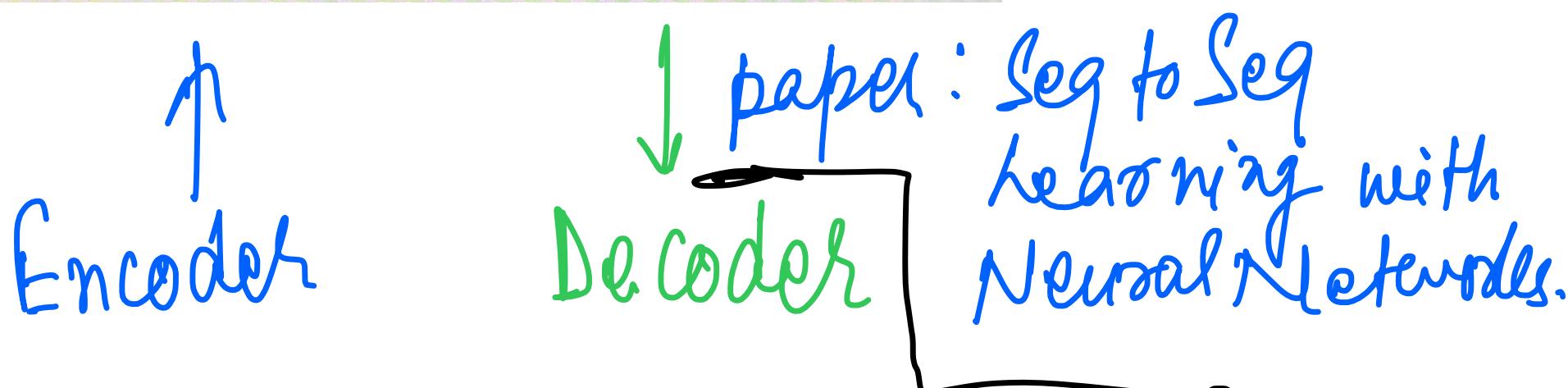
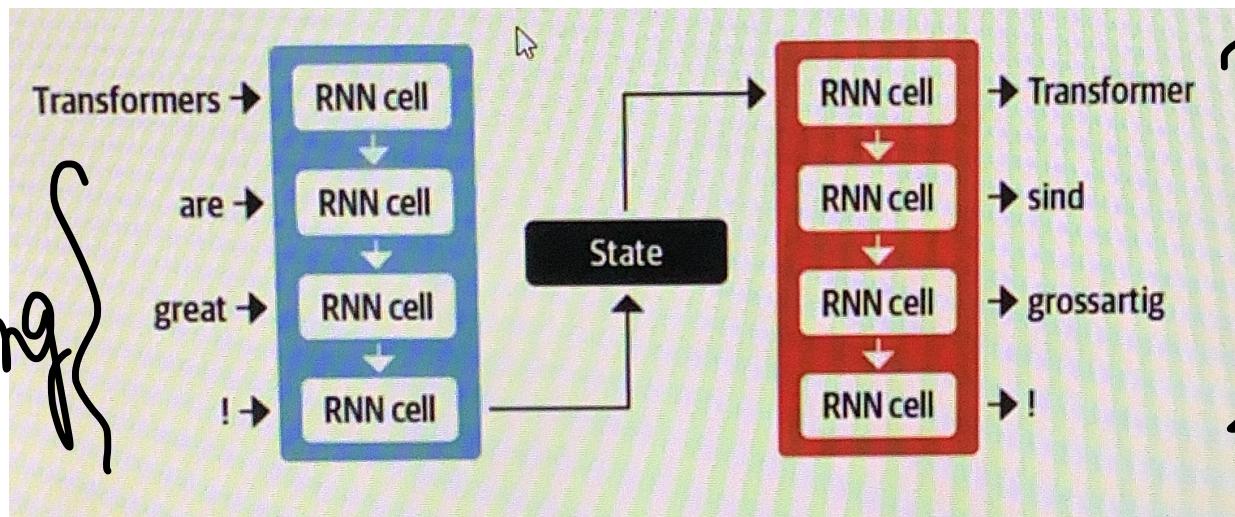
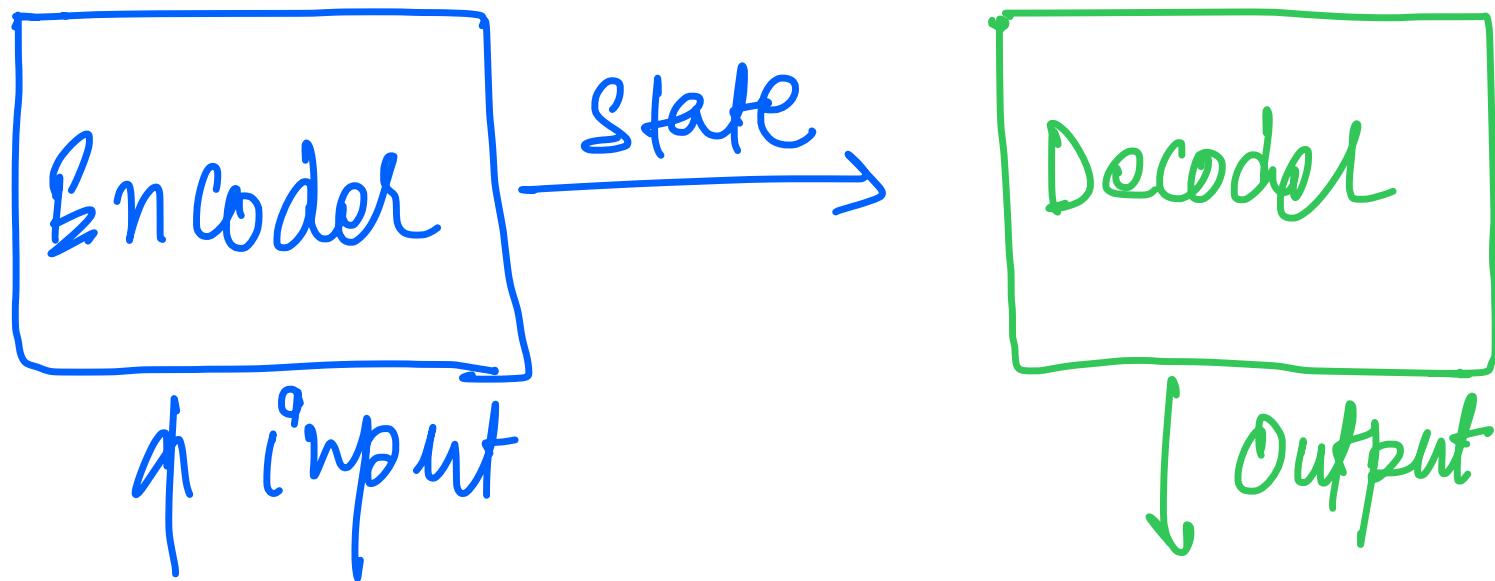


Transformers
Attention Mechanism

Encoder-Decoder-Architecture



Problem with this architecture is it kind of starts missing initial words for long sentences. (30 words)

which leads to "Attention Mechanism".

SELF Attention

NLP → WORDS → Numbers (Vectorization)

Methods : One Hot Encoding
(OHE) :
Mat | Cat | Rat

	mat	cat	rat	
mat	1	0	0	<Vector><mat>
cat	0	1	0	<Vector cat>
rat	0	0	1	<Vector Rat>

TF-IDF

④ Word2Vec

word embeddings { king → [0.6, 0.1, 1.0, 0.9]
 Queen → [0.3, 0.2, 0.9, 1.0]

④ Problem with Word Embedding

Word Embedding captures "Average Meaning" of a word. Depending on the Data set. Ex:

- 1) An apple a day keeps the doctor away
- 2) Apple is healthy
- 3) Apple is better than orange
- 4) Apple makes great phones
- .
- .

[faste, Technology]

} → [0.8, 0.1]

Let say we 10K sentences f 9K as fruit f 1K as Technology. Hence

Vectors will be something like

[0.9, 0.1] tilted towards food

(Average weighted), hence

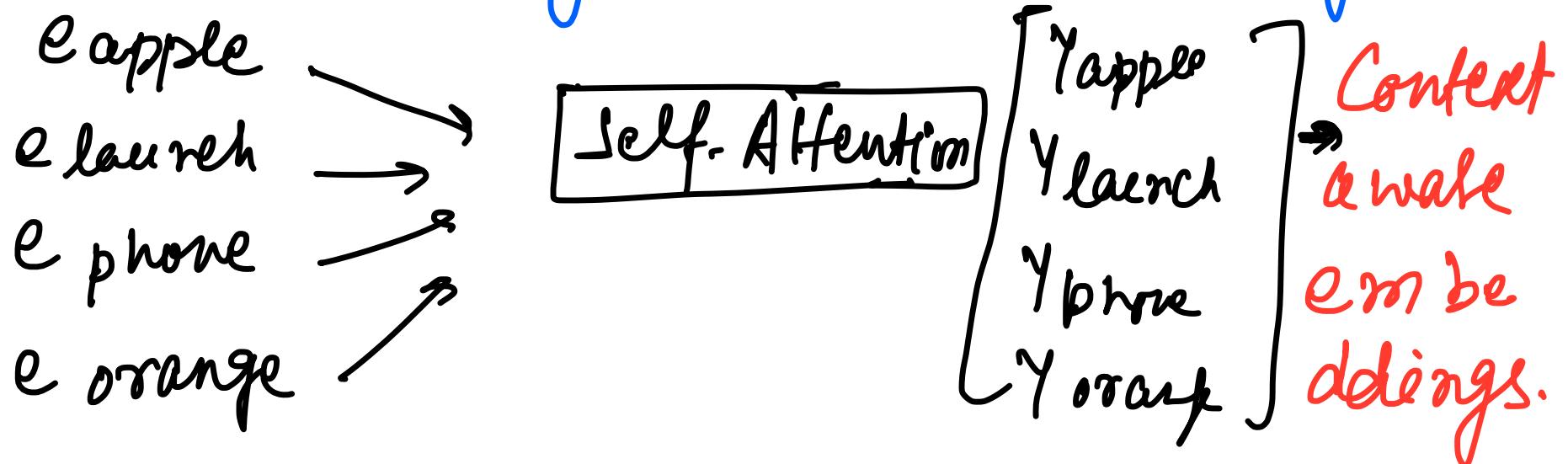
word embeddings lot depend upon

Data set. & embeddings are STATIC.

But we need contextual embeddings which can generate embedding for the context. Ex: Apple as Tech

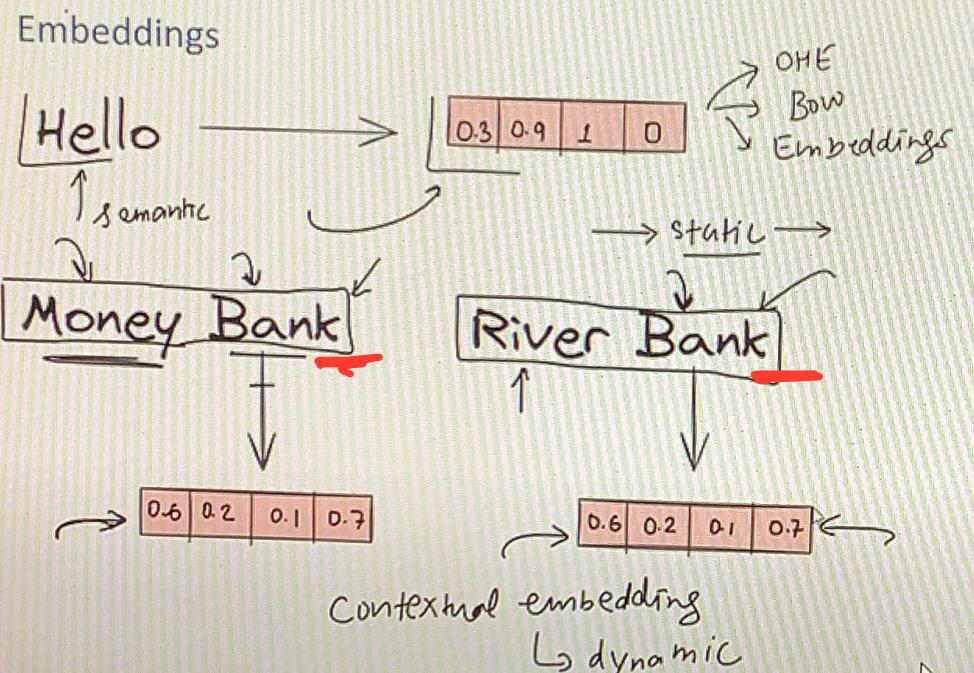
Apple as food both should have different embeddings.

* Self-Attention helps us to achieve Dynamic Embeddings.

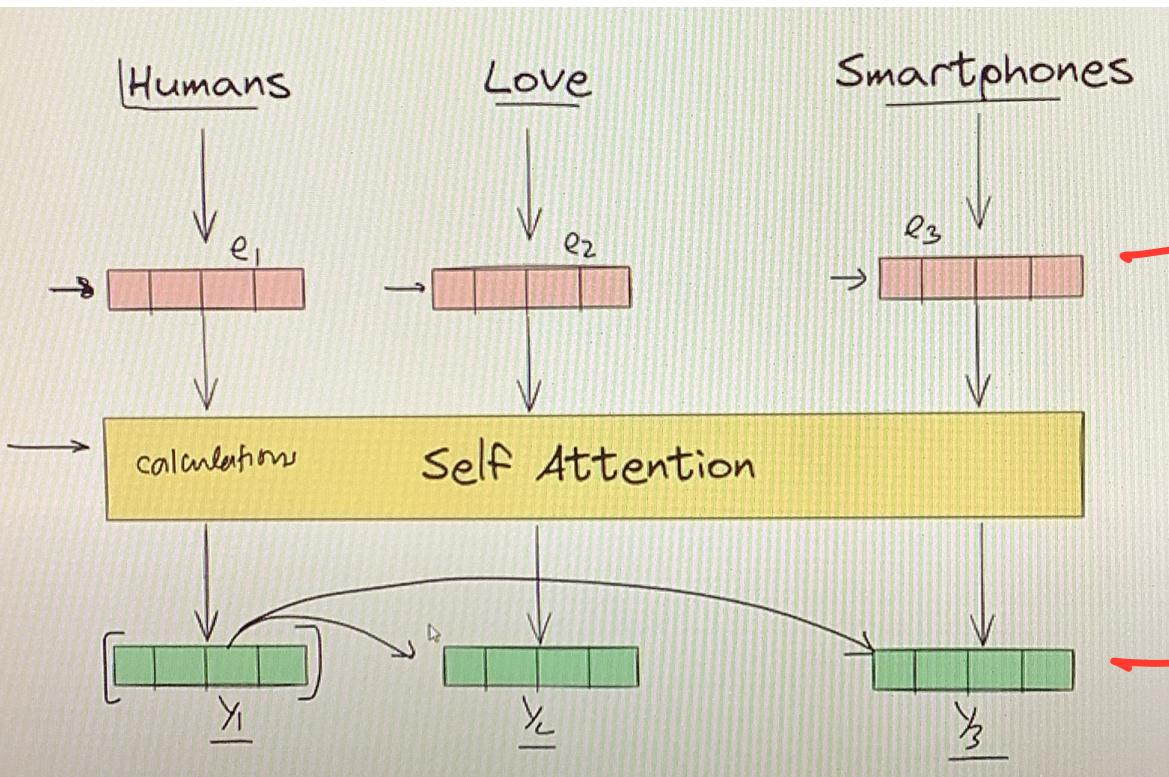


Static Embeddings

Word "Bank"



has very different
meaning in both
context but embed
- dings is same.



→ Static Embeddings
(Word2Vec --)

→ Context
Aware Embedding.

Context Embedding understand the context of word (ex: Apple as a fruit or apple as Tech). Now we

will understand "Self-Attention"

money bank grow → (SBS, HDFC etc..)

river bank flows → (River Bank)

Will think "bank" is a combination

Money , Bank, Grow)

money bank grows

$$\text{bank} = 0.2 * \text{Money} + 0.7 * \text{Bank} + 0.1 * \text{grows}$$

river bank flows (river, bank, flows)

$$\text{bank} = 0.5 \text{ river} + 0.4 \text{ bank} + 0.1 \text{ flows}$$

Now bank has become context

dependent. Below is the detail for

both sentences. Each sentence is

S1

$$\begin{aligned}\text{money} &= 0.7 \text{ money} + 0.2 \text{ bank} + 0.1 \text{ grows} \\ \text{bank} &= 0.25 \text{ money} + 0.7 \text{ bank} + 0.05 \text{ grows} \\ \text{grows} &= 0.1 \text{ money} + 0.2 \text{ bank} + 0.7 \text{ grows}\end{aligned}$$

S2

$$\left. \begin{aligned}\text{river} &= 0.8 \text{ river} + 0.15 \text{ bank} + 0.05 \text{ flows} \\ \text{bank} &= 0.2 \text{ river} + 0.78 \text{ bank} + 0.02 \text{ flows} \\ \text{flows} &= 0.4 \text{ river} + 0.01 \text{ bank} + 0.59 \text{ flows}\end{aligned} \right\}$$

presented as combination of other words

in that sentences.

Now will replace words with their embeddings. Ex:-

$$\left[\begin{array}{l} e_{\text{money}}^{(\text{new})} = 0.7 e_{\text{money}} + 0.2 e_{\text{bank}} + 0.1 e_{\text{grows}} \\ e_{\text{bank}}^{(\text{new})} = 0.25 e_{\text{money}} + 0.7 e_{\text{bank}} + 0.05 e_{\text{grows}} \\ e_{\text{grows}}^{(\text{new})} = 0.1 e_{\text{money}} + 0.2 e_{\text{bank}} + 0.7 e_{\text{grows}} \end{array} \right]$$

e_{money}
new
is weighted

average of (e_{money} , e_{bank} , e_{grows})

these weights (0.7, 0.2, 0.1) represents

similarity e_{money} new & (e_{money} ,

e_{bank} , e_{grows}).

Similarity between embedding (=Vector) can be obtained using

Dot Product). Hence weights (0.7, 0.2, 0.1)

can be consider as Dot Product between those embeddings.

ndim

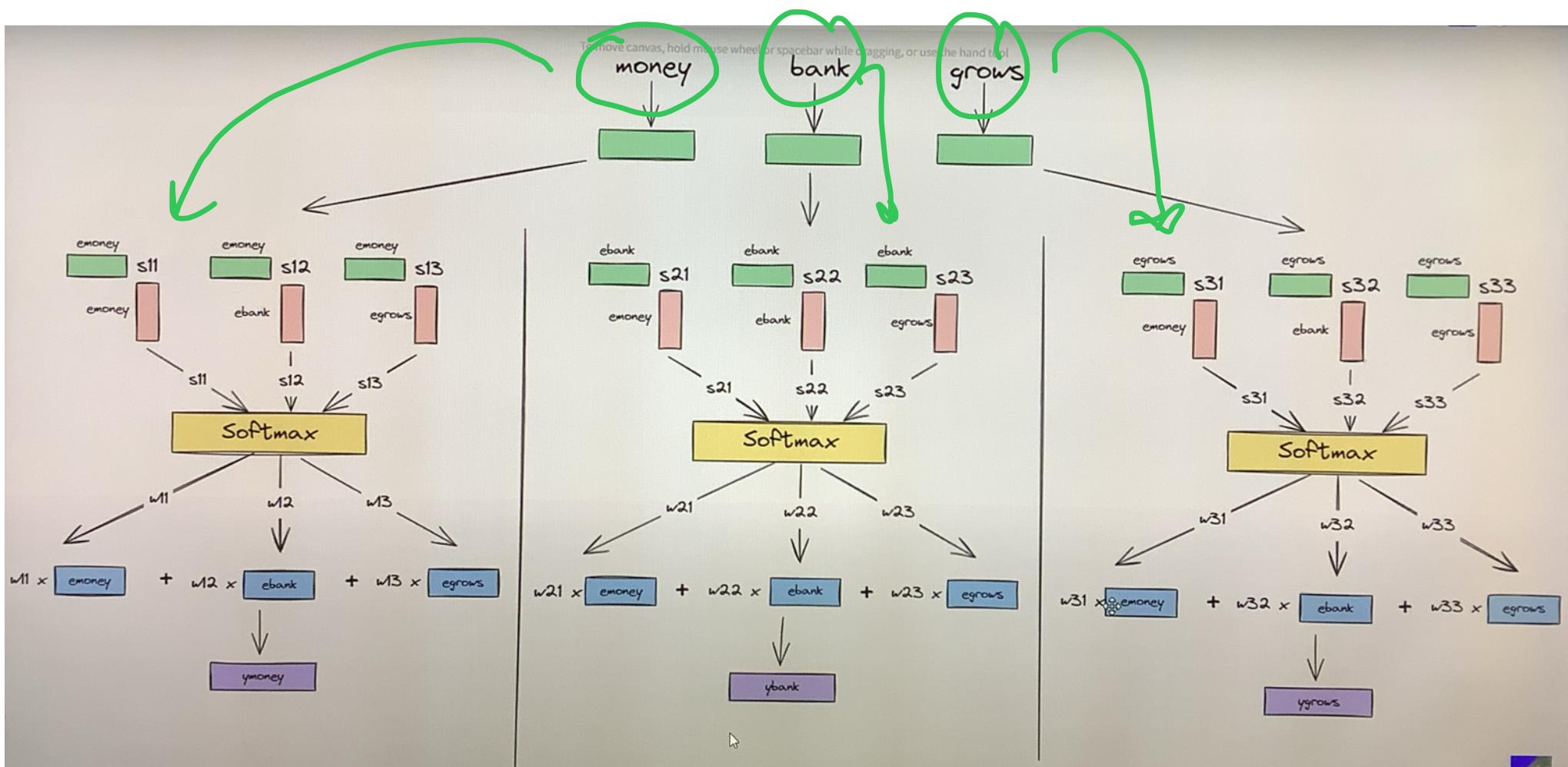
$$\begin{aligned} e_{\text{money}}^{(\text{new})} &= 0.7e_{\text{money}} + 0.2e_{\text{bank}} + 0.1e_{\text{grows}} \\ e_{\text{bank}}^{(\text{new})} &= 0.25e_{\text{money}} + 0.7e_{\text{bank}} + 0.05e_{\text{grows}} \\ e_{\text{grows}}^{(\text{new})} &= 0.1e_{\text{money}} + 0.2e_{\text{bank}} + 0.7e_{\text{grows}} \end{aligned}$$

similarity

$$\begin{array}{c} e_{\text{money}} \quad e_{\text{money}} \\ e_{\text{money}} \quad e_{\text{bank}} \\ e_{\text{money}} \quad e_{\text{grows}} \\ e_{\text{grows}} \quad e_{\text{money}} \end{array}$$

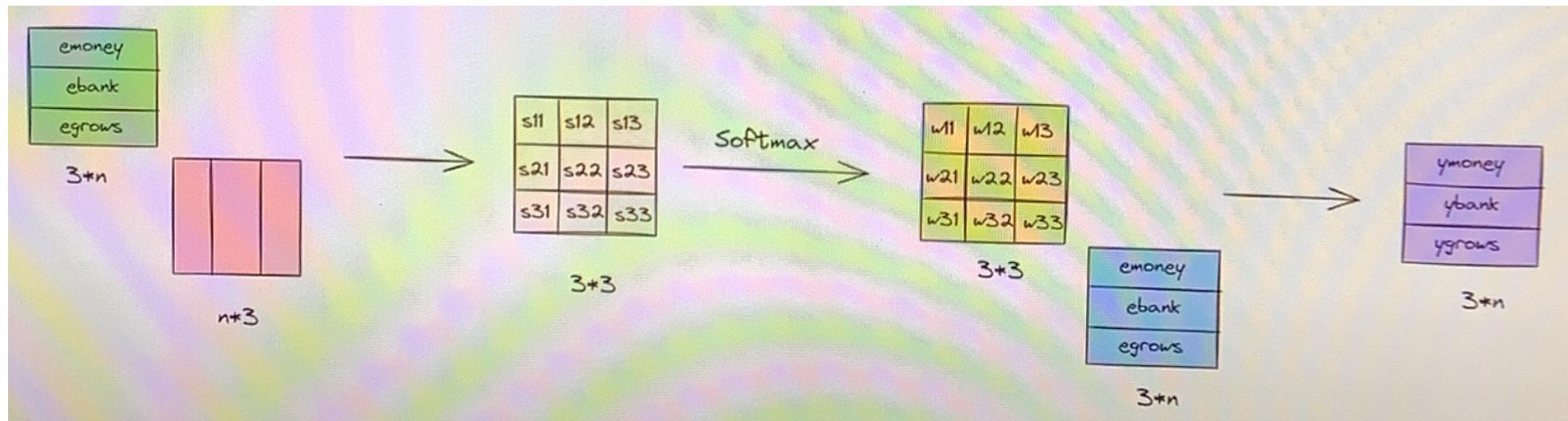
dot products

$$e_{\text{bank}} \cdot e_{\text{money}}^T \quad - \quad w_t$$
$$e_{\text{bank}} \cdot e_{\text{bank}}^T \quad - \quad w_t$$
$$e_{\text{bank}} \cdot e_{\text{grows}}^T \quad - \quad w_t$$
$$e_{\text{grows}} \cdot e_{\text{money}}^T \quad - \quad w_t$$
$$e_{\text{bank}}^{(\text{new})} = [e_{\text{bank}} \cdot e_{\text{money}}^T] e_{\text{money}} + [e_{\text{bank}} \cdot e_{\text{bank}}^T] e_{\text{bank}} + [e_{\text{bank}} \cdot e_{\text{grows}}^T] e_{\text{grows}}$$



* The above operation is parallel
 embeddings for each word can be
 calculated parallelly they don't
 dependency of previous step's word.

Below is the Parallel Approach.



⊗ There is no learning parameter.

what it means our embeddings is

not dependent on the data. Ex: =
eng → hindi

$\langle \text{ENGLISH} \rangle$ How are you	$\langle \text{HINDI} \rangle$ कैसे होते हैं
---	---

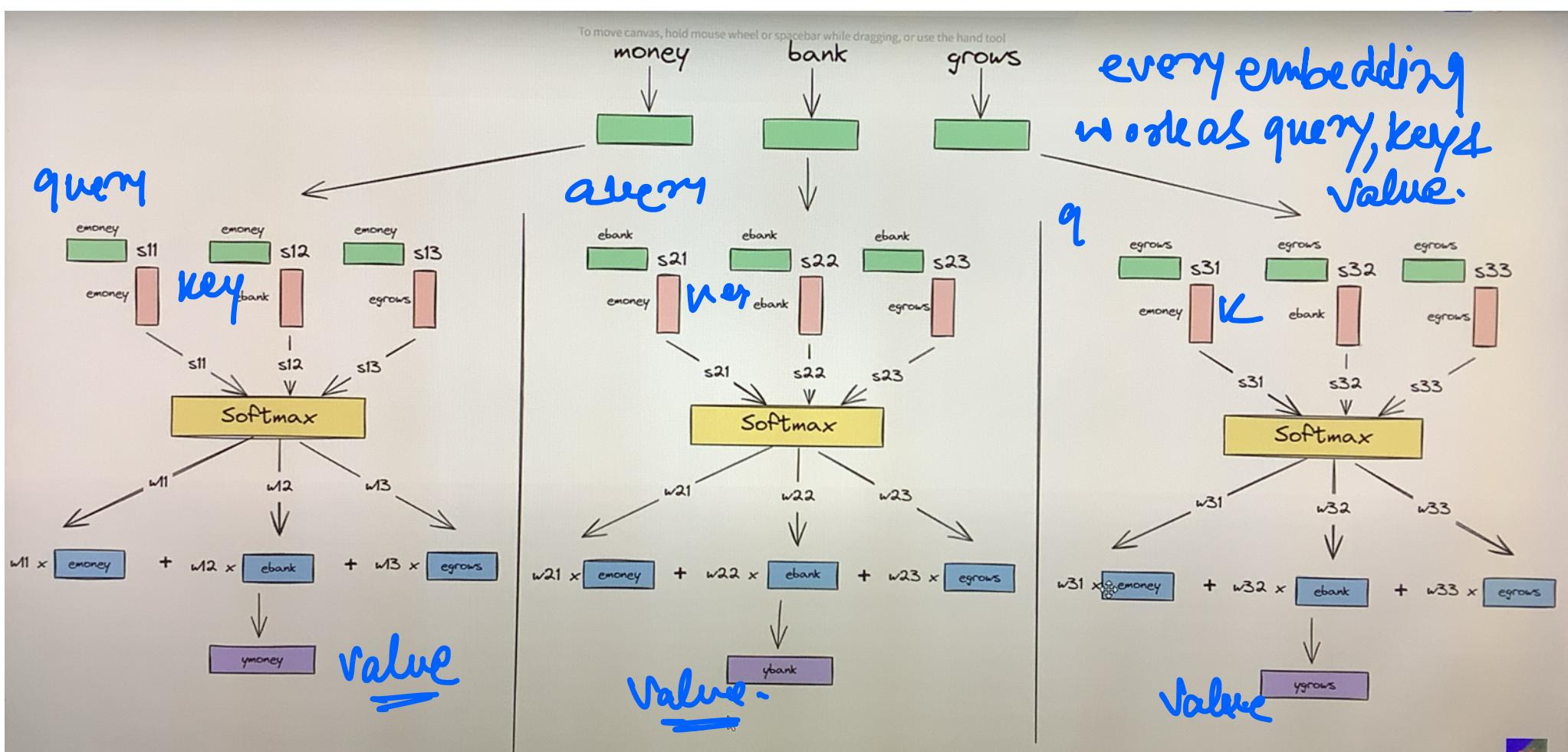
Translation is not TASK specific.
problem

BR: piece of lake → આણાની કામ] dataset.

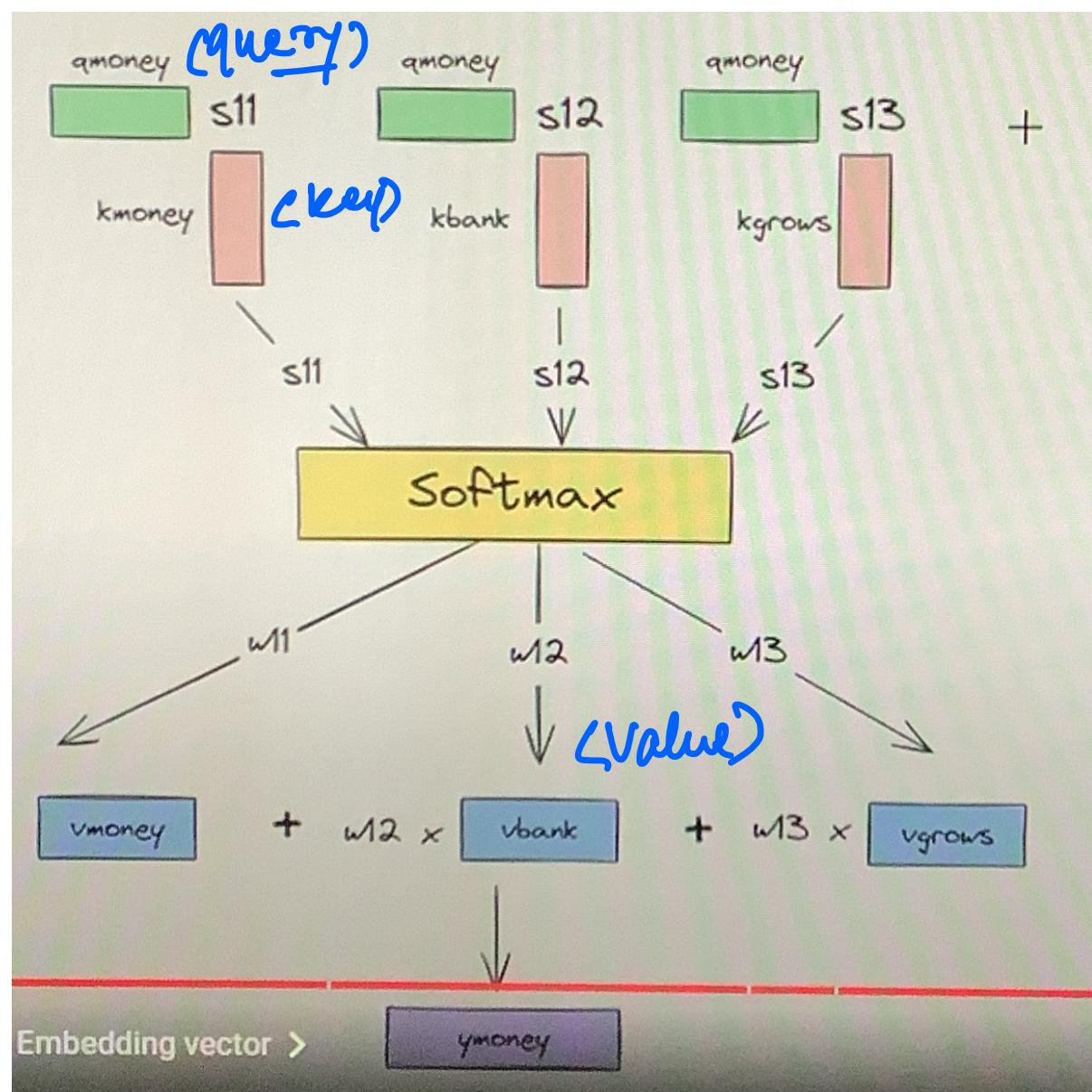
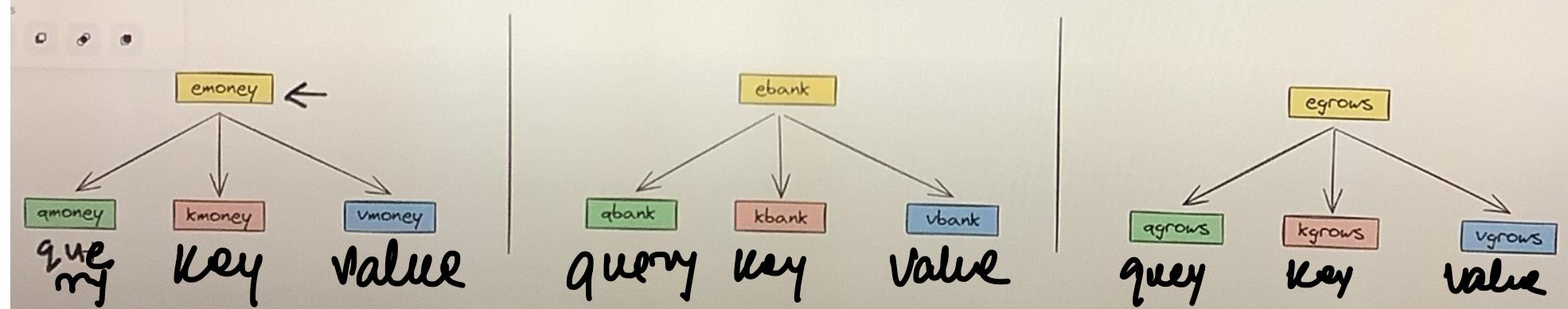
but it will generate કેવું એવું ગણિત

because there is no learning from
DATASET. Hence we need

TASK specific embeddings. Hence we
need Weights & Biases into process.



Now will take embedding create 3 embeddings out of that.



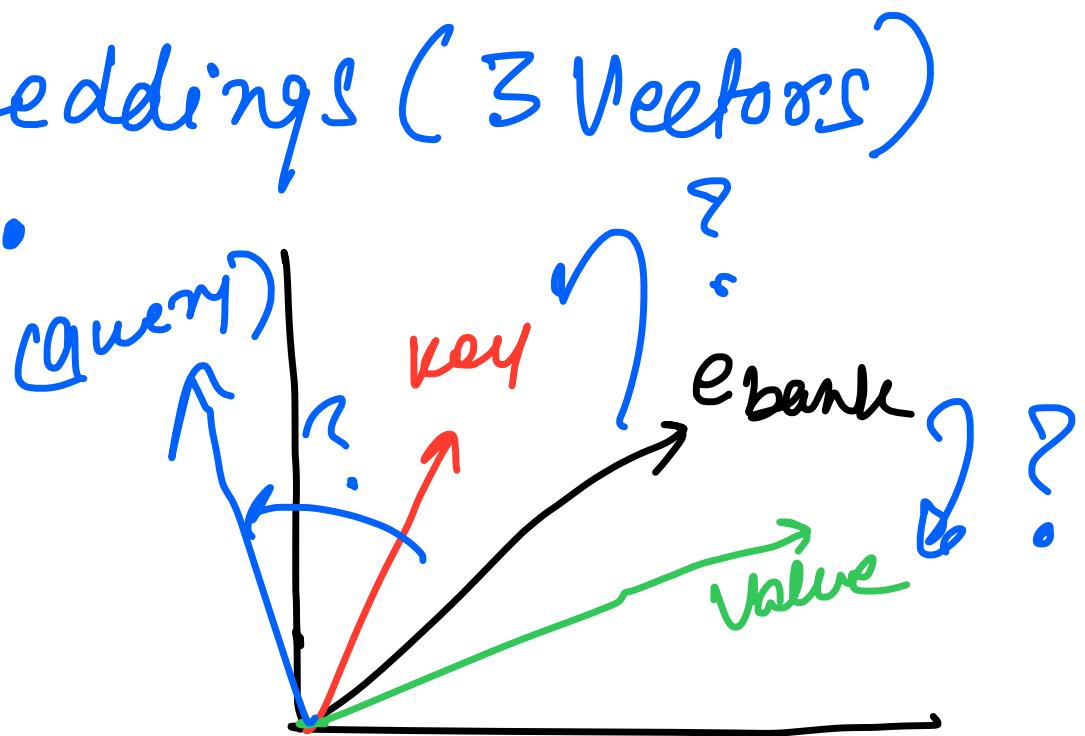
Will replace embeddings generated in this Archite - chure

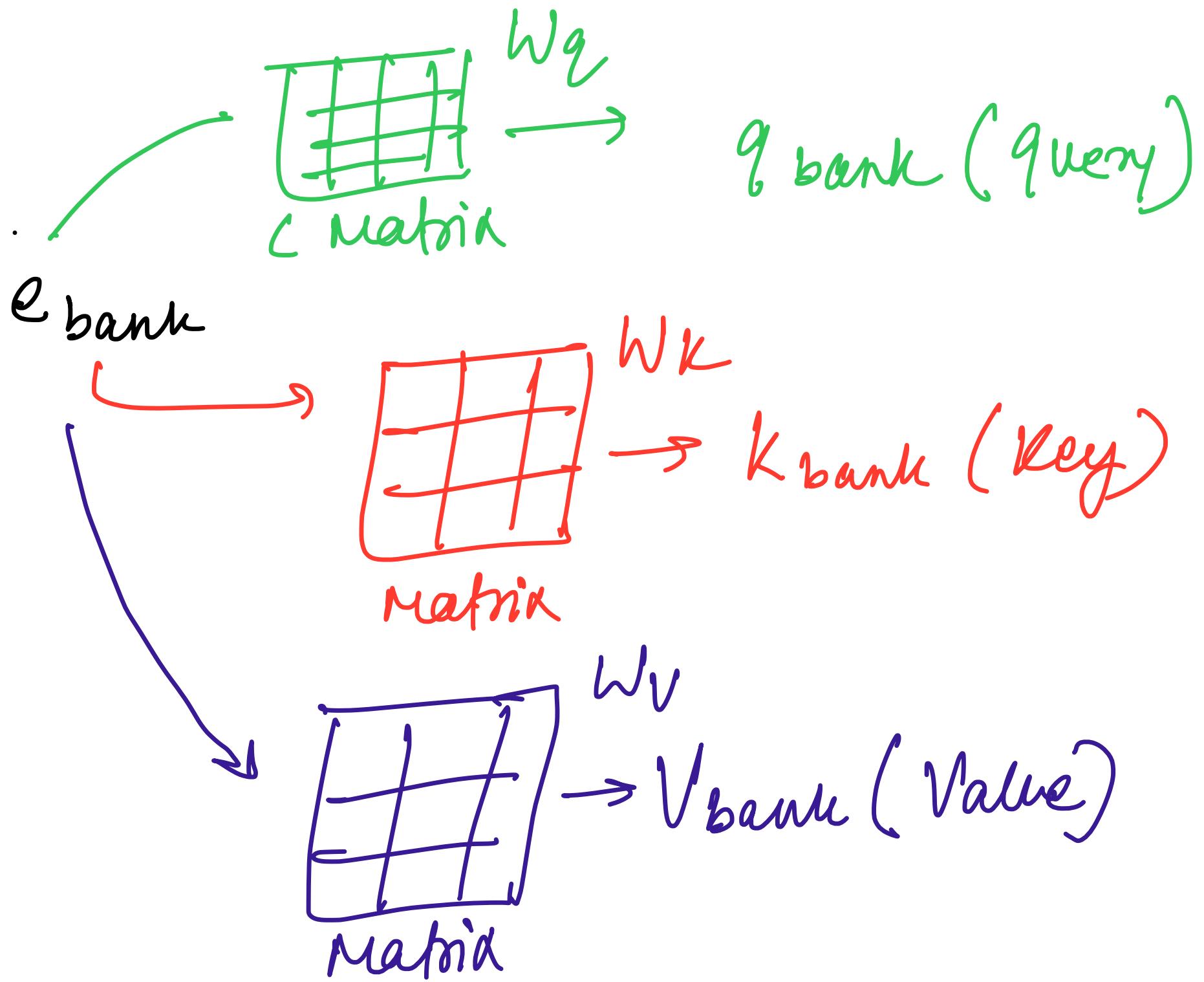
Now we have to figure out how to make 3 embeddings (3 Vectors) out of ONS.

In Algebra we

have limited way to do this

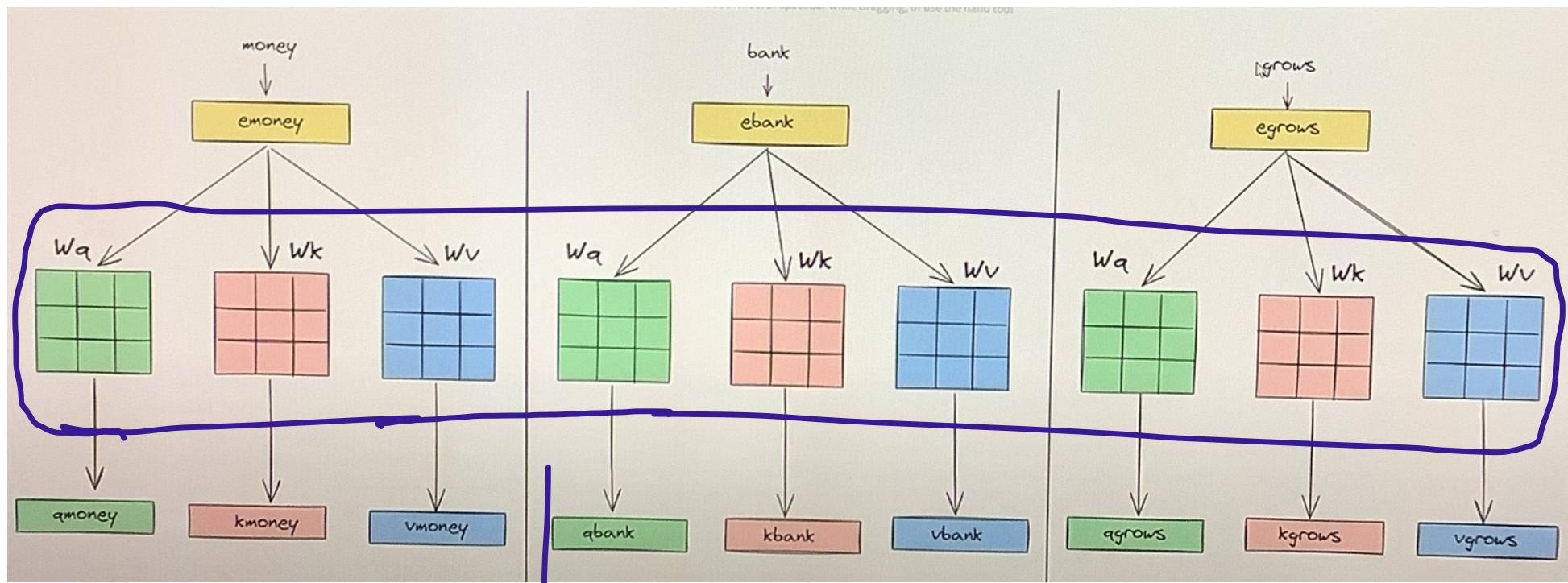
in Linear Algebra Magnitude Transform or Linear Transform.



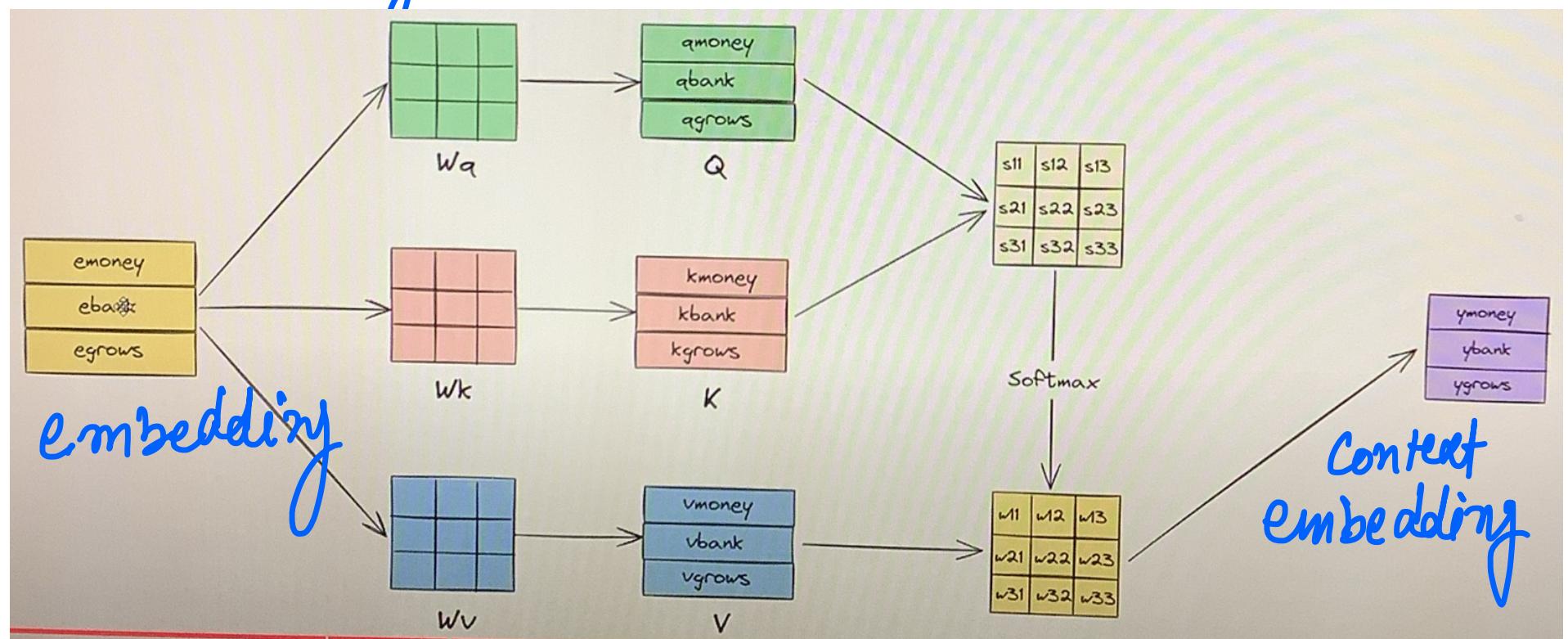


Now how to decide W^q , W^k , W^v .

Here comes to learning from Data.
will start with "Random Weights".



Creating $(Q_{ui}, \text{Key}, \text{Value})$ from
embeddings. \hookrightarrow (Same matrix for all words)



< Parallelization >

Now Attention can be put mathematically like below:-

$$\text{Attention}(Q, K, V) = \text{Softmax}\left(\frac{QK^T}{\sqrt{d_k}}\right)V$$

query key Value.

Same formula from paper

$$\text{Attention}(Q, K, V) = \text{Softmax}\left(\frac{QK^T}{\sqrt{d_k}}\right)V$$

We are scaling with $\sqrt{d_k}$. Hence in paper we called it Scaled Dot product.

Now why Scaling? Why $\sqrt{d_k}$

d_k = dimension of context embeddings.

Why Scaling: Nature of Dot Product

low Dim \rightarrow Dot Product \rightarrow low Variance

High Dim \rightarrow Dot Product \rightarrow high Variance

1 dim \rightarrow $\text{Var}(X)$

2 dim \rightarrow $2 \cdot \text{Var}(X)$

3 dim \rightarrow $3 \cdot \text{Var}(X)$

· · · · · · ·

d dim $\rightarrow d \cdot \text{Var}(X)$

$$\left. \begin{array}{l} \text{if } X \rightarrow \text{Var}(X) \text{ if } Y = CX \\ Y \rightarrow \text{Var}(Y) = C^2 \text{Var}(X) \end{array} \right\} \text{⊗}$$

Now we want to retain $\text{Var}(X)$ independent of vector dimension. So if we

do $\frac{Y}{\sqrt{2}} \rightarrow \frac{1}{2} \text{Var}(Y) \rightarrow \frac{1}{2} Z \text{Var}(X)$

\approx it will retain $\approx \text{Var}(X)$.

hence divide by $\sqrt{d_k}$

$$\text{Softmax} \left(\frac{QK^T}{\sqrt{d_k}} \right) \underset{\sim}{\ast} N \equiv \text{Attention.}$$

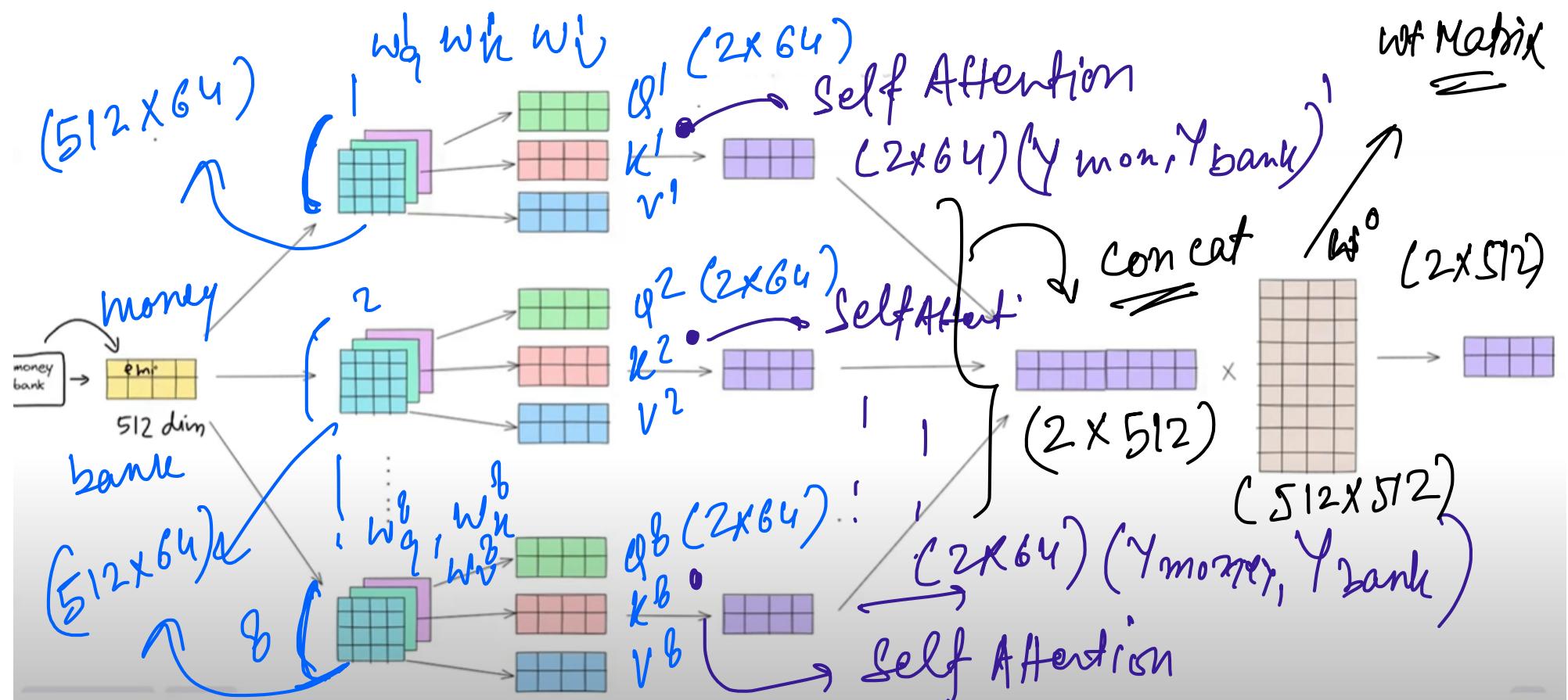
MultiHead Attention $\stackrel{?}{=}$ Single

Head Attention capture only ONE perspective off Text. Ex: Summarize document can be done in multiple ways but Self-Attention capture only one perspective.

[The man saw the astronomer with]
a telescope.

Sentence have Two meanings.
So we use more than one Self-Attention

to capture different meanings. Ex
Self-attention. 512 dim vector



We reduce $512 \rightarrow 64$ to handle
computation.

Positional Encoding \Rightarrow Self

attention has no mechanism to
understand the order of words.

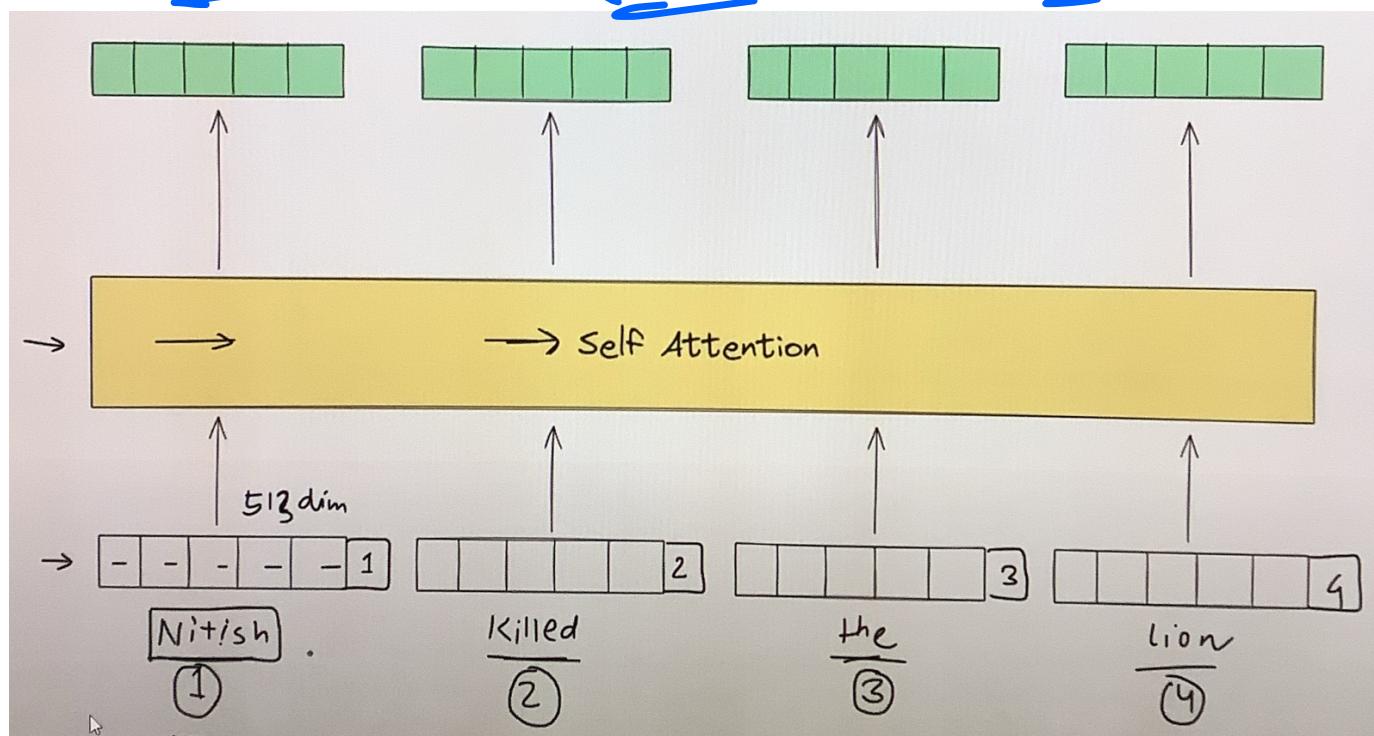
① Nithish killed lion } $\rightarrow s_1$

② lion killed Nithish } $\rightarrow s_2$

Hard to differentiate s_1 & s_2 .

* Because RNN is One Word at a time so this problem is not faced.

Simple Solution Word Count :=



Problems : 1. Unbounded

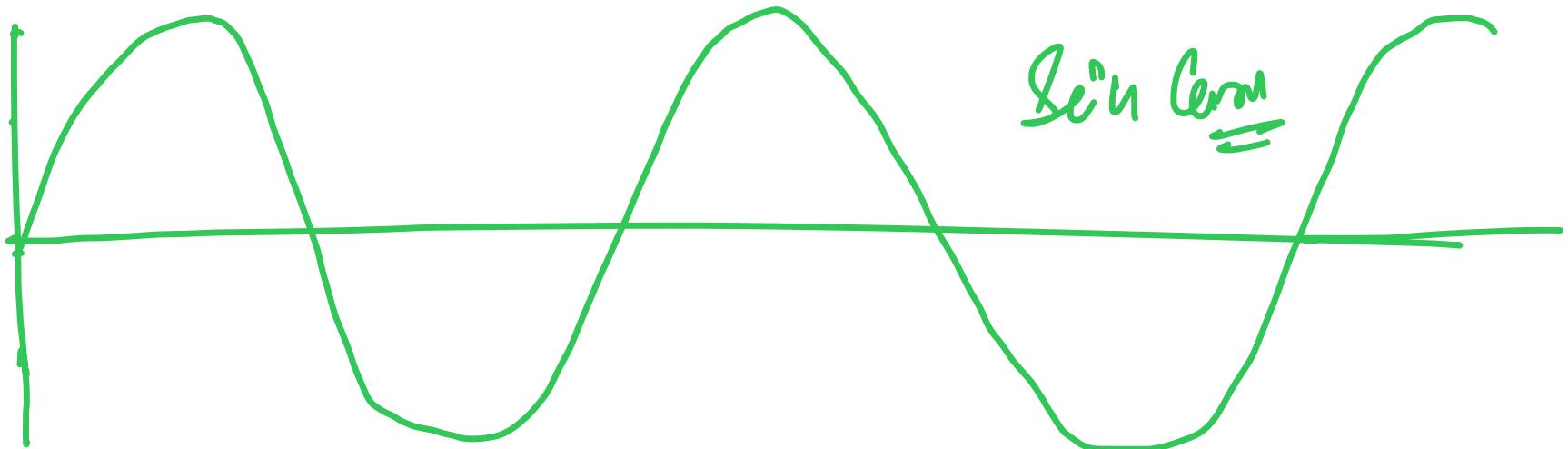
2o) Discrete positions not smooth transition

3o) No relative positioning.

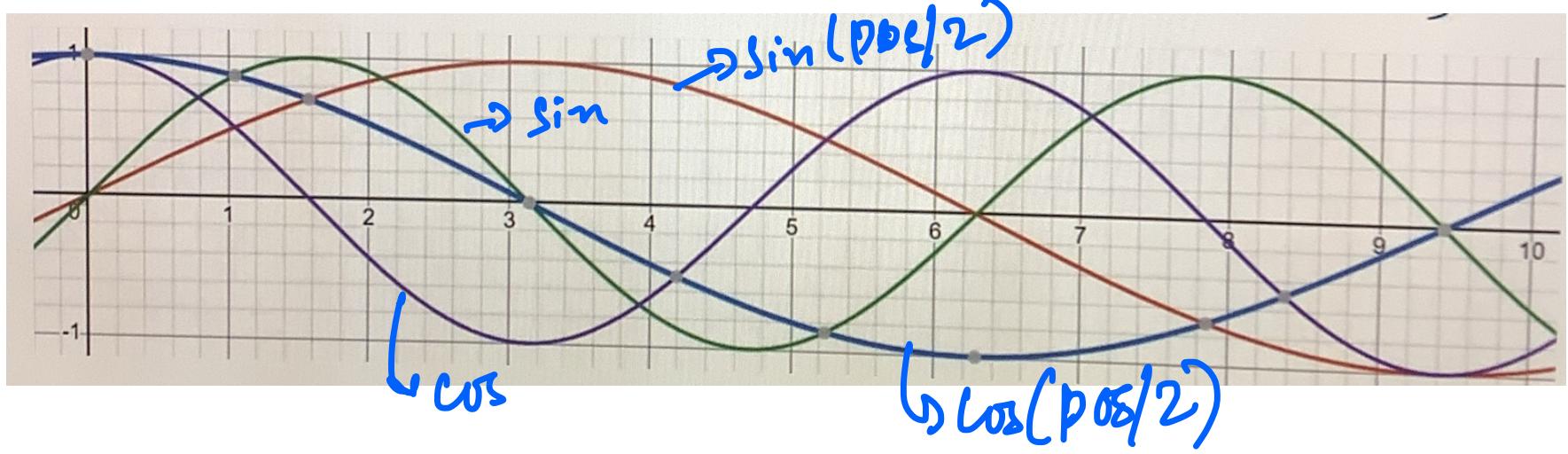
What we want a function to be

1. bounded
2. continuous
3. periodic

Options are Trigonometric functions



$y = \sin(\text{position})$ as \sin is periodic we can have same y for different position. Will use multiple functions.



So the positional embedding will be a Vector. chances of having same embedding is less. We can keep on increasing number of function .

- ④ Size of Word embedding & Positional encoding Vector will be same size.

embedding = (Word Embedding + Positional Embedding)

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

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

pos = position of the word

d_{model} = Dim of embedding

$$i = \begin{pmatrix} 0 \rightarrow d_{model}/2 - 1 \\ 1 \end{pmatrix} \underbrace{\begin{pmatrix} 0 \rightarrow 2 \\ 0, 1, 2 \end{pmatrix}}_{=}$$

river bank

$$= \langle pos=0, i=0 \rangle$$

$$P(0, 0) = \sin\left(0 / 10000^0/d_m\right) = 0$$

$$P(0, 1) = \cos\left(0 / 10000^0\right) = 1$$

$\langle \text{Pos}=0, i=1 \rangle$

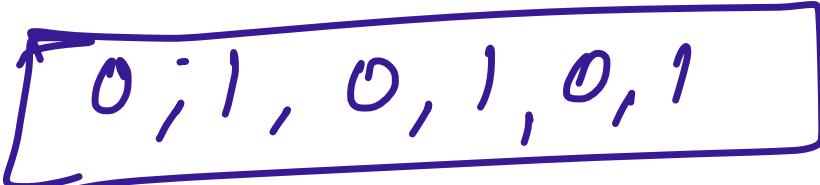
$$f(0, 2) = \sin(0 / 10000^{1/3}) = 0$$

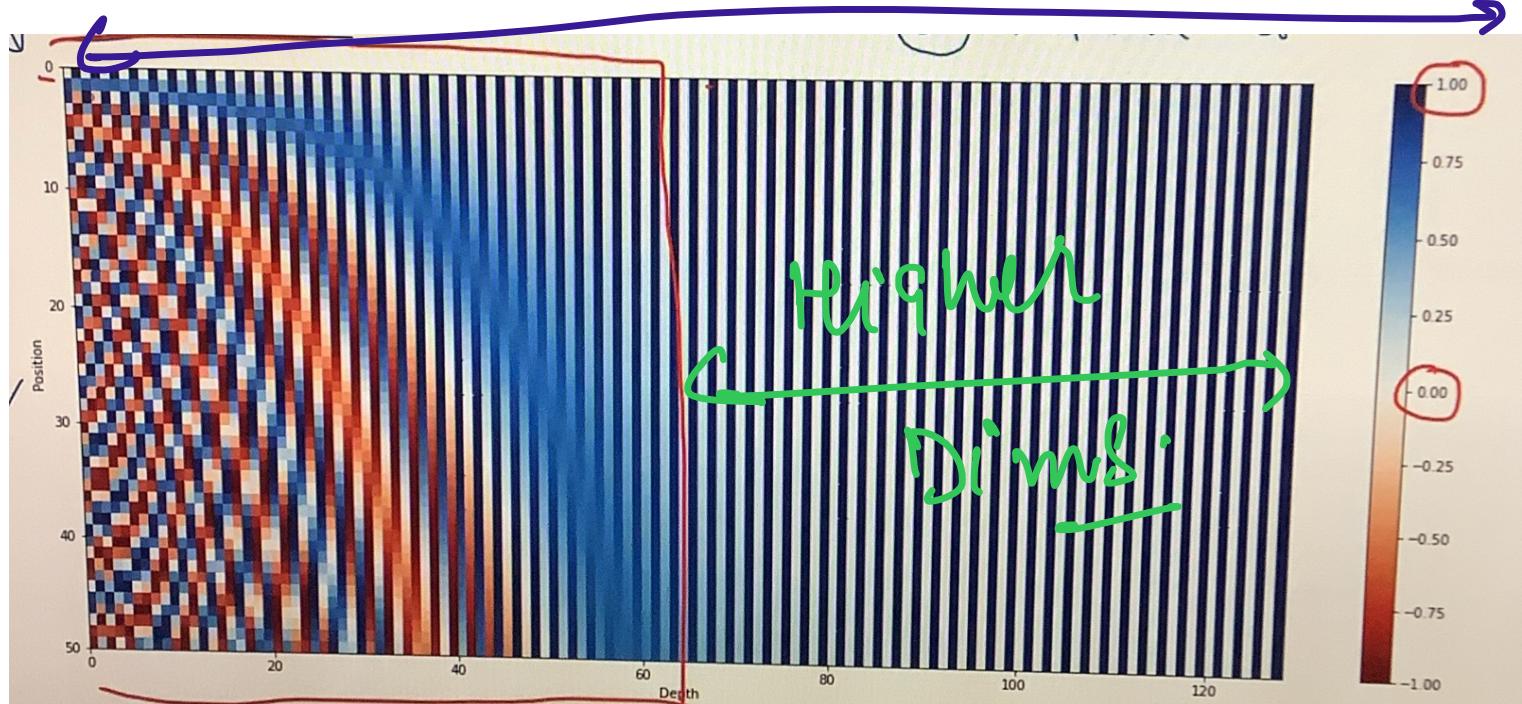
$$P(0, 3) = \cos(0 / 10000^{1/3}) = 1$$

$\langle \text{Pos}=0, i=2 \rangle$

$$f(0, 4) = \sin(0 / 10000^{2/3}) = 0$$

$$P(0, 5) = \cos(0 / 10000^{2/3}) = 1$$

river \Rightarrow 
0; 1, 0, 1, 0, 1



128 dim

50 words sentence.

As we can see $\text{pos}=0$, have 0, 1, 0, 1 ...

as we have found for word rivers.

④ So as we can see higher Dimension is almost same.

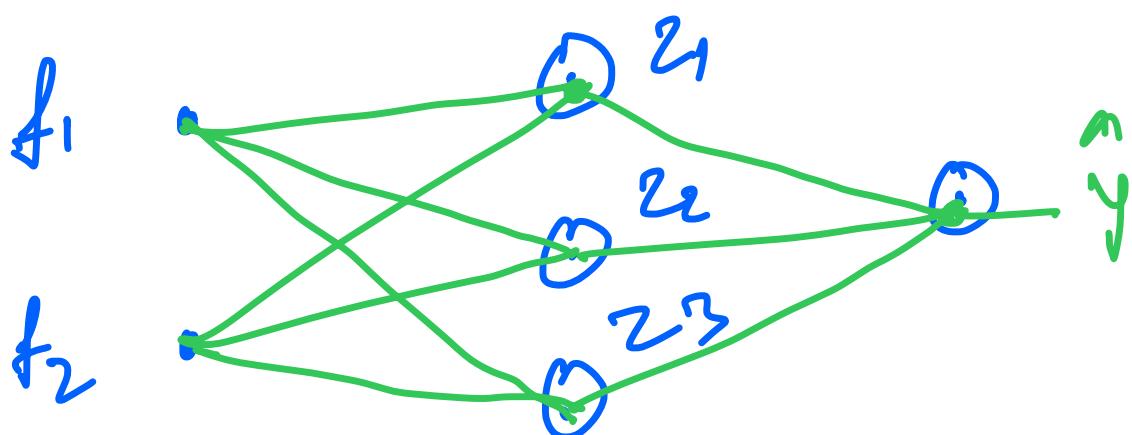
Analogy to binary encodings.

100	binary encodings
0:	0 0 0 0
1:	0 0 0 1
2:	0 0 1 0
3:	0 0 1 1
4:	0 1 0 0
5:	0 1 0 1
6:	0 1 1 0
7:	0 1 1 1
8:	1 0 0 0
9:	1 0 0 1
10:	1 0 1 0
11:	1 0 1 1
12:	1 1 0 0
13:	1 1 0 1
14:	1 1 1 0
15:	1 1 1 1

As we can see
higher bit level less
frequently it changes
(say 2^{-1}) frequency.

Layer Normalization :-

Revise Batch Normalization :-



BATCH NORM

S	f_1	f_2	z_1	z_2	z_3	
s_1	2	3	7	5	4	
s_2	1	1	2	3	4	
s_3	5	4	1	2	3	
s_4	6	1	7	5	6	
s_5	7	1	3	3	4	

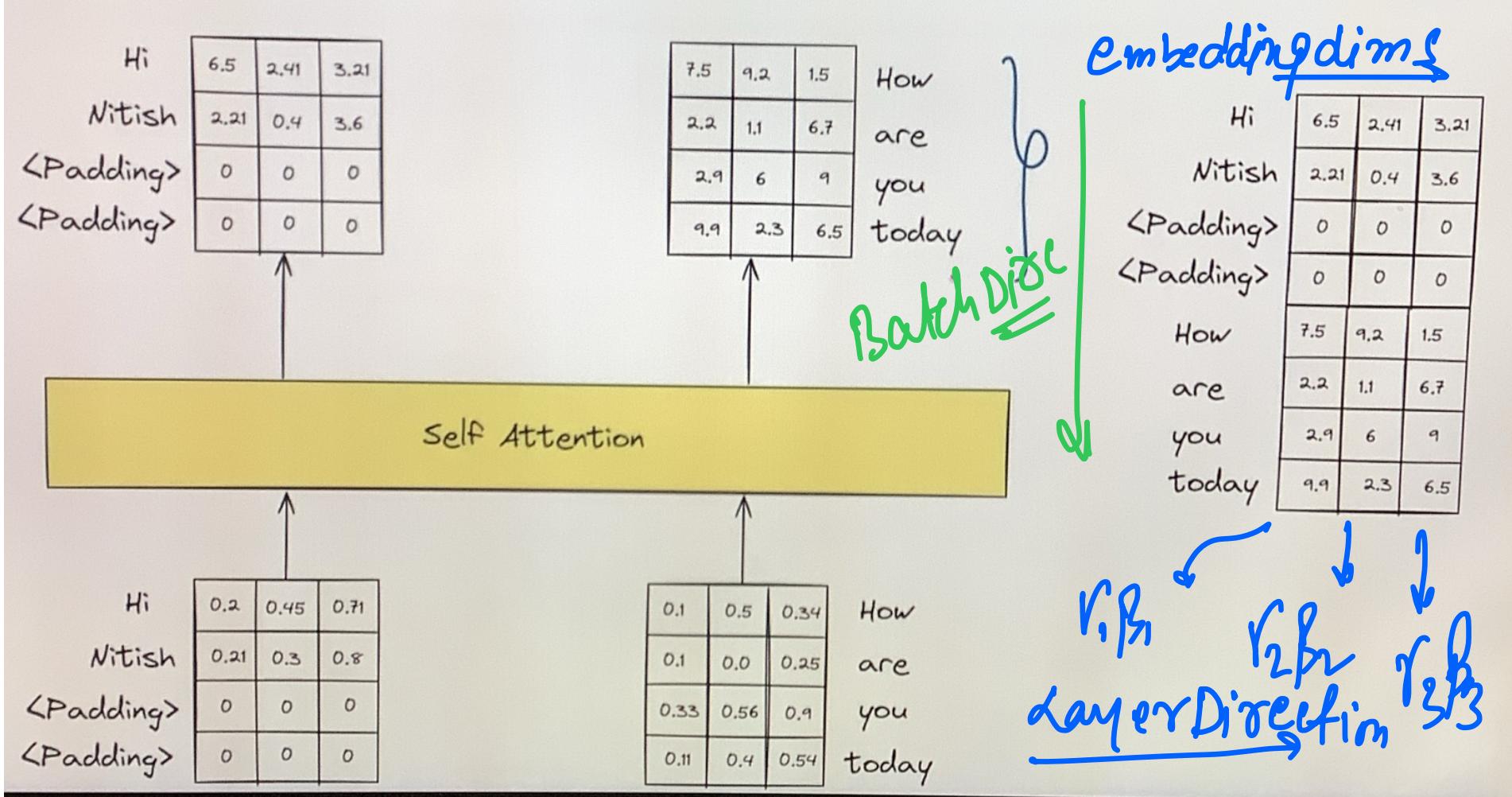
Across Batch
 i's Batch Norm
 Batch Direction
 Column Wise

Across Features → (Layer Norm)

S	f_1	f_2	z_1	z_2	z_3	
s_1	2	3	7	5	4	
s_2	1	1	2	3	4	
s_3	5	4	1	2	3	
s_4	6	1	7	5	6	
s_5	7	1	3	3	4	

Across Features
 Layer Norm
 Layer Norm Direction

Why Layer Norm In Transformers?

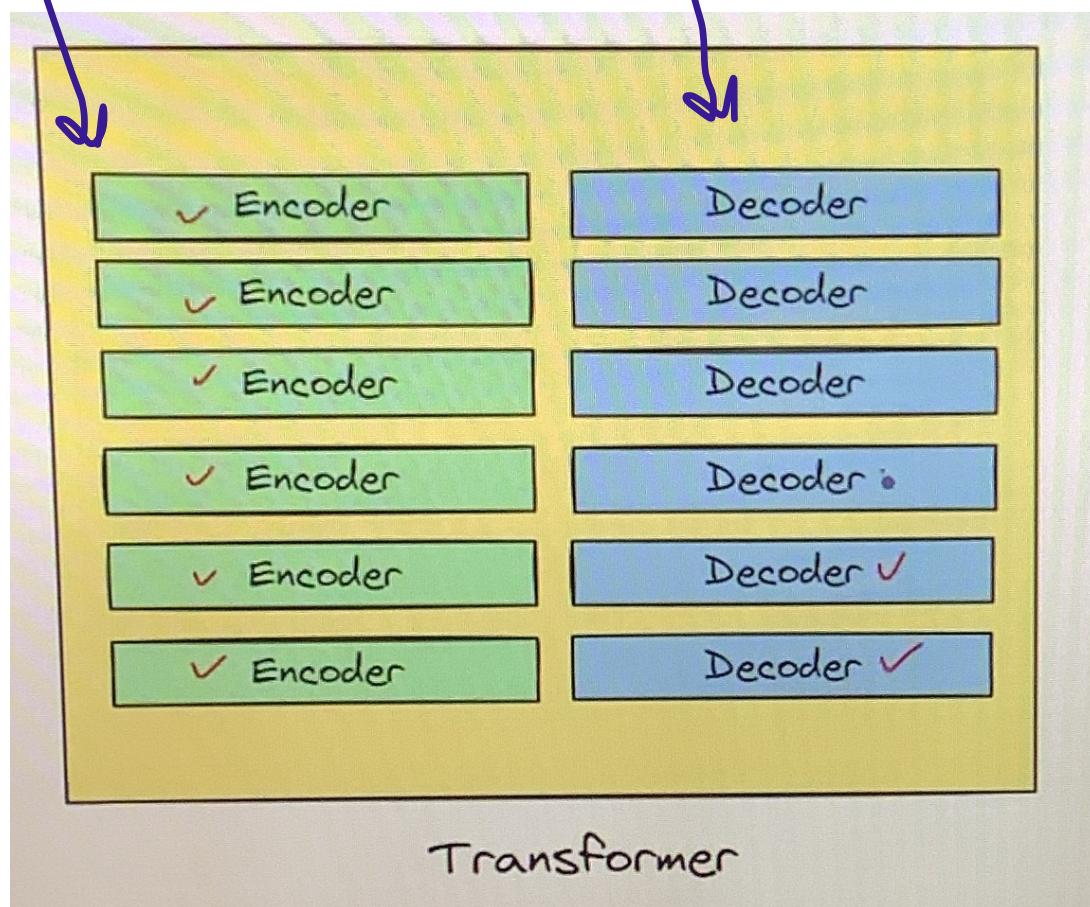
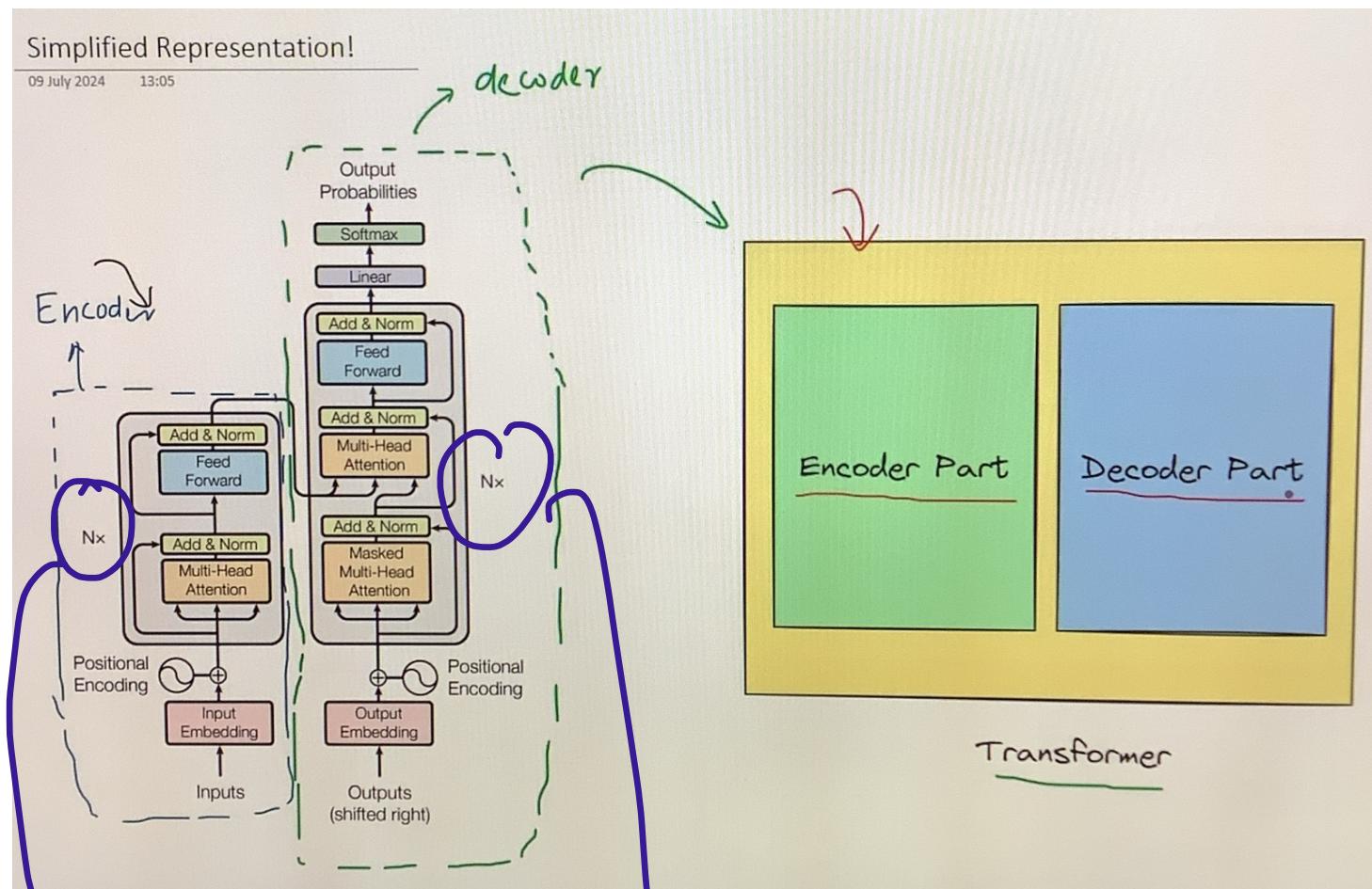


As we can see in Batch Norm (Padding)

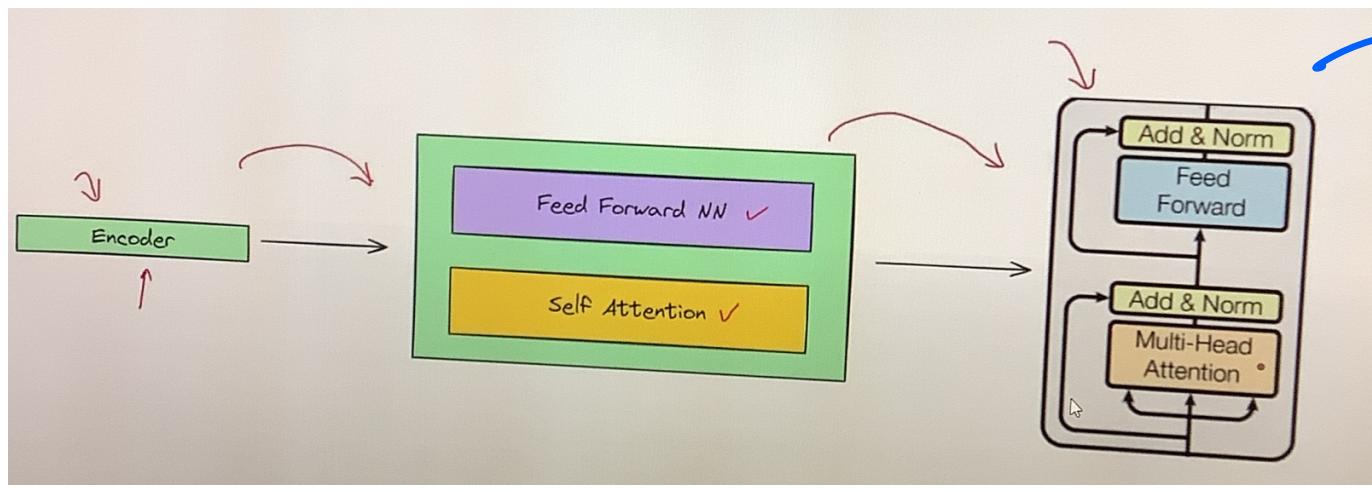
zeros will effect over all calculation.

Assume a sentence with 3 words &
>100 trailing zeros (padding)

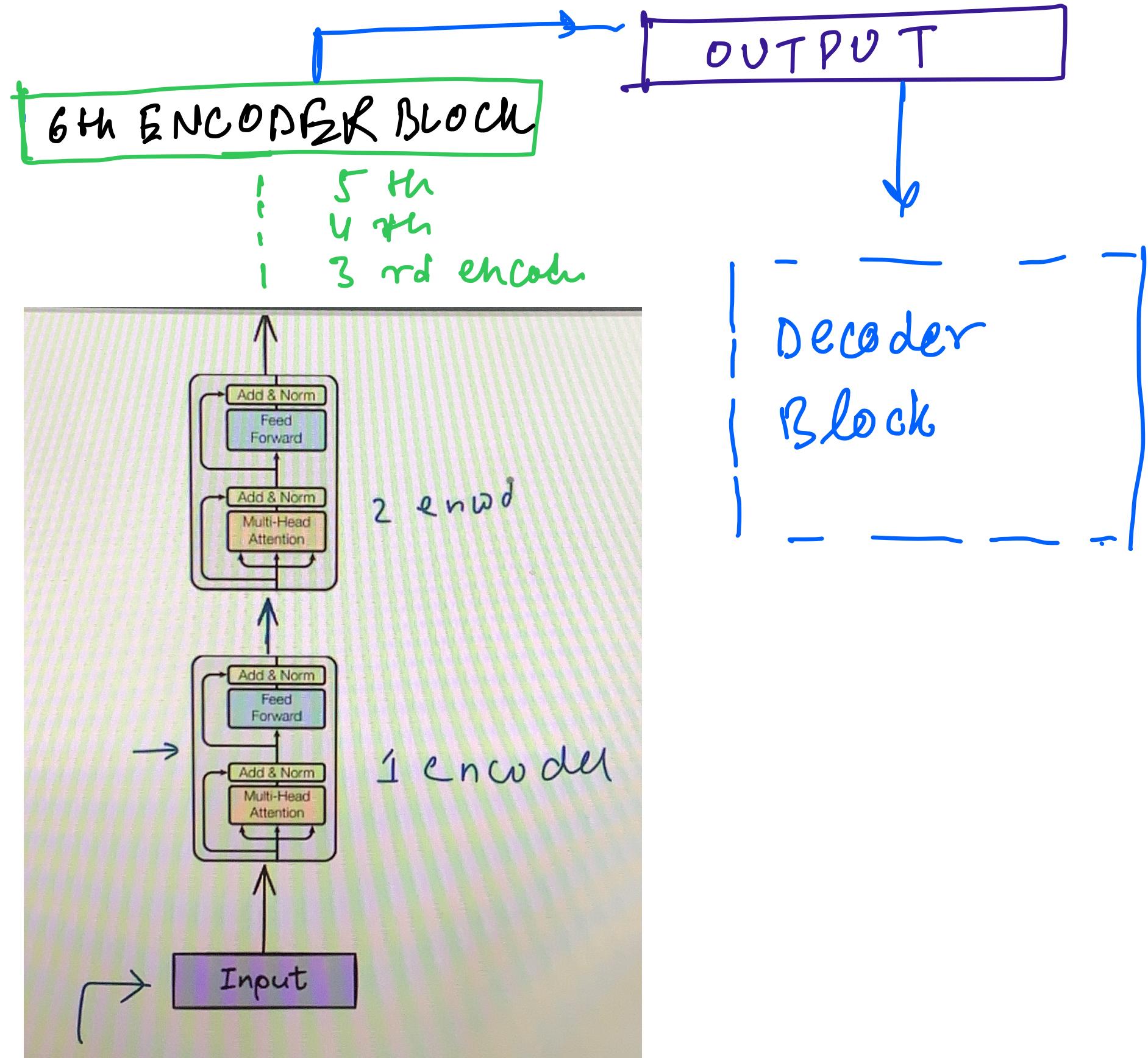
Transformer Architecture

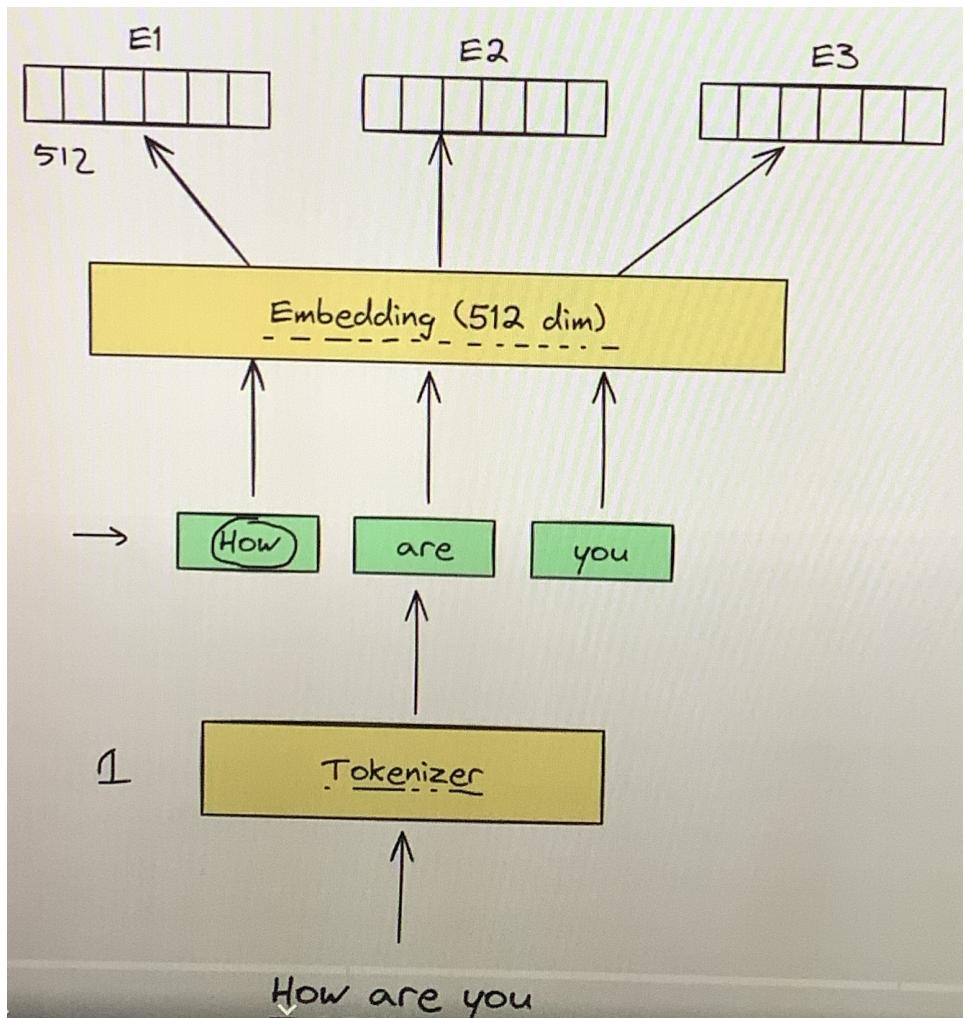


Nx means there
are multiple
Encoder & Decoder
blocks.



ON E
encoder
block.



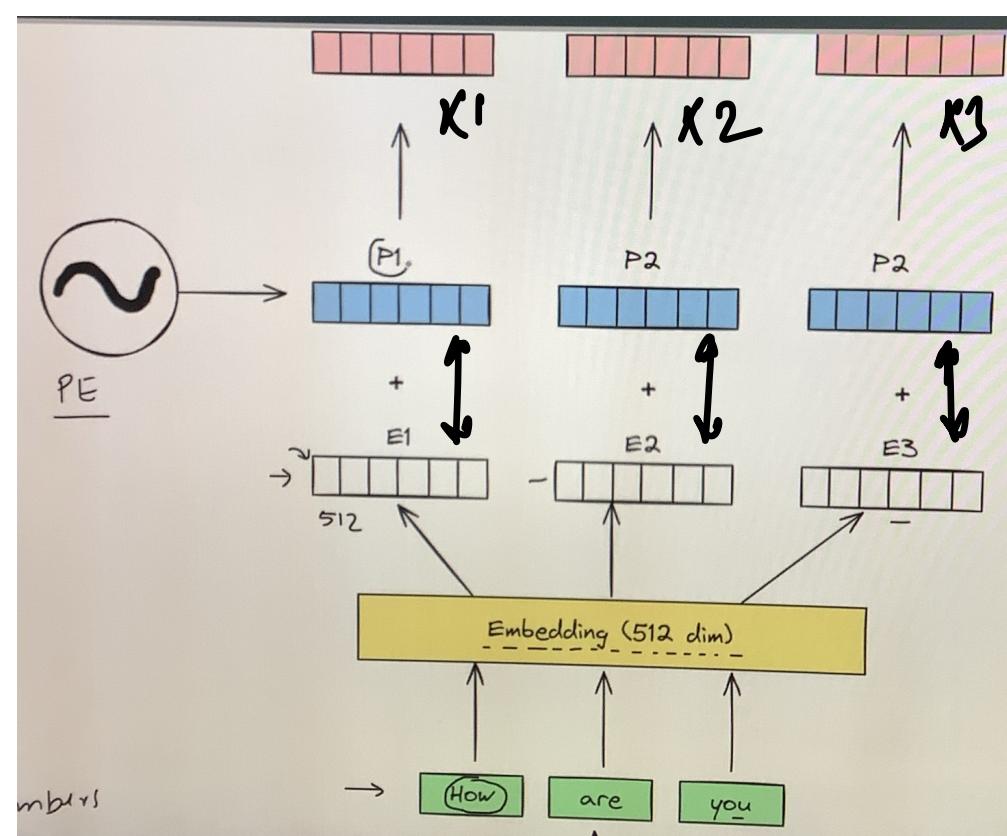


Step 1: Text = "How are You"

get Tokenize. That Tokens

goes to embedding
(Word2Vec) etc

& create E1, E2 & E3.

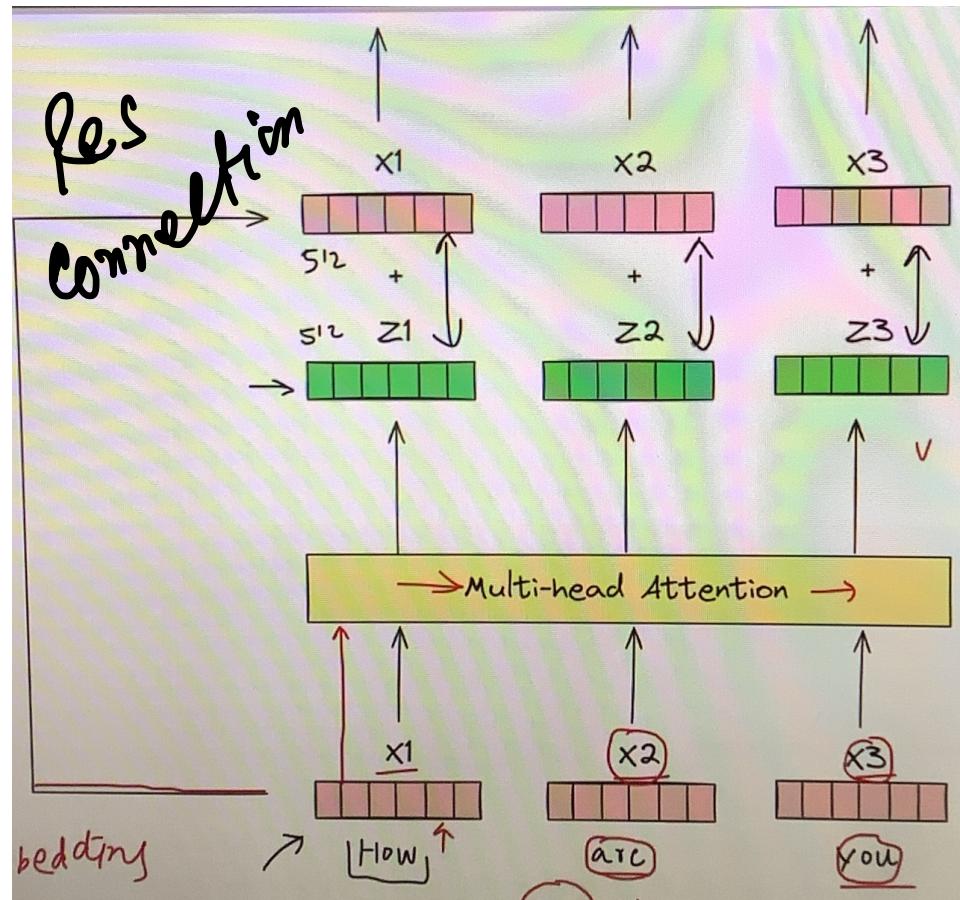


Step 2: Incorporate
positional encoding

$(P_1, P_2 \& P_3) + (E_1, E_2 \& E_3)$

generate $(X_1, X_2 \& X_3)$

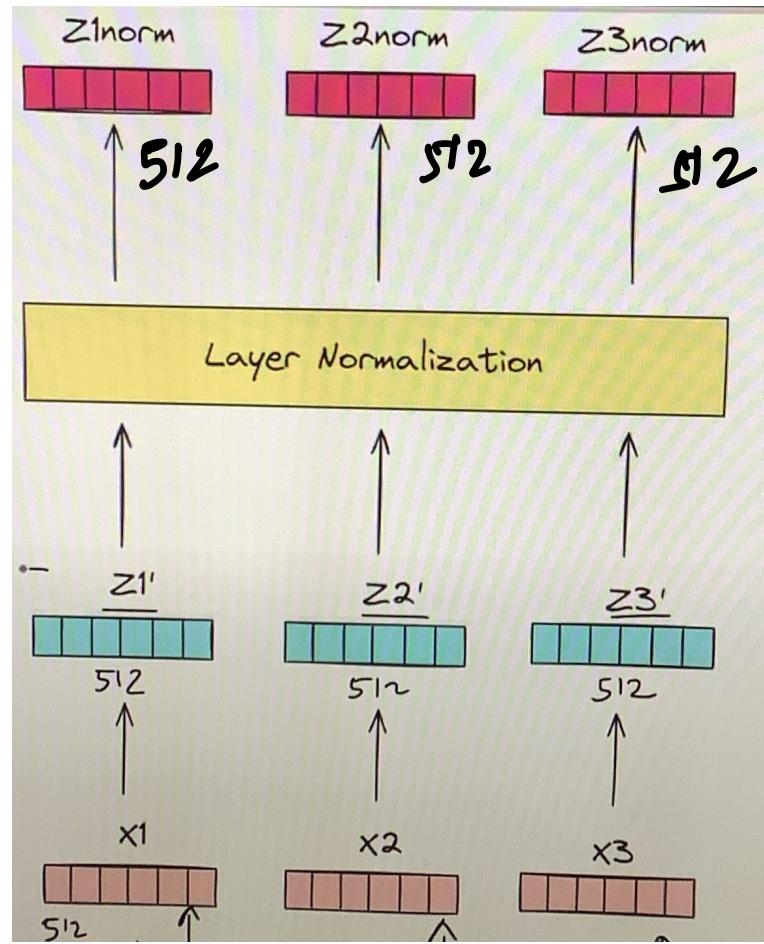
Now Input is ready =



Step 3: With
Self-Attention or
Multi-Head Attention

We create "CONTEXT"

Context-aware embedding ($\vec{z}_1, \vec{z}_2, \vec{z}_3$)

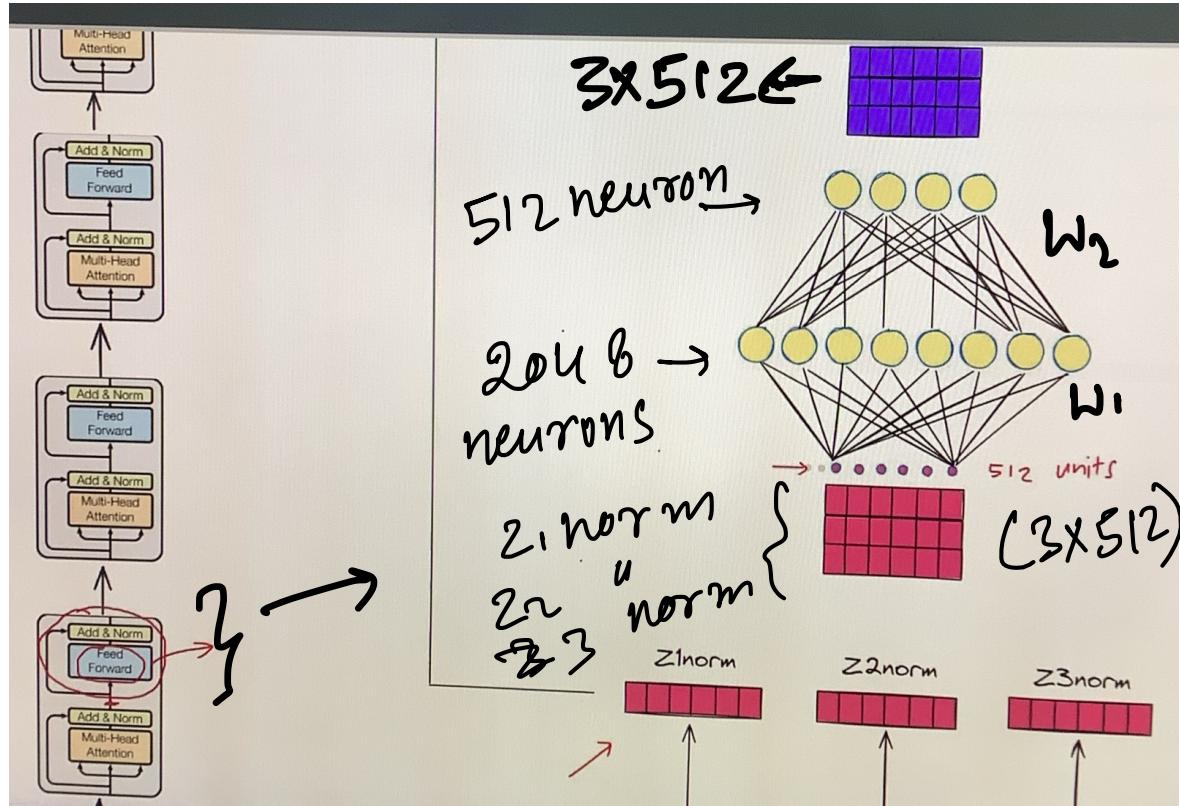


$$\begin{aligned} \text{Step 4!} \quad & z'_1 = z_1 + r_1 \\ & z'_2 = z_2 + r_2 \\ & z'_3 = z_3 + r_3 \end{aligned}$$

→ Residual Connection

Then perform layer
Norm.

$z_i \rightarrow \underline{z_i \text{ norm}}$



Step 3: $W1 = 512 \times$

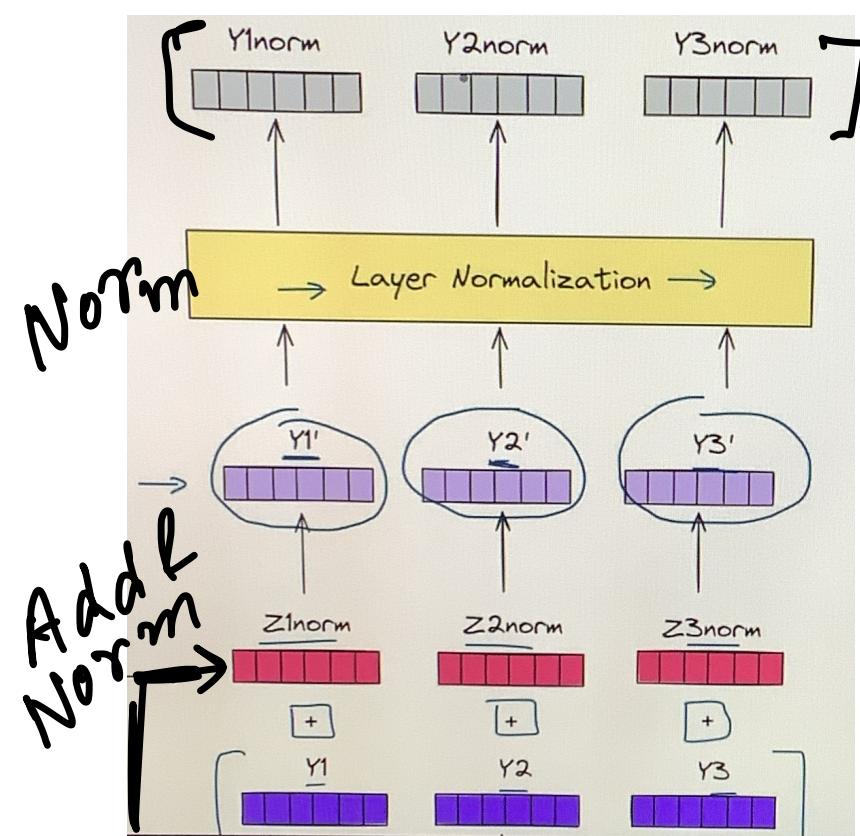
2048.

$W2 : 2048 \times 512$

$$[(3 \times 512) \times (512 \times 2048)] + b_1 = (3 \times 2048)$$

$$(3 \times 2048) \times (2048 \times 512) + b_2 = (3 \times 512)$$

⑧ Added Some Non-linearity.



Output will be
Input for next Encoder

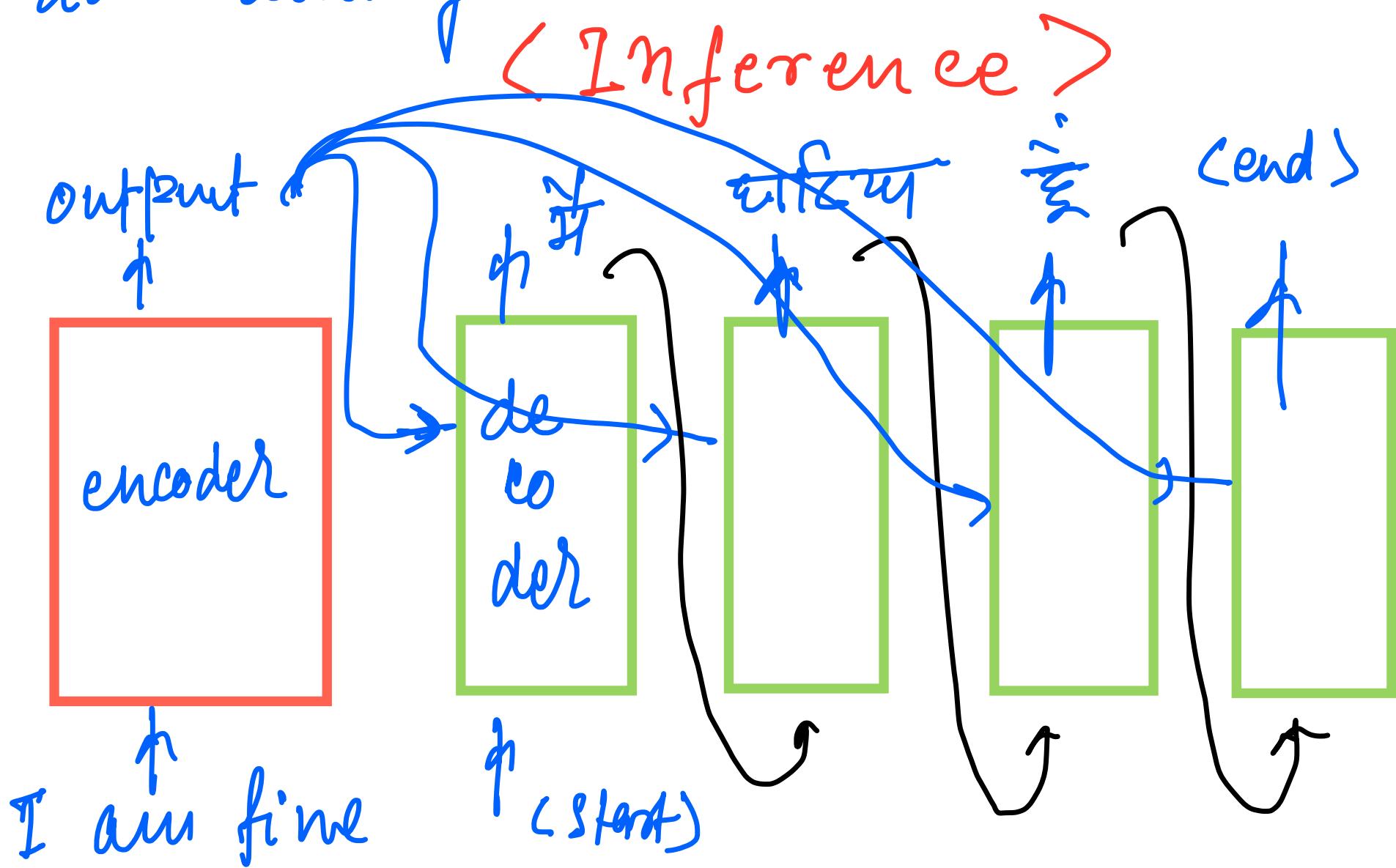
block. Each Encoder block

have their own Wts & Biases.

Masked Self Attention:

Transformer Decoder is auto regressive

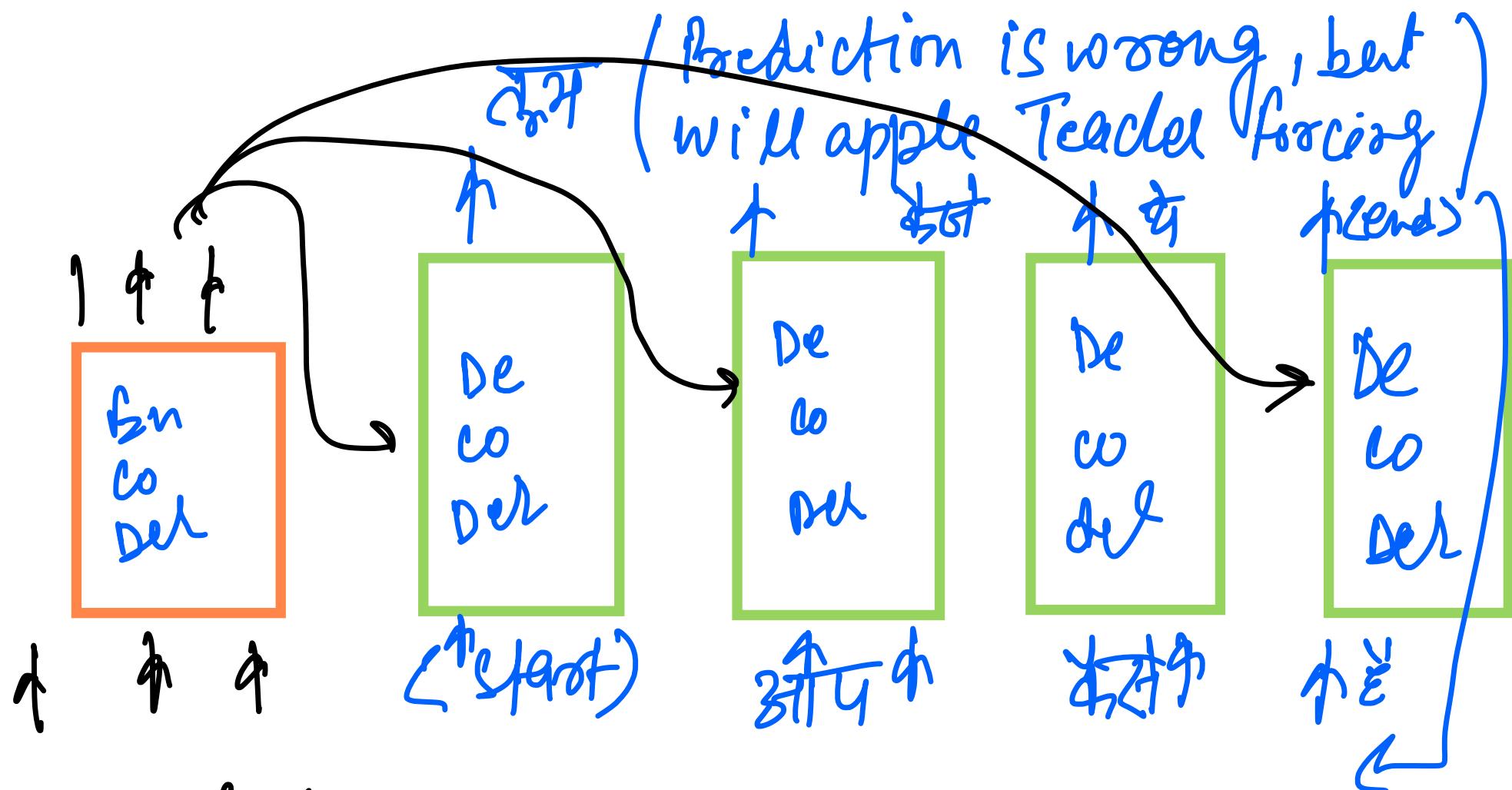
at inference time & non auto regressive
at training time.



< Training >

S.No	English Sentence	Hindi Sentence
1	How are you?	आप कैसे हैं
2	Congratulations	बधाई हो
3	Thank you	धन्यवाद

< Data set >



How are you

at training time will
send the correct next prediction

from our data set.

True Label	Prediction	
अप्प आहे वा	अमी आही पा	Will calculate the loss & backprop w/ update.

The above process in auto-regressive steps are involved. Which makes the training process very slow.

Making Parallel Training: There is no reason we should wait for previous step as the next Token coming from "DATASET". Hence "अप्पे न कै", can be passed parallelly.

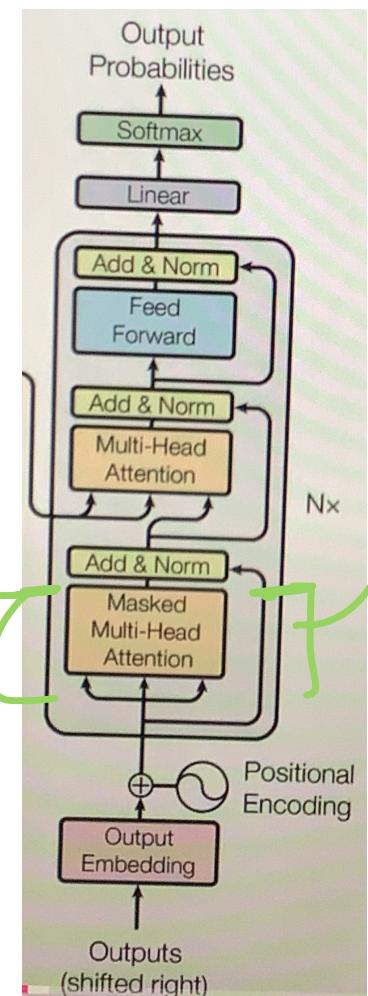
Masked Multi-Head Attention

Just for the time being

let's assume

Masked
Multi Head
Attention }
Multi head
Attention

Self Attention.



Word embedding → Self Attention → context embedding

$$\begin{aligned} \text{输出}_{ce} &= 0.8 \text{ 会议} + 0.1 \text{ 咖啡} + 0.1 \text{ 茶} \\ \text{咖啡}_{ce} &= 0.15 \text{ 会议} + 0.75 \text{ 咖啡} + 0.1 \text{ 茶} \\ \text{茶}_{ce} &= 0.1 \text{ 会议} + 0.2 \text{ 咖啡} + 0.7 \text{ 茶} \end{aligned}$$

It's possible we can get some - thing like

from eq ①, ② & ③ to represent

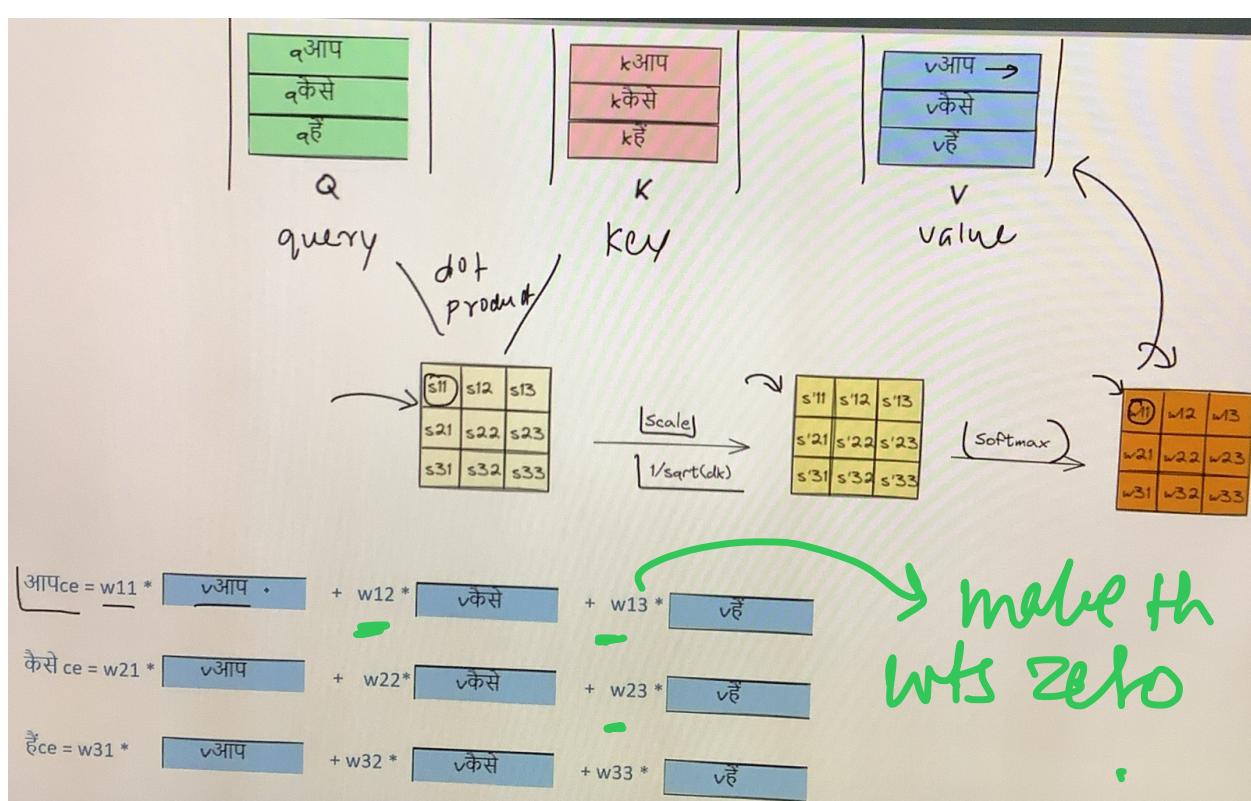
a word we need all the words

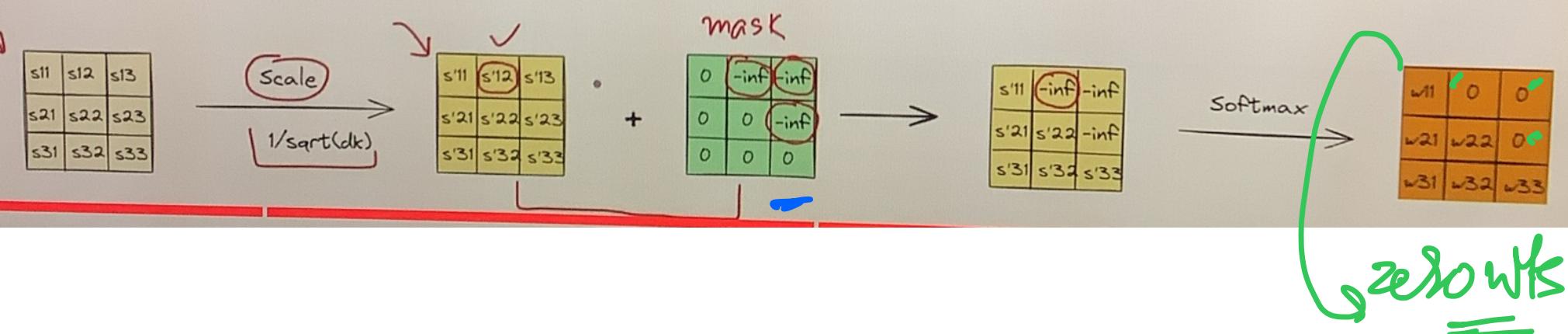
in sentence which is available at

"Train Time" but not on "Inference"

time. Hence become case of "Data

leak"





CROSS ATTENTION

This first instance

where Encoder &

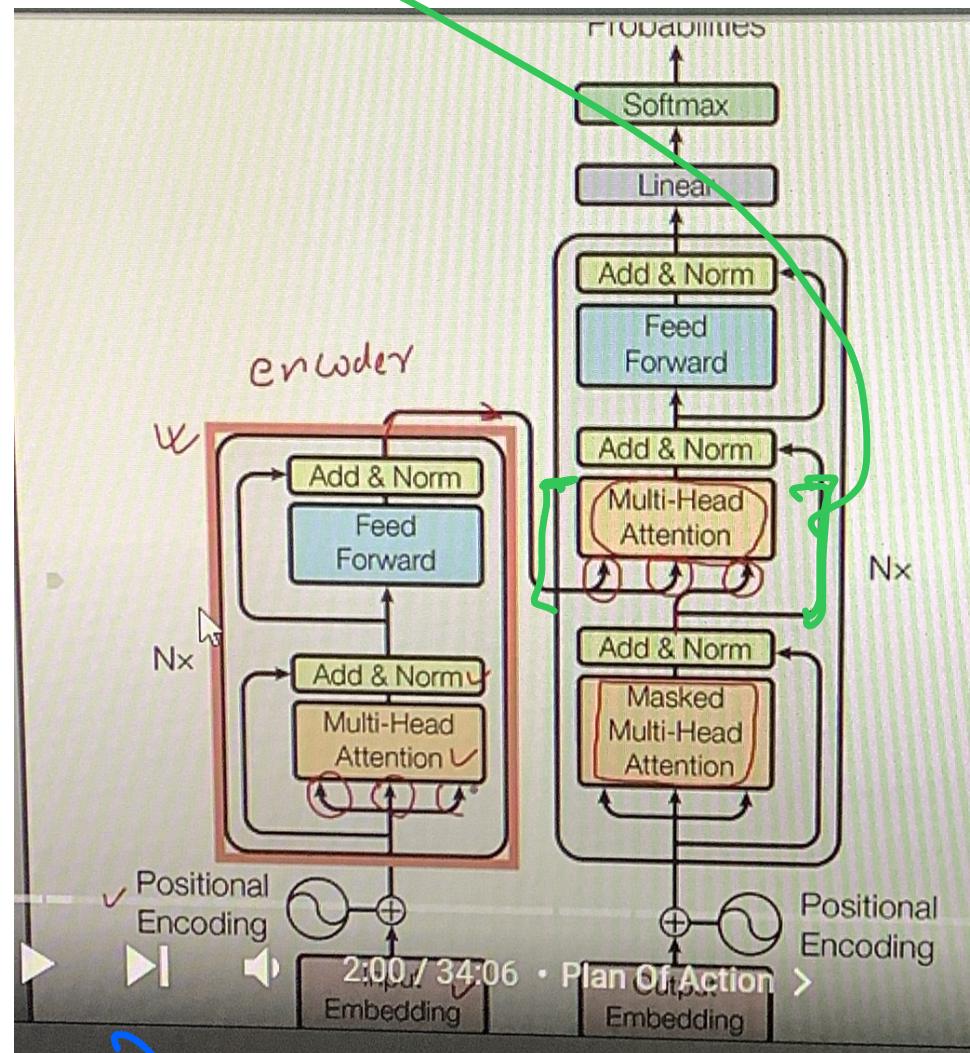
Decoder input

meets. (Cross Attention).

Let's take an example translation:-

I like eating ice-cream

→ मुझे ज्ञानकीम दोनों पसंग है।



When we are looking for
correlation between "Two Different"

sequences
Ex: →

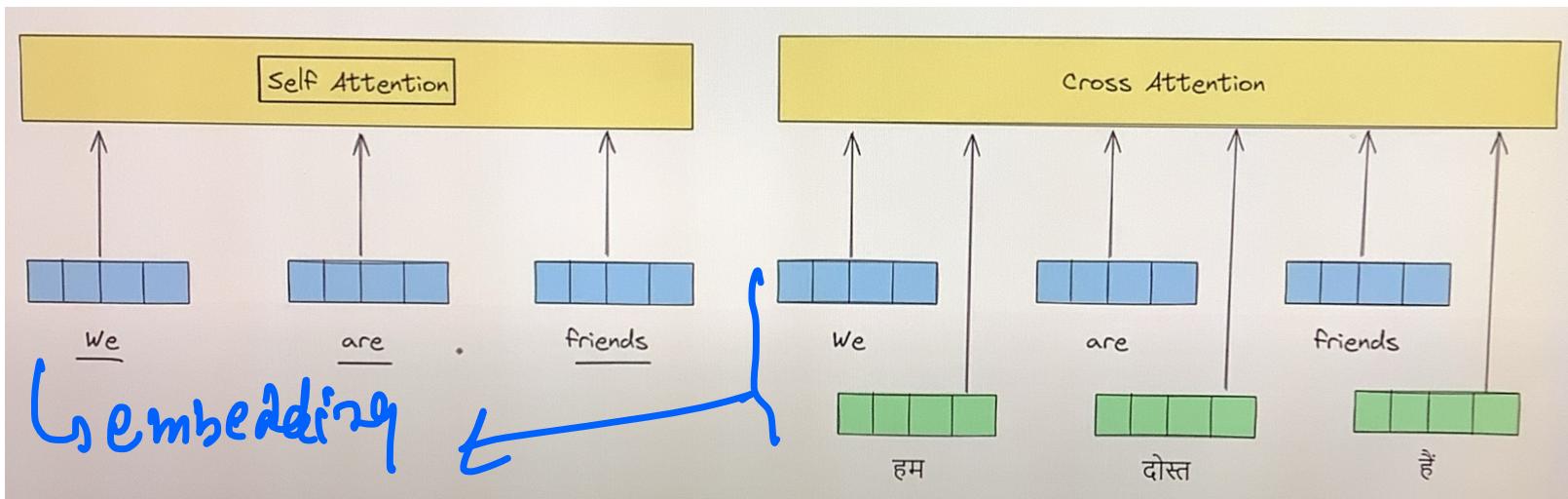
This mechanism
is known

	मुझे	आइसक्रीम	खाना	पसंद	है
I	●	•	•	●	•
like	●	●	●	●	●
eating	●	●	●	●	●
ice-cream	●	●	●	●	●

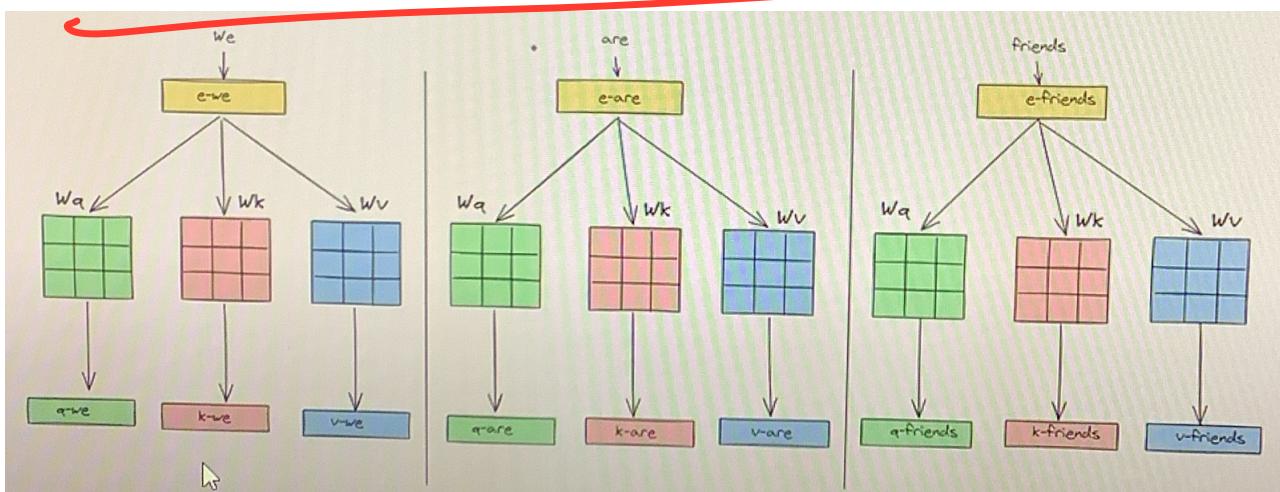
as "Cross Attention".

Cross attention is a mechanism used in transformer architectures, particularly in tasks involving sequence-to-sequence data like translation or summarization. It allows a model to focus on different part of an input sequence when generating output seq.

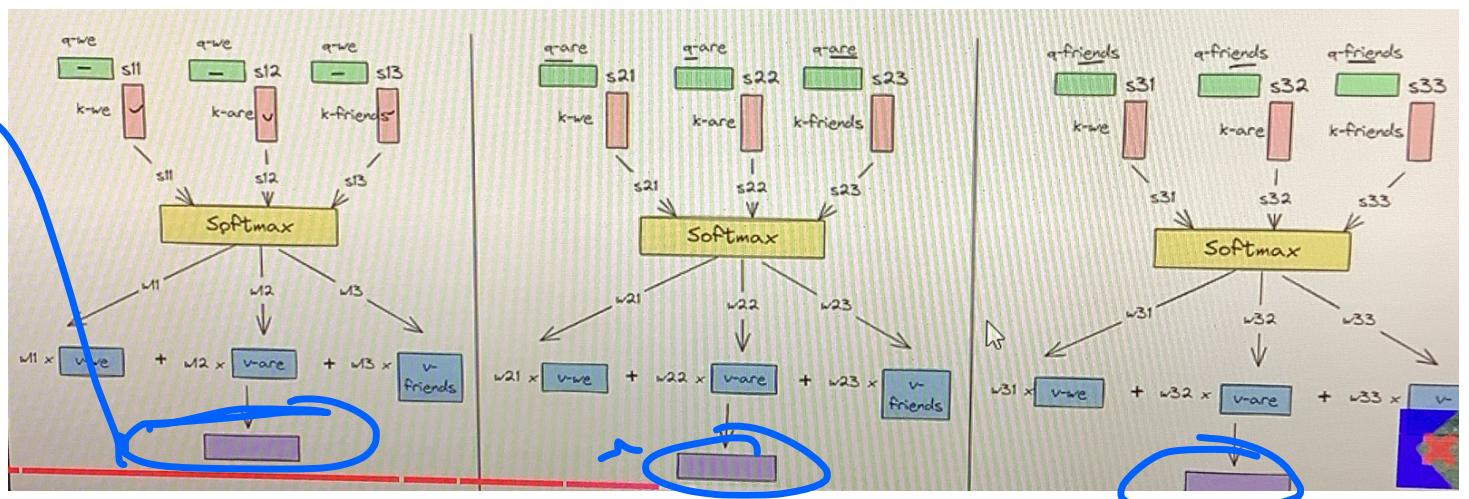
Self VS Cross Attention = A high level view.



Self Attention Quick View =



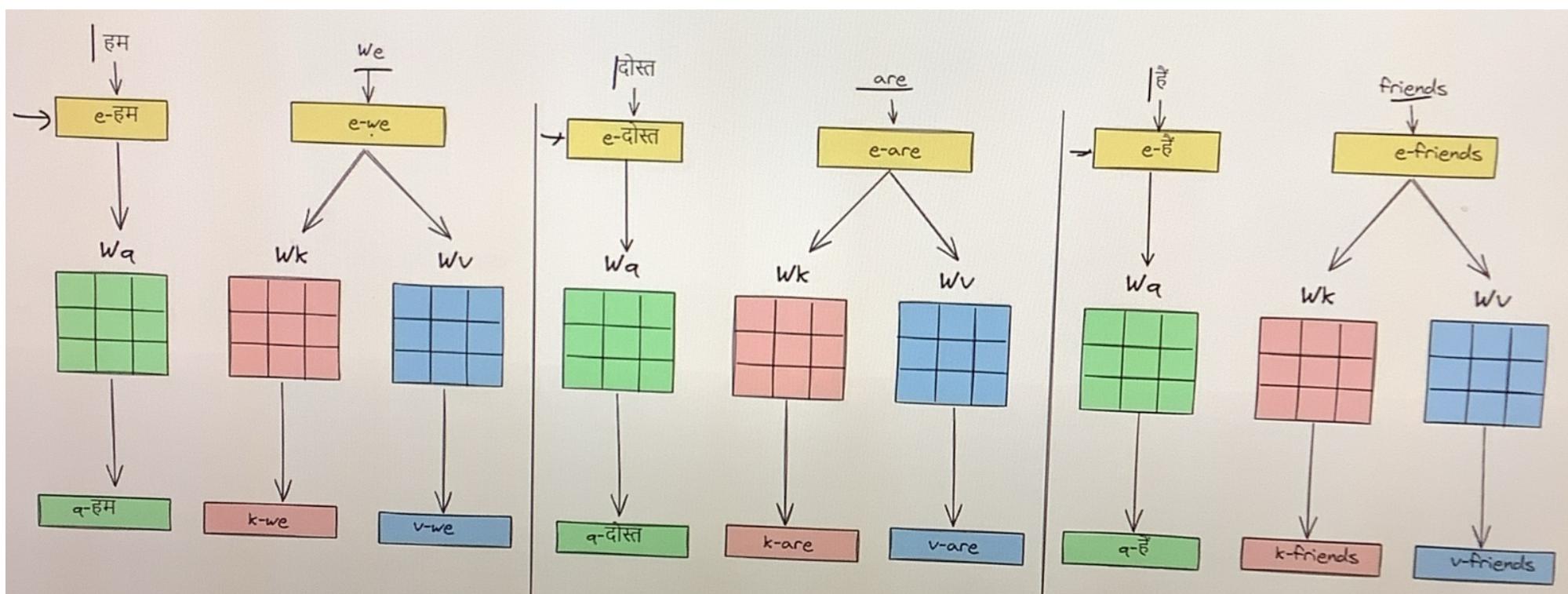
Context
embedding

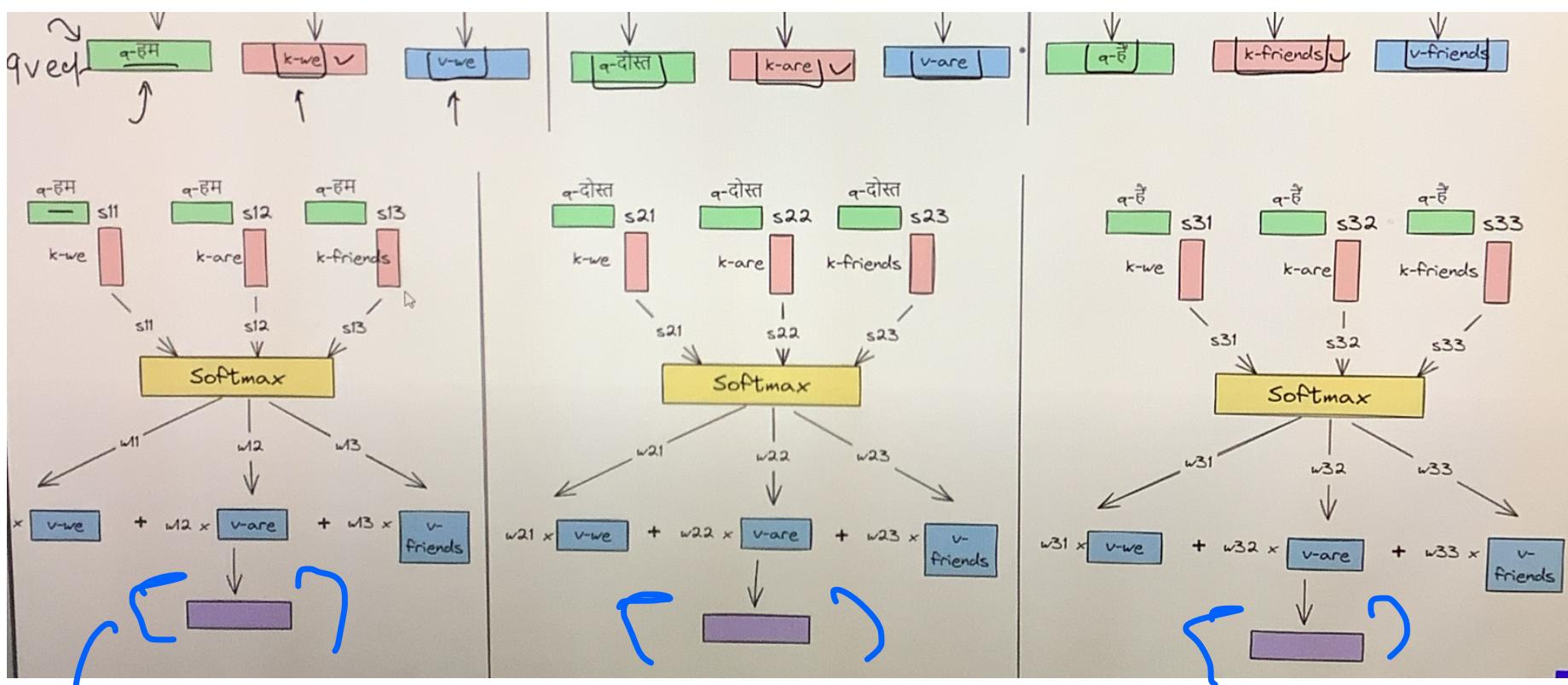


We are doing translation

English (Input) \rightarrow Hindi (Output)

* query (come from output sequence)
Key & Value (Input sequence).





→ final context embedding

④ Number of output from "Cross Attention"

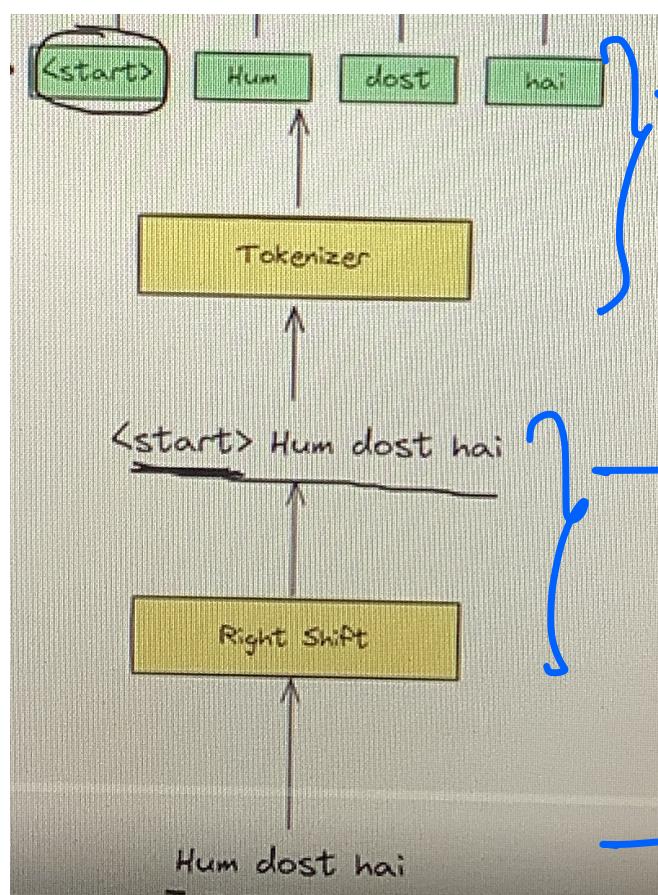
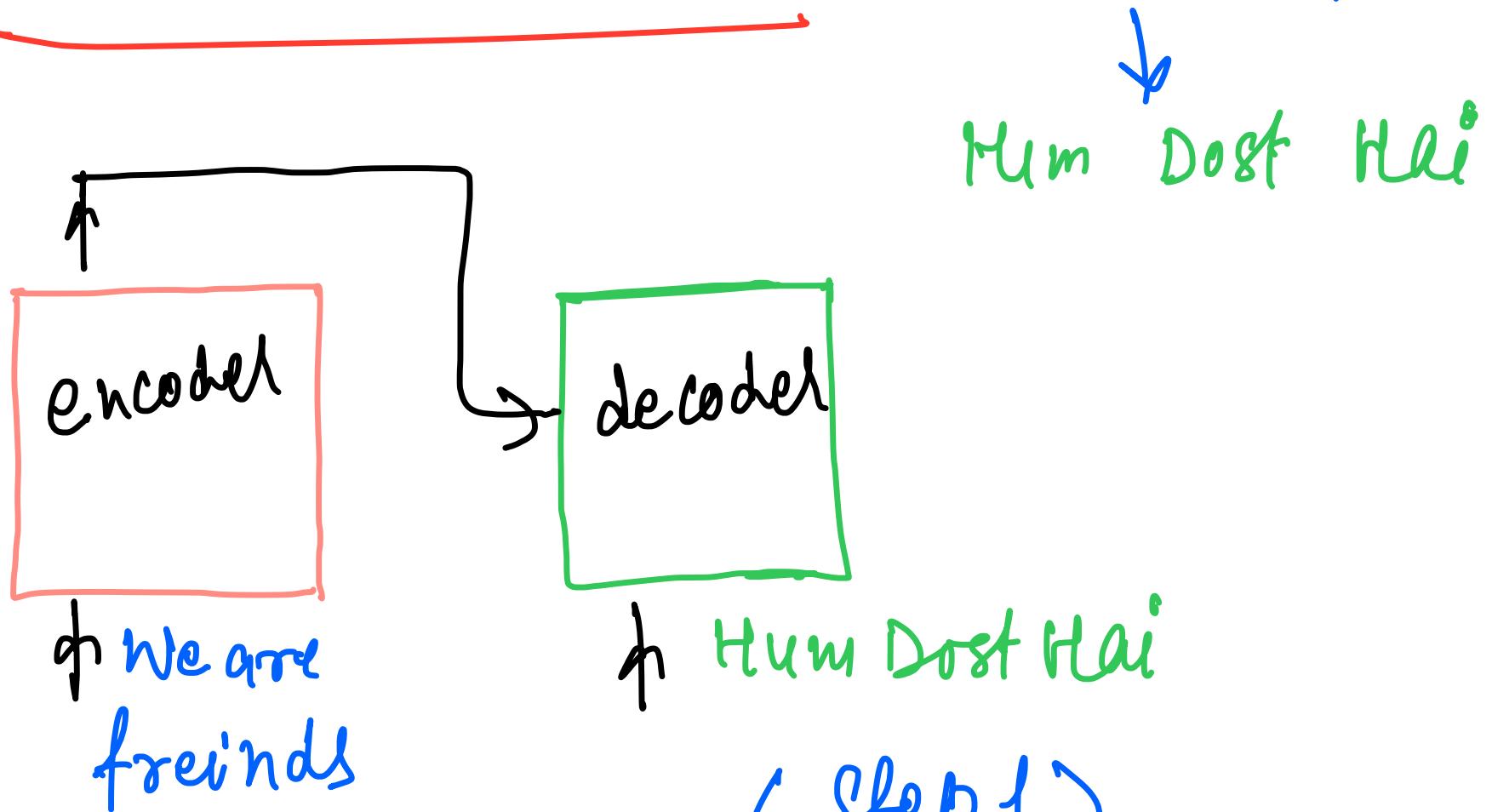
is == No. of output sequences

④ Translation / ④ Multi-model

(Image captioning, text to speech)

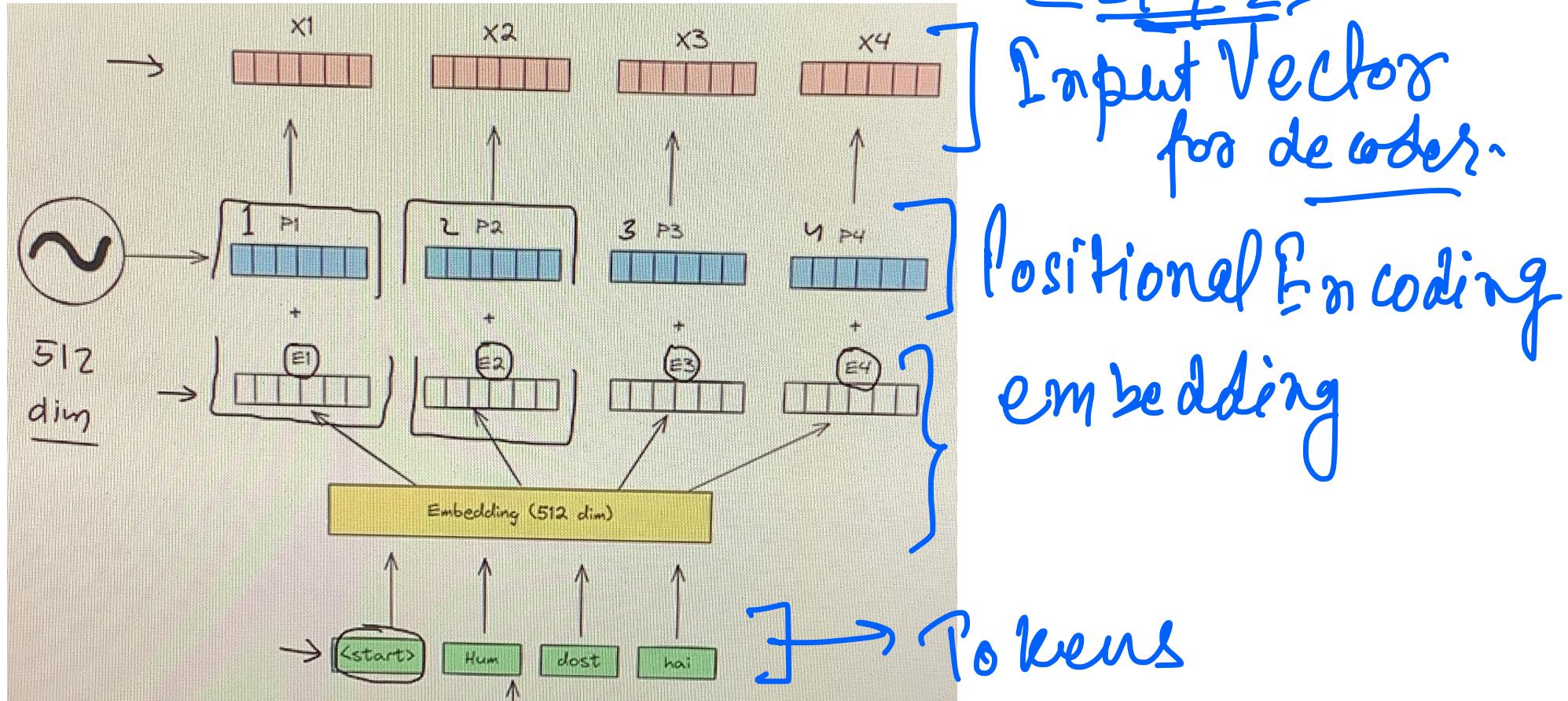
(text to image).

Decoder Architecture of "We are friends"

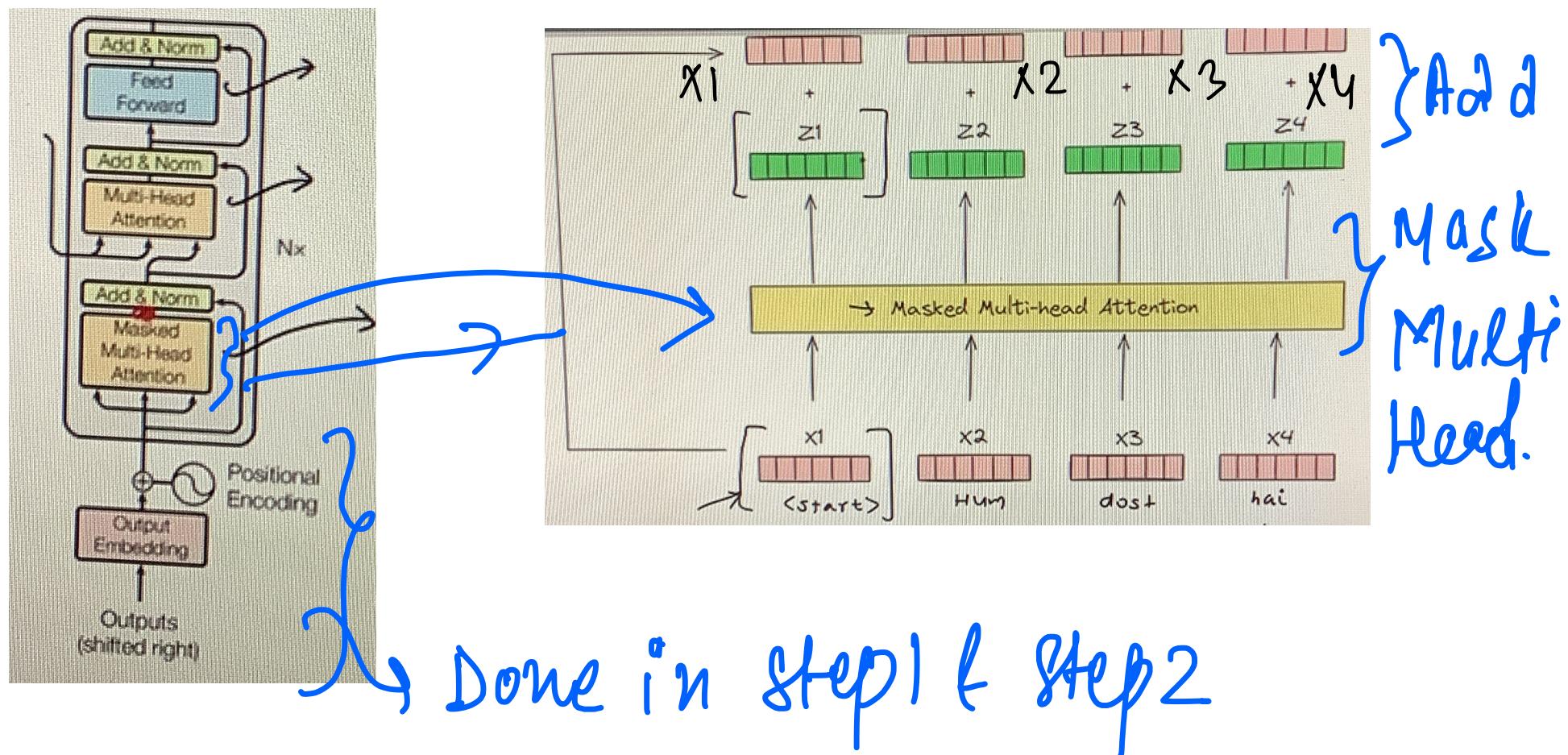


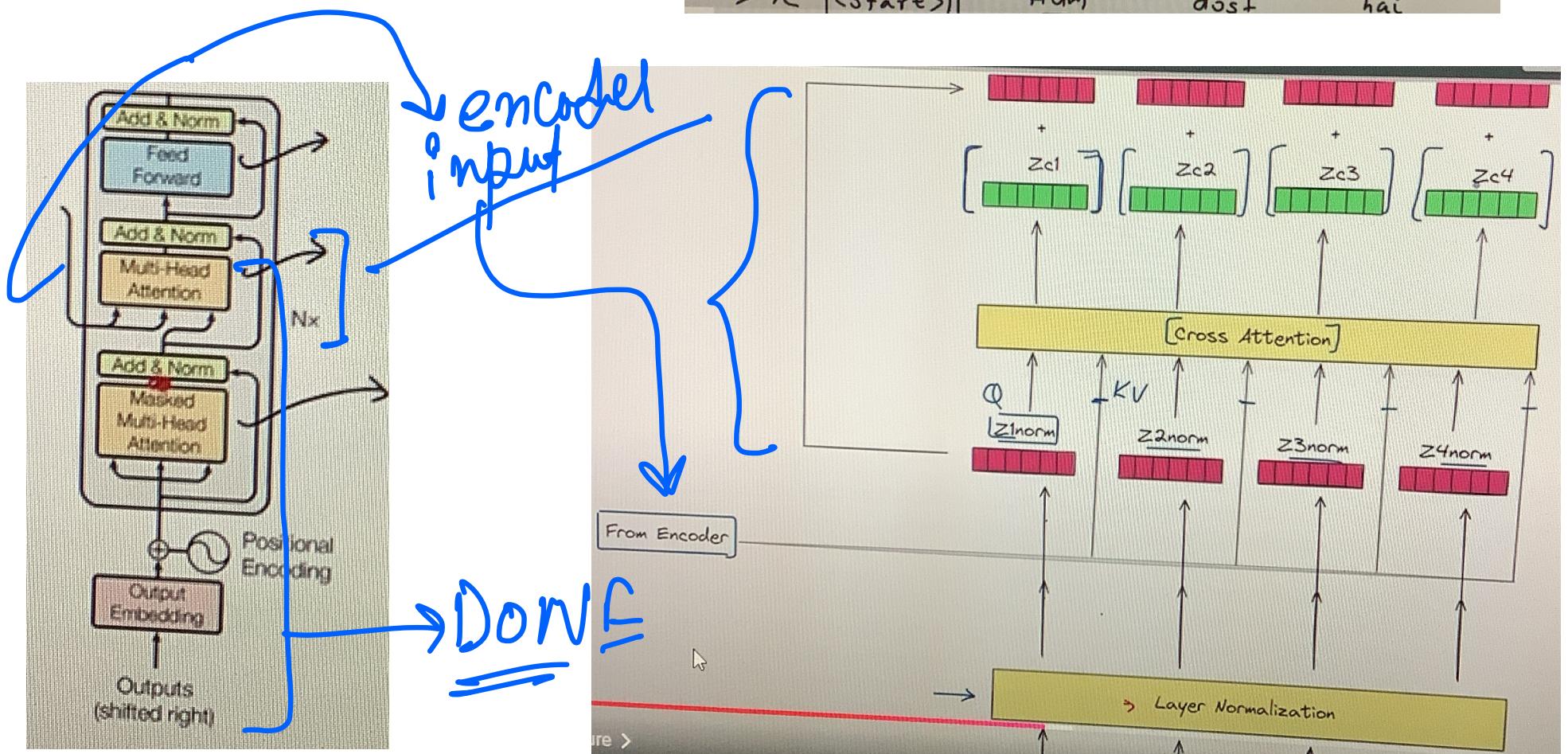
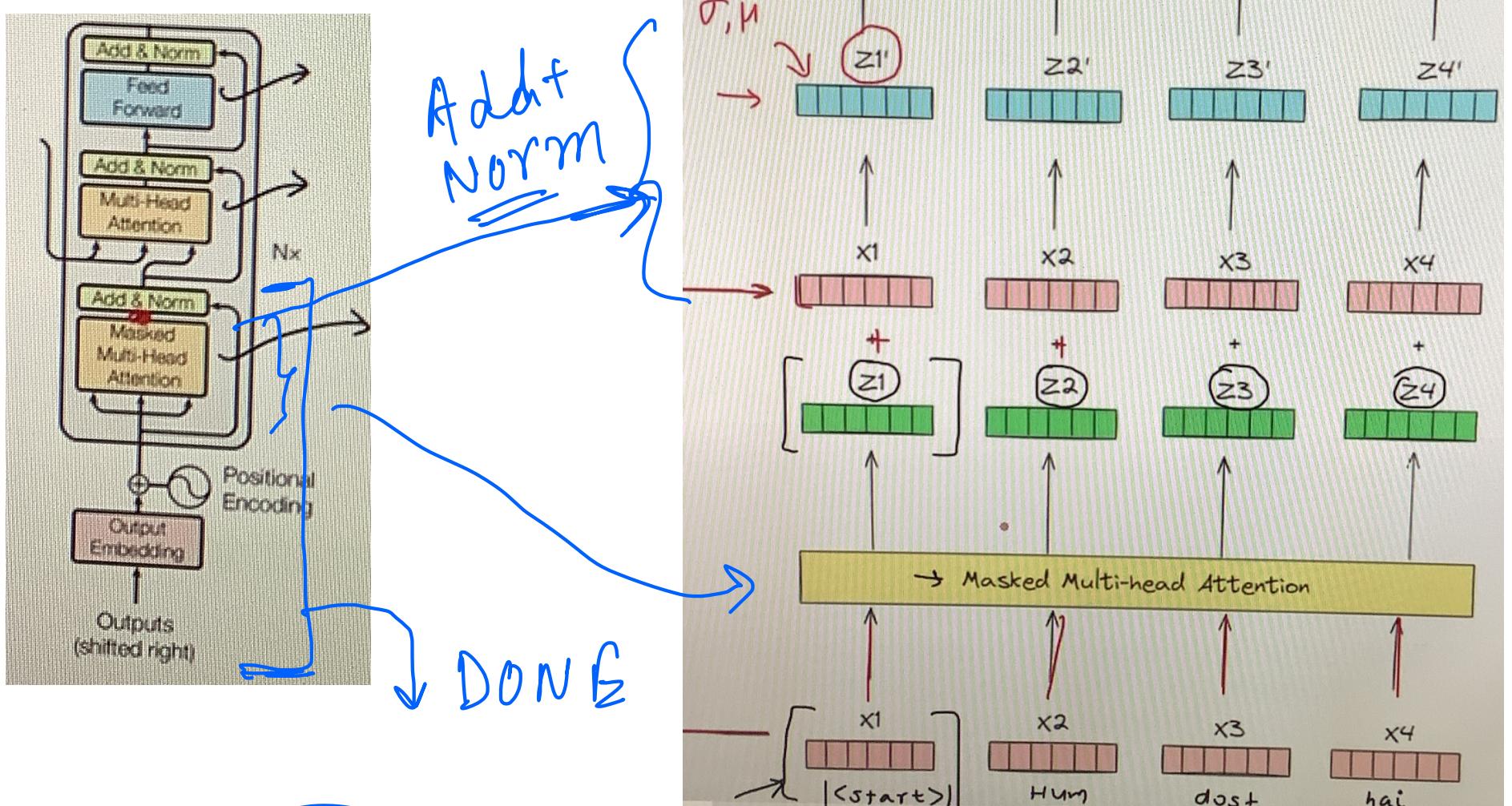
→ Right shift by append
< start > token.

→ Output Sequence



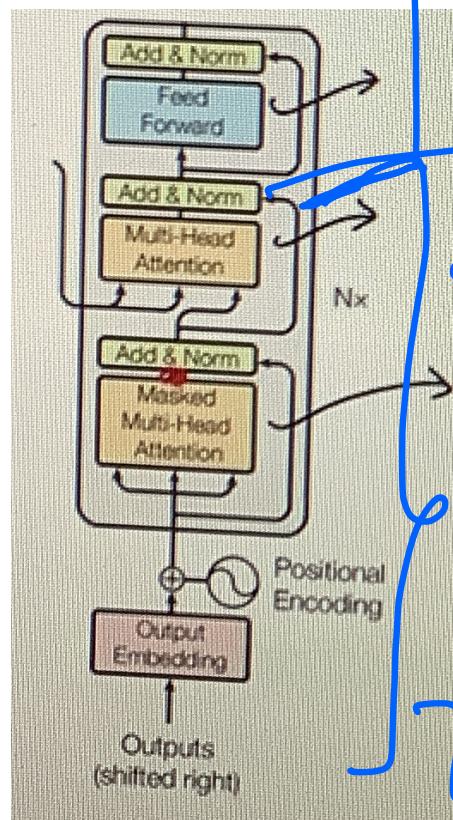
<Step 3>



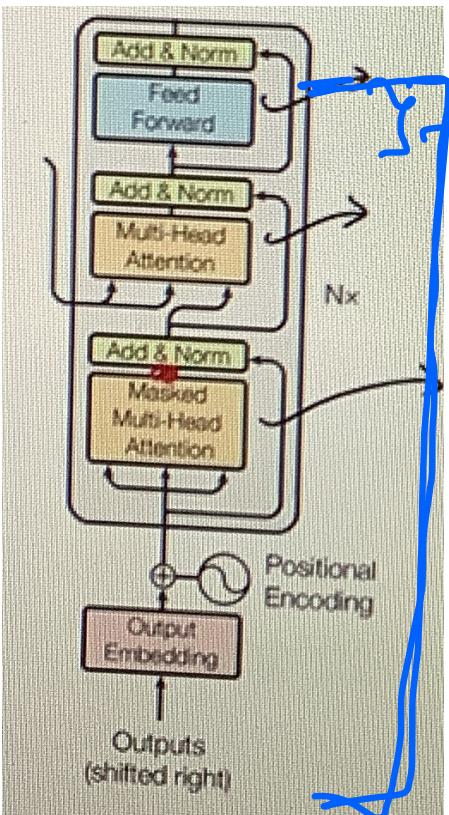
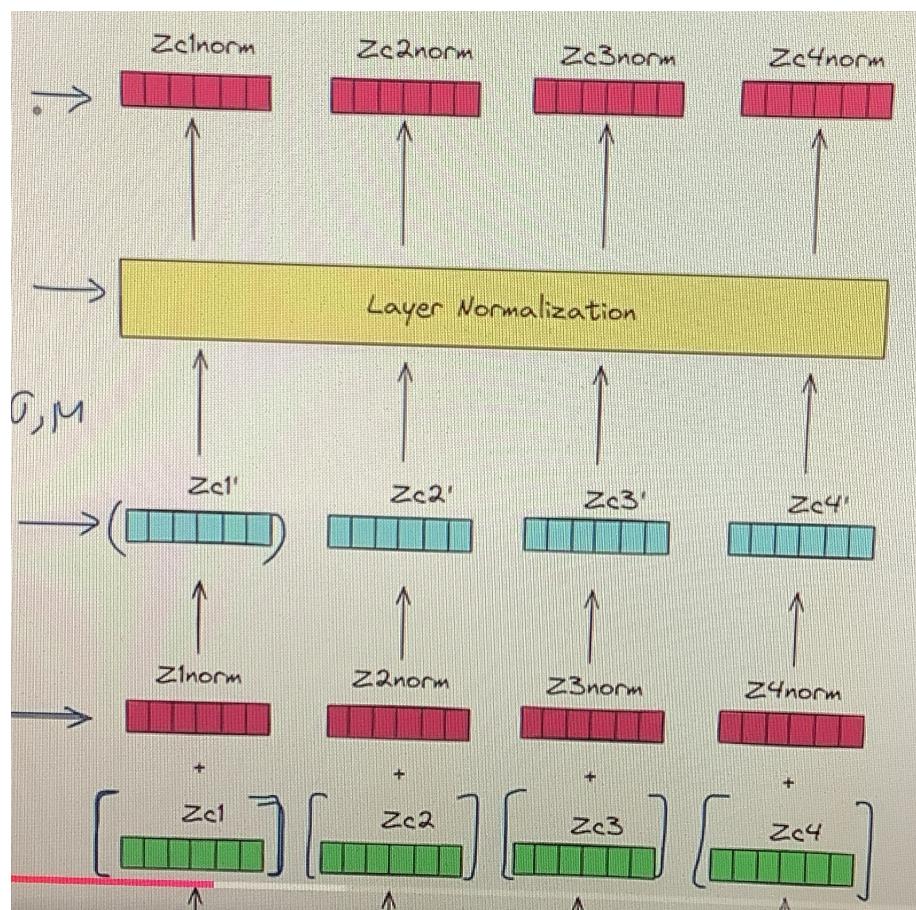


Output sequence will make query
K & V will be from Encoder!

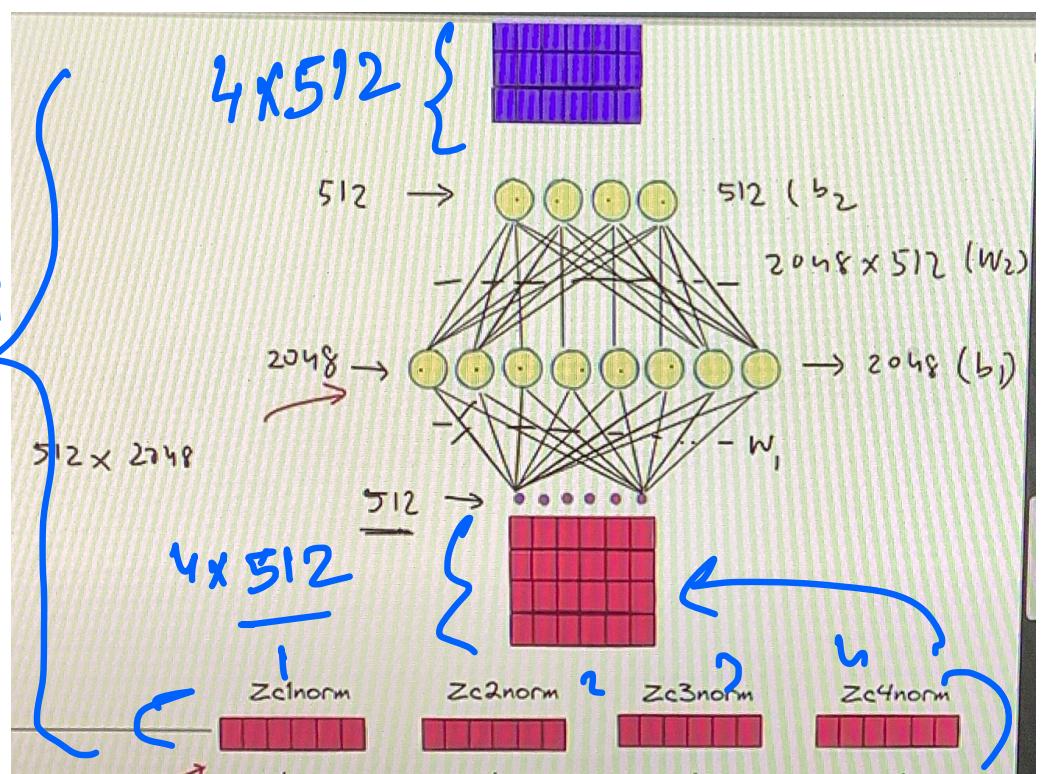
Now Add & Norm:

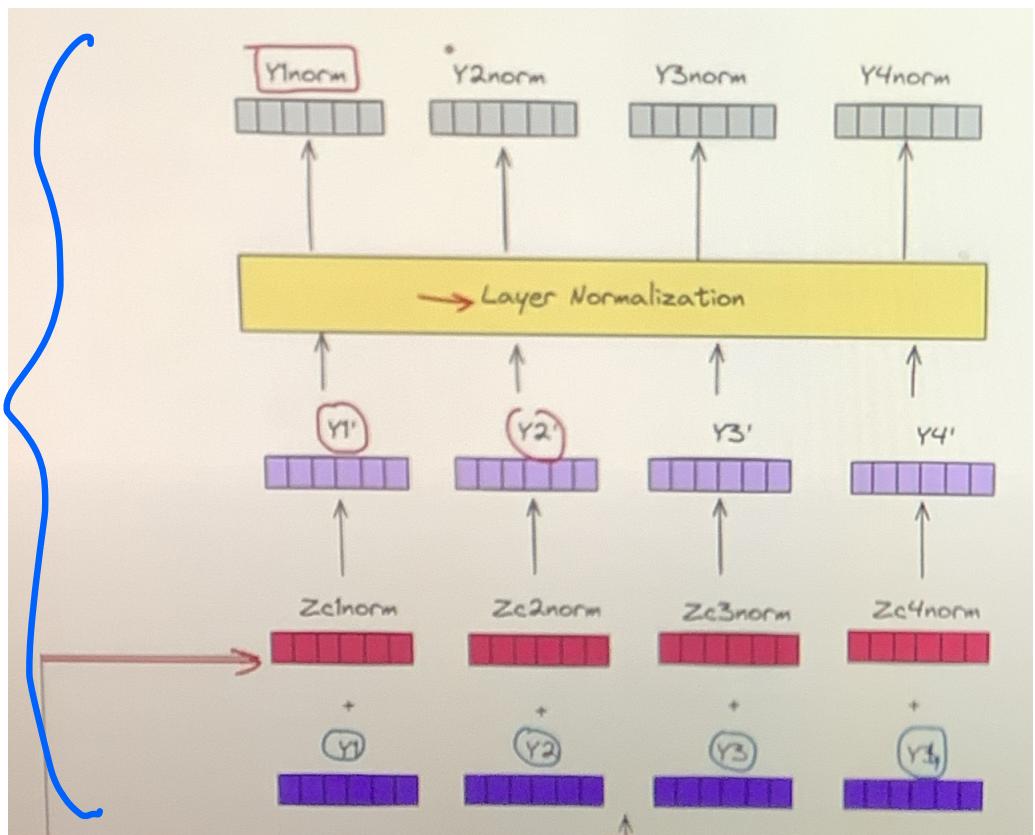
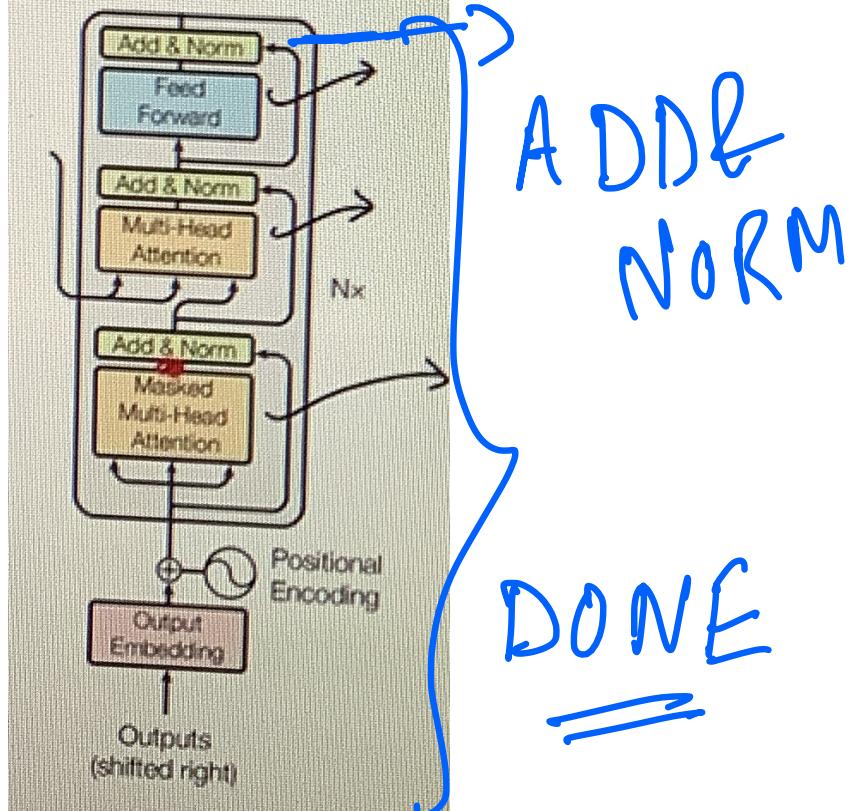


Add & Norm
DONE



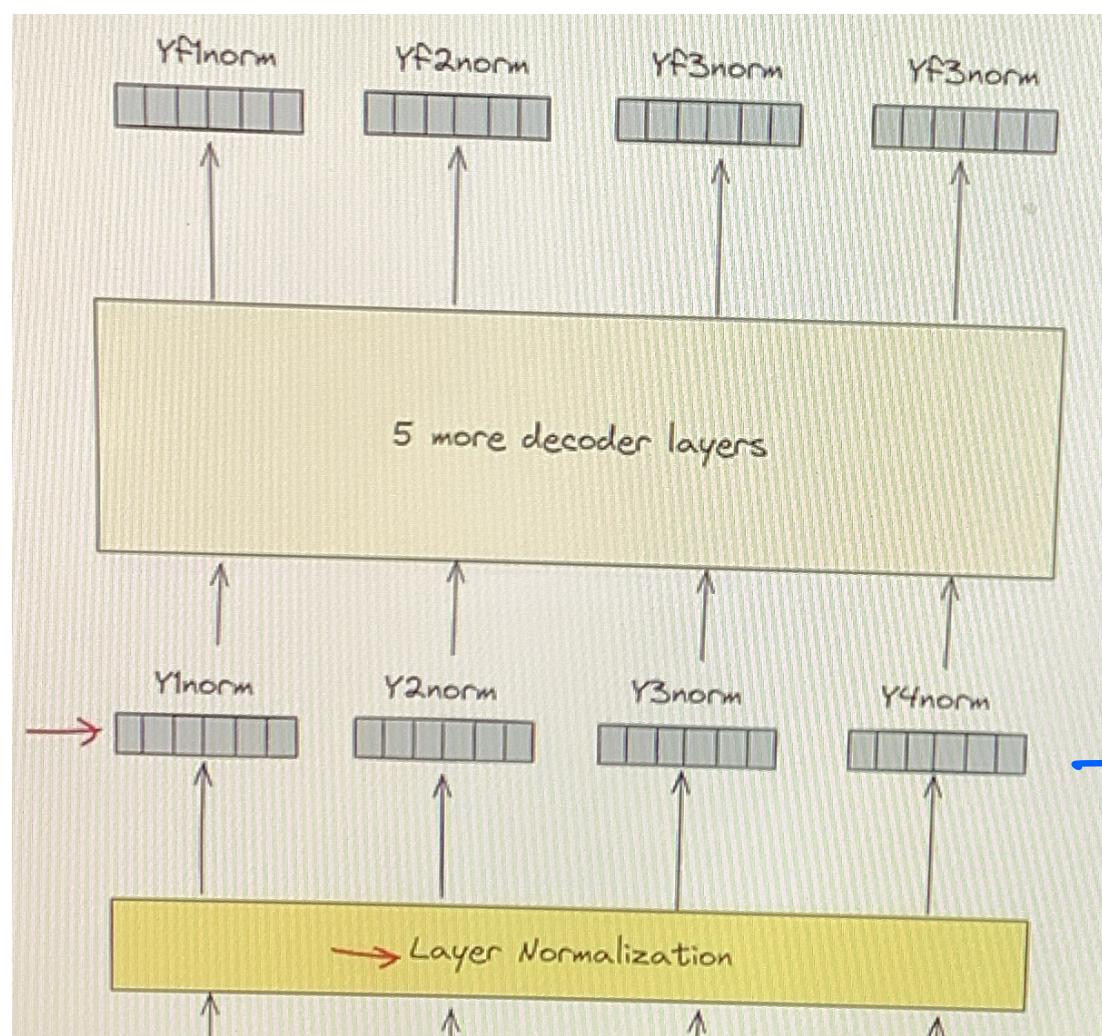
Feed forward
DONE





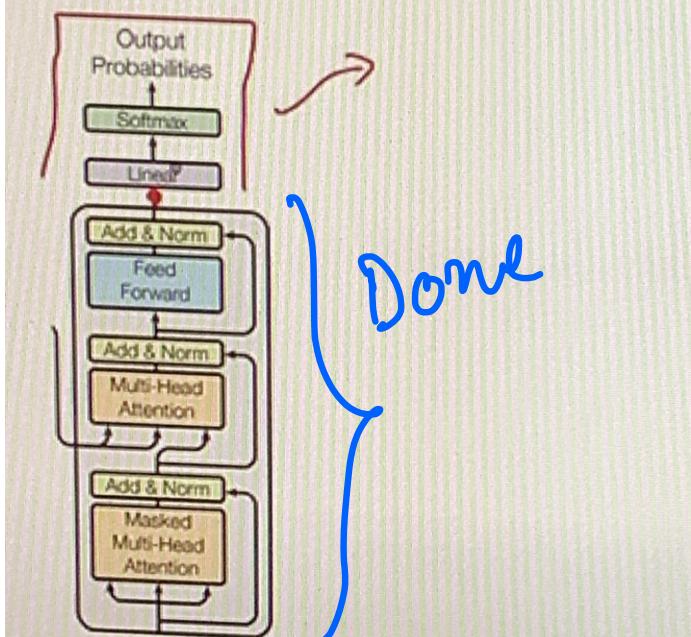
$(Y_{1norm}, Y_{2norm}, Y_{3norm} \text{ f } Y_{4norm})$

is Output from FIRST decoder.



Same operation
will repeat.

→ Output from first
decoder



Done

Now last linear & softmax.

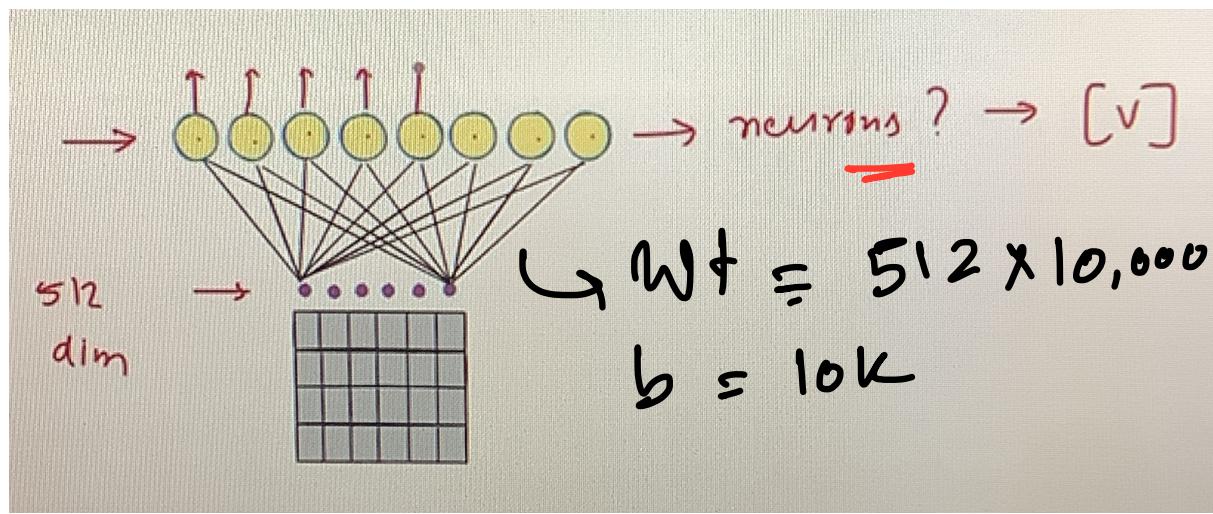
Dataset:-

I am good
Mai badiya
hoon
:
:
:
:
5000 samples.

We count unique words in Output seq.

say in 5K hindi samples we have
10K hindi words. So the

linear layer will be having 10K
neurons.



$$V = \text{No. of neurons} \\ \equiv 10K.$$

each index of neurons represent a word.

Transformer During Inference :-

