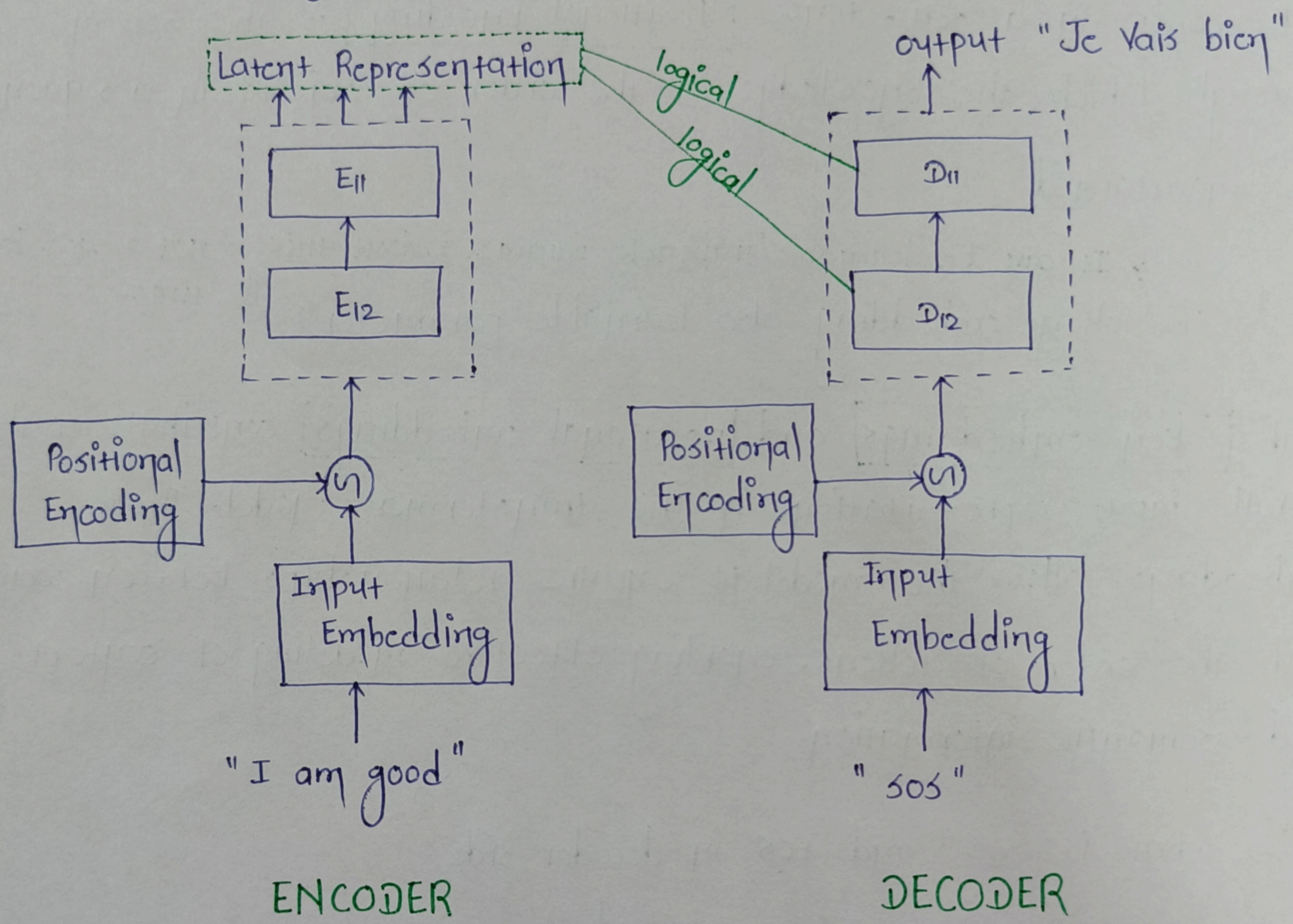
it is a special token or symbol used in sequence-to-sequence models (RNNs or transformer) to indicate the beginning of sequence. The "SOS" token is prepended to the input sequence when generating the output, serving as an indicator for the decoder to start generating the target sequence.

When decoder encounters the "SOS" token, it initializes its internal states and begins generating the output tokens one by one until it reaches an EOS token or a predefined maximum length.

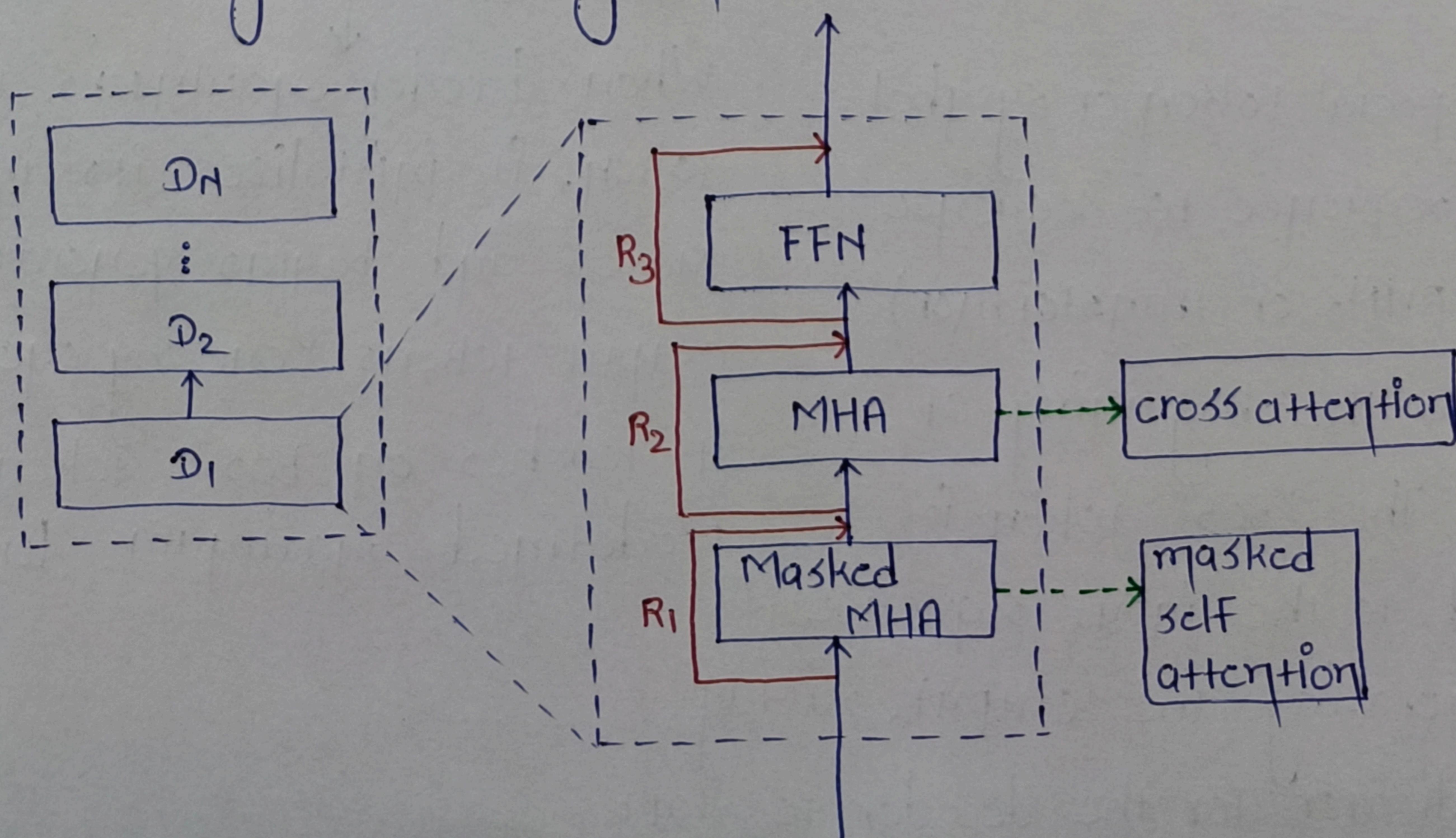
→ Latent representation:

↳ it refers to the intermediate representation of input data at a particular layer within the model.

→ Top-level diagram of transformers:

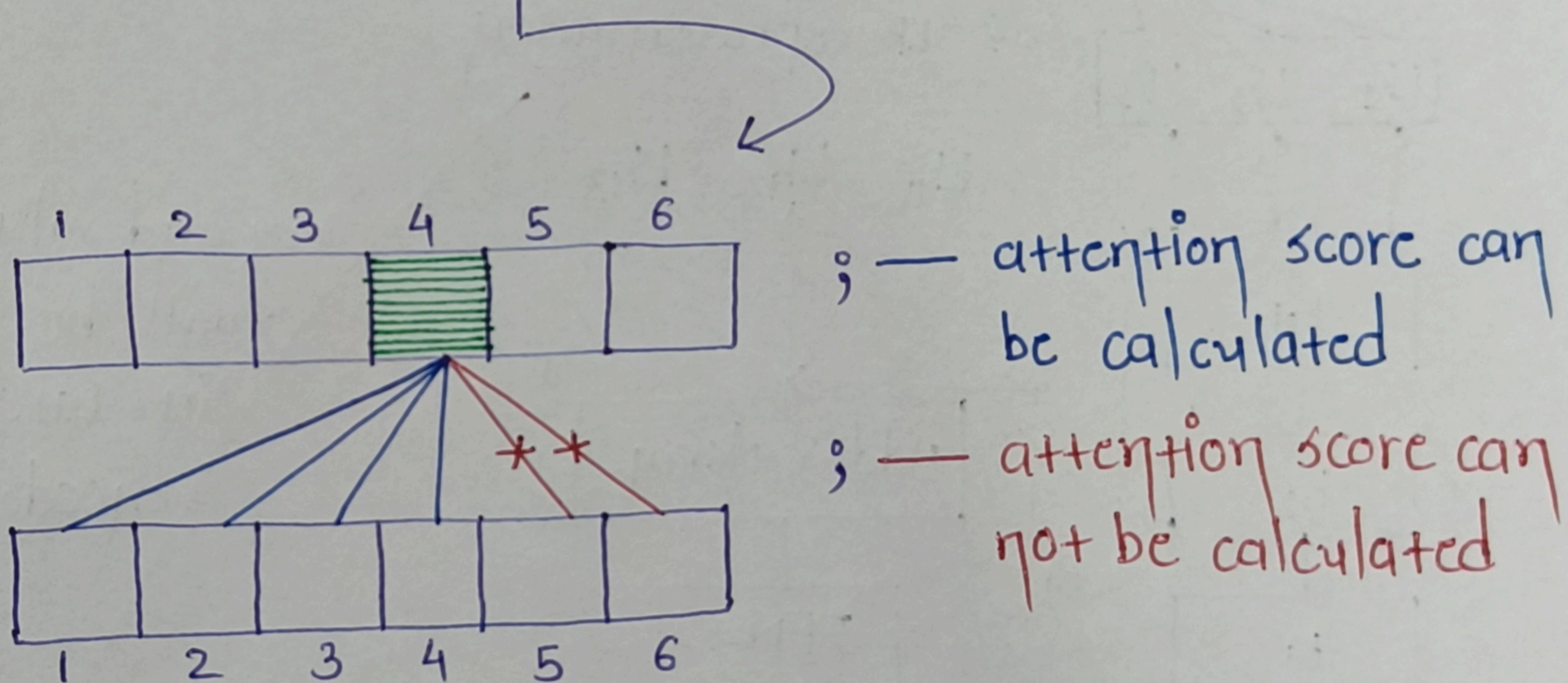


→ Decoder single block diagram:



→ Here in the decoder side output tokens of sequence will not be generated parallel. But encoder side "I am good" sequence given as an input and its all token processed parallel. French translation of "I am good" is "Je vais bien", so decoder side first "Je" token will be generated, then it will feed as an input to decoder again, and decoder will predict next word i.e. "vais". Again "vais" is will feed as an input to decoder and it will predict next word i.e. "bien". Here based on next word prediction capability of decoder, its performance will measure. To achieve this type of mechanism "MASKED SELF ATTENTION" is used.

consider  
input sequence  
with 6 tokens



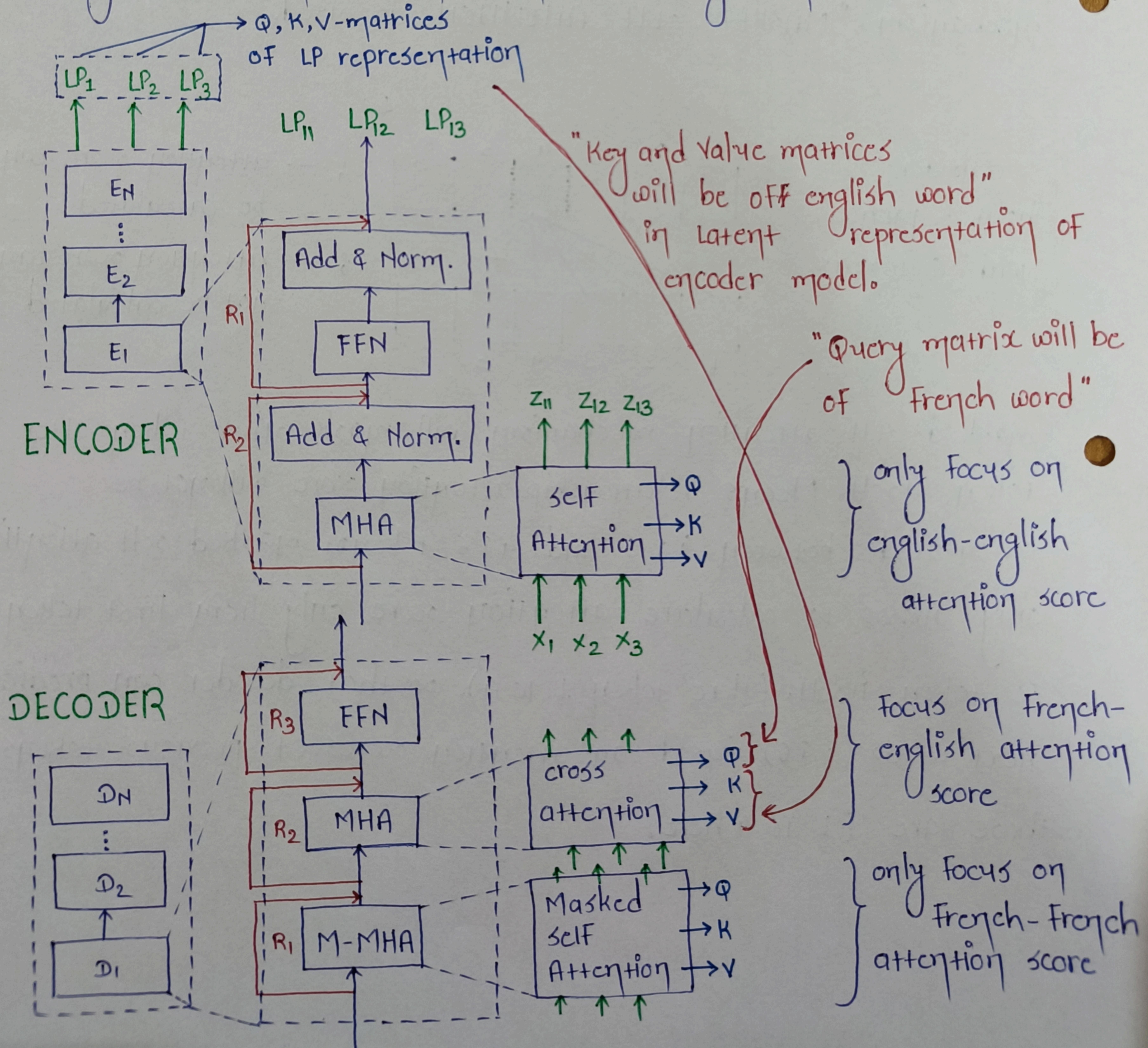
masked self attention mechanism will mask other tokens after token itself. Means relationship/attention score will not be calculated between 4 & 5 and 4 & 6 tokens, masked self attention only allows to calculate attention score only from first token to token itself. (Here token 1 to 4). So that decoder can predict next words (5 & 6) based on attention scores of previous tokens those are {1 to 4 here.}

## → cross attention:

Here there is an example ["I am good."] so whenever translation will occur from English to French i.e. ["Je Vais Bien"] then in the cross attention there is a relationship between French word query and English word key at decoder side.

and value also

Cross-attention score is calculated in same way as self-attention score is calculated. But the difference is only with the Q-matrix and V-matrix. Let's consider an example of English to French translation. (specifically consider at decoder side.)



in the cross-attention, it tries to learn relationship between French query and English key and value matrices.

In the cross-attention French to English mapping is there at decoder side because cross-attention tries to learn that [how much I am close to given English representation] i.e. latent representation of encoder.

→ Cross-attention also known as encoder-decoder attention.

→ Final output of decoder of input ["I am good"] to encoder will be

- ["Je Vais Bien"] input "I am good" is also represented by  $x$ . So, as seen earlier  $x$  will be,

$$x = \begin{bmatrix} x_I \\ x_{am} \\ x_{good} \end{bmatrix}; x_{\{I, am, good\}} \text{ is word embedding of dimensions } 512.$$

Let's consider output "Je Vais Bien" represented by  $y$ . So  $y$  will be,

$$y = \begin{bmatrix} y_{je} \\ y_{vais} \\ y_{bien} \end{bmatrix}; y_{\{je, vais, bien\}} \text{ is also word embedding of dimensions } 512.$$

(1) what is word embedding?

word embedding refers to a technique used to represent words in a continuous vector space. Main idea behind word embedding is to capture the semantic and syntactic relationships b/w words based on their context in a given corpus. Word2Vec, Glove, and FastText all are algorithms for word embedding.

(2) what is token embedding?

token embedding is a broader term that encompasses the representation of any individual token in a sequence of text. A token can be word, subword, character or any other unit that the model operates on. Token embedding captures the meaning and context of a token within a specific sequence.

Additional:

Word embedding is a specific type of token embedding that focuses on representing individual words.

(3) what is truncated backpropagation through time (TBPTT)?

It is a technique used in recurrent neural networks (RNNs) to address the vanishing or exploding gradient problem during training. In standard backpropagation, gradients are computed and propagated through the entire sequence of inputs and hidden states, which can be computationally expensive and can lead to the vanishing or exploding gradient problem. To mitigate this issues, truncated

backpropagation through time breaks the sequence into smaller subsequences or chunks. Instead of propagating gradients through the entire sequence, TBPTT limits the no. of time steps over which gradients are computed. The sequence is divided into chunks of a fixed length, and gradients are computed and updated within each chunk. The hidden state at the end of each chunk is then used as the initial state for the next chunk. By truncating the backpropagation through time, TBPTT reduces the computational burden and addresses the vanishing or exploding gradient problem. However, it also introduces a limitation in capturing long-term dependencies that span across multiple chunks. The length of the chunks is typically chosen based on a trade-off between computational efficiency and model's ability to capture long-term dependencies.

(4) What is vanishing gradient and exploding gradient problem?

The vanishing gradient problem occurs when the gradients shrink exponentially over time, making it difficult for the model to learn long-term dependencies. The exploding gradient problem occurs when the gradients grow exponentially, leading to instability during training.

(5) what is gradient clipping?

It is a commonly used technique to overcome exploding gradient problem in neural networks. Gradient clipping involves capping the magnitude of gradients during training to prevent them from becoming too large. Two commonly used approaches are;

(a) L2-Norm clipping:

if  $\|\text{grad}\| > \text{threshold}$ :  
 $\text{grad} = (\text{threshold} / \|\text{grad}\|) * \text{grad}$   
; where  $\|\text{grad}\|$  represents the L2-Norm of gradient vector.

→ consider a vector,  
Also known as Euclidean  $x = [x_1, x_2, \dots, x_n]$   
↑ Norm  
L2-Norm is calculated as:  
 $\|x\|_2 = \sqrt{(x_1^2 + x_2^2 + \dots + x_n^2)}$

AND

L1-Norm is calculated as:

$$\|x\|_1 = |x_1| + |x_2| + \dots + |x_n|$$

(b) value clipping:

Individual elements of the gradient vector are checked against a predefined threshold. If any element exceeds the threshold, it is set to the threshold value. This method is less commonly used compared to L2-Norm clipping.

$$\text{grad} = [x_1, x_2, x_3, x_4, x_5], \text{ threshold} = y$$

$x_1 > y$     $x_2 > y$     $x_3 > y$     $x_4 > y$     $x_5 > y$