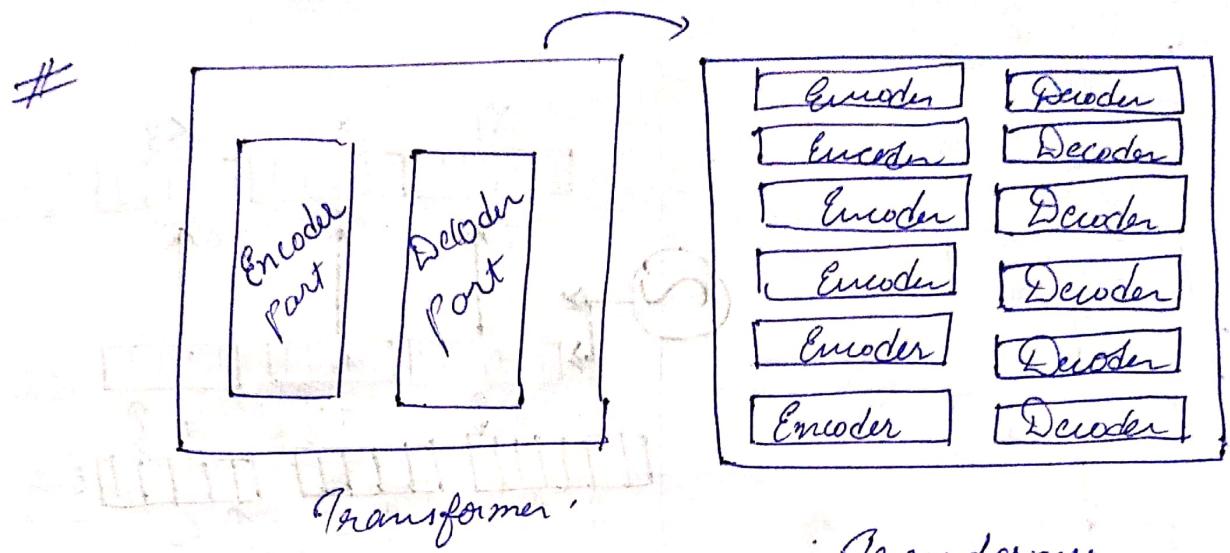
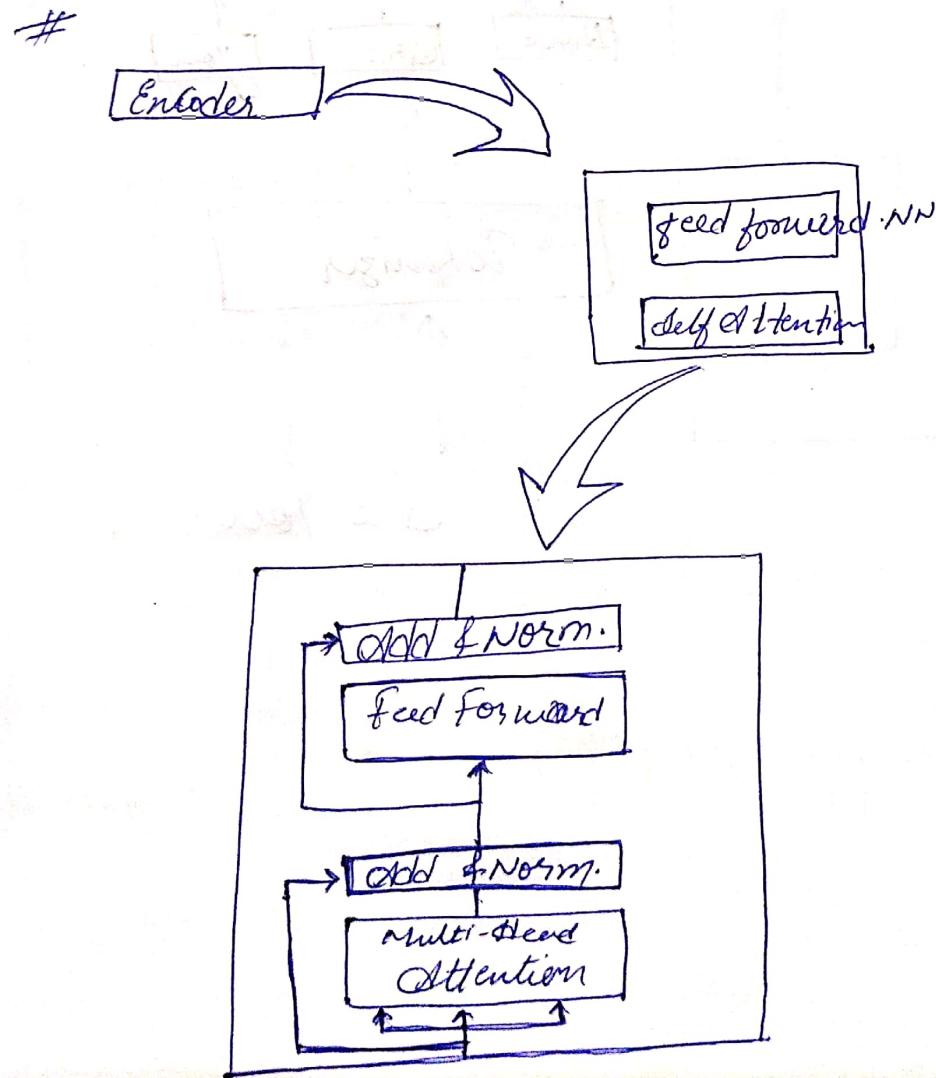
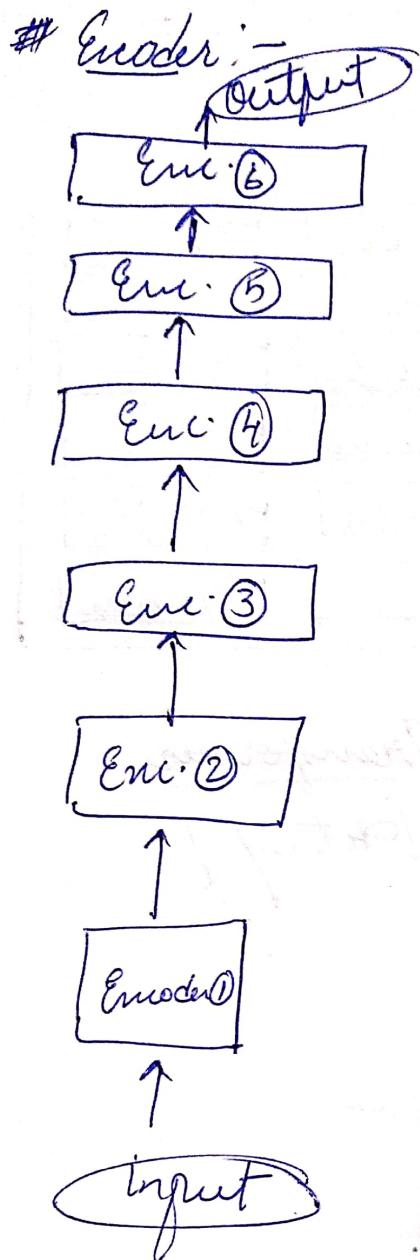


Transformer Architecture

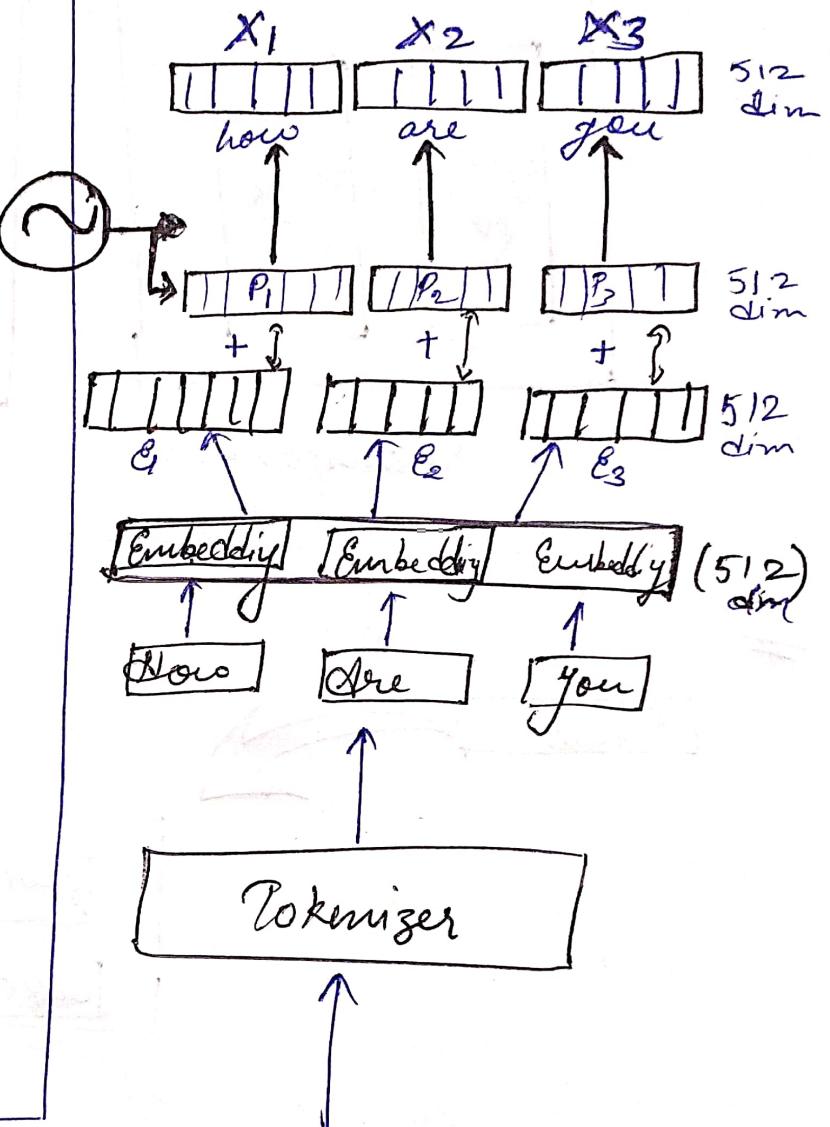


Transformer
[Actual]



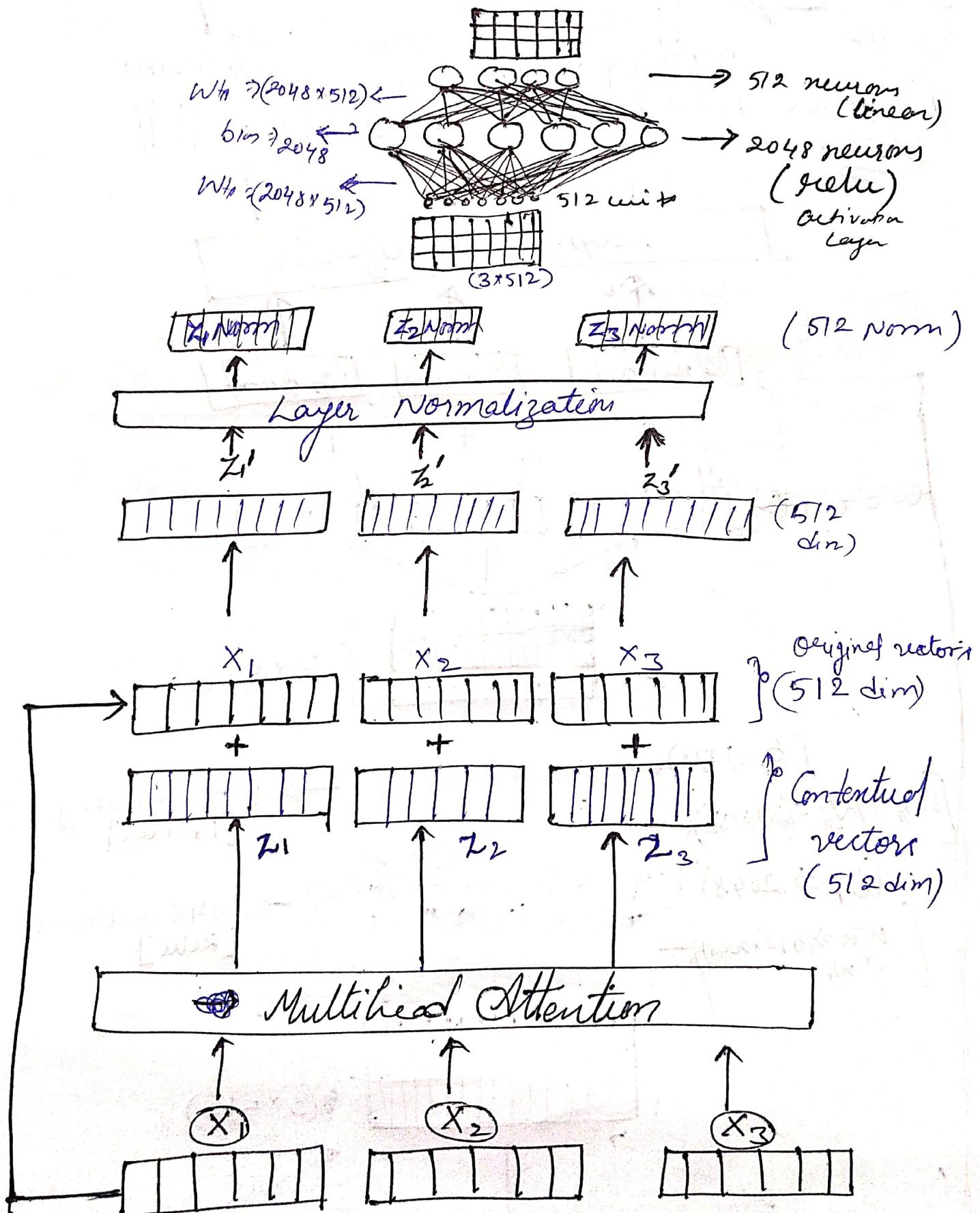


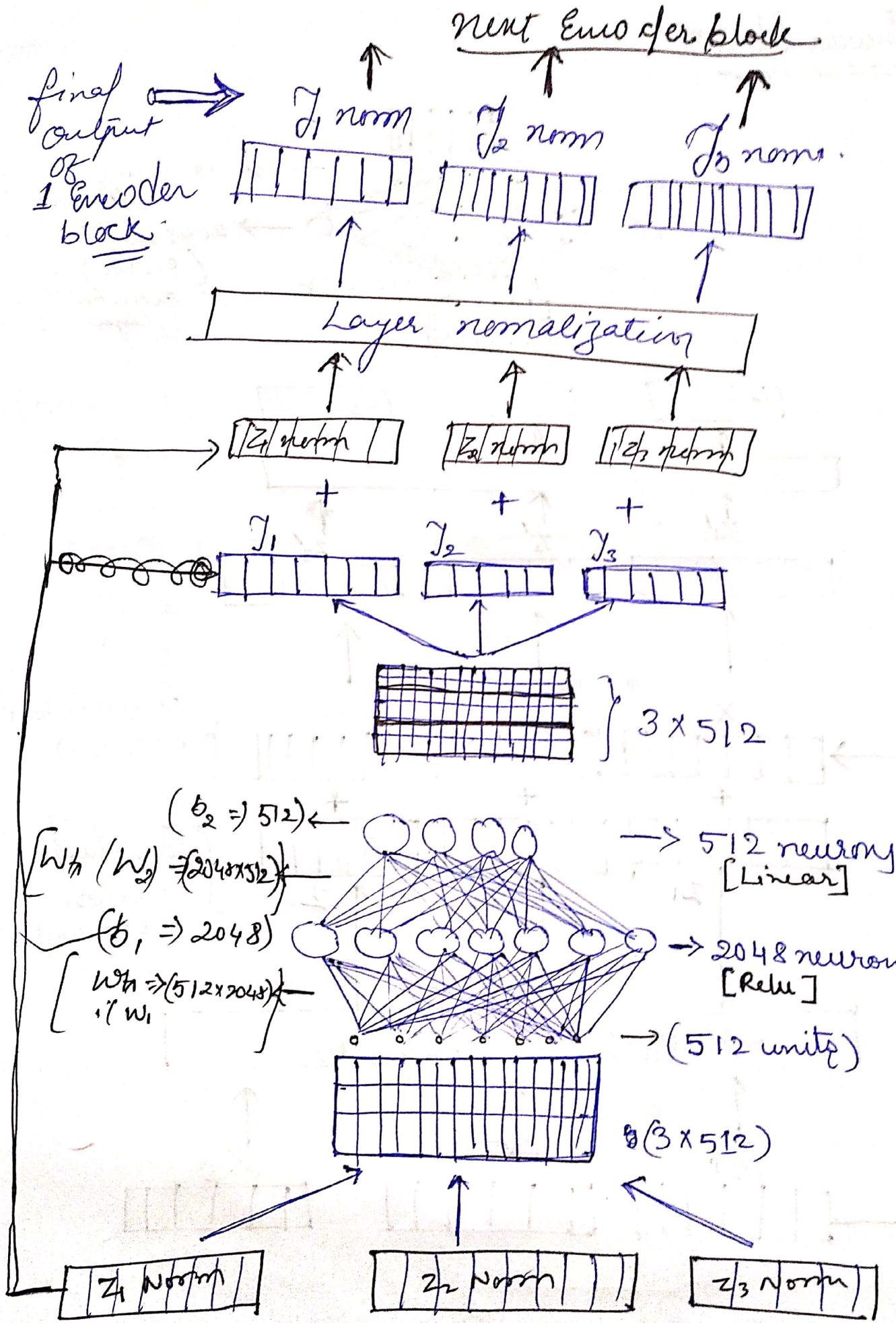
Input block :-



How are you.

Encoder block





Equations :- Neural Network :-

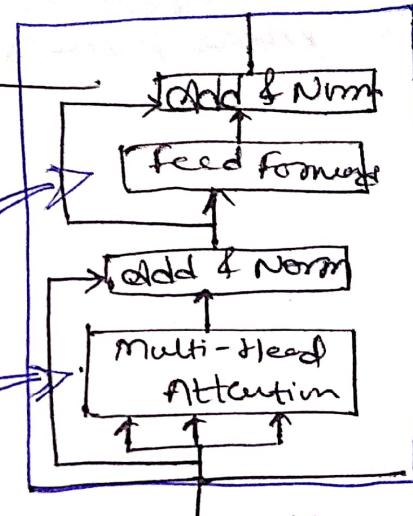
$$[\text{relu}(\text{Z}_{\text{norm}} \cdot W_1 + b_1)] \cdot W_2 + b_2$$

Why use residual connections?

Reasons to use :-

1) Stable training.

2) Original features are always preserved, if Multihead attention does not work well.



Why to use FFNN??

* Self-Attention has all the linear operations in it, that is why it doesn't capture non-linear relations that well,

so, relations well To capture the non-linear a FFNN is used.

Q Why use more than 1 encoder block?
Why 6 encoder blocks,

→ ∵ language is a very complex
to understand that's why more
Encoder blocks help in understanding
more & more deeper patterns of languages.

Masked Self Attention

- * Auto-regressive Models :-

- * The Transformer decoder is autoregressive at inference time & non-autoregressive at training time.

What is Autoregression?

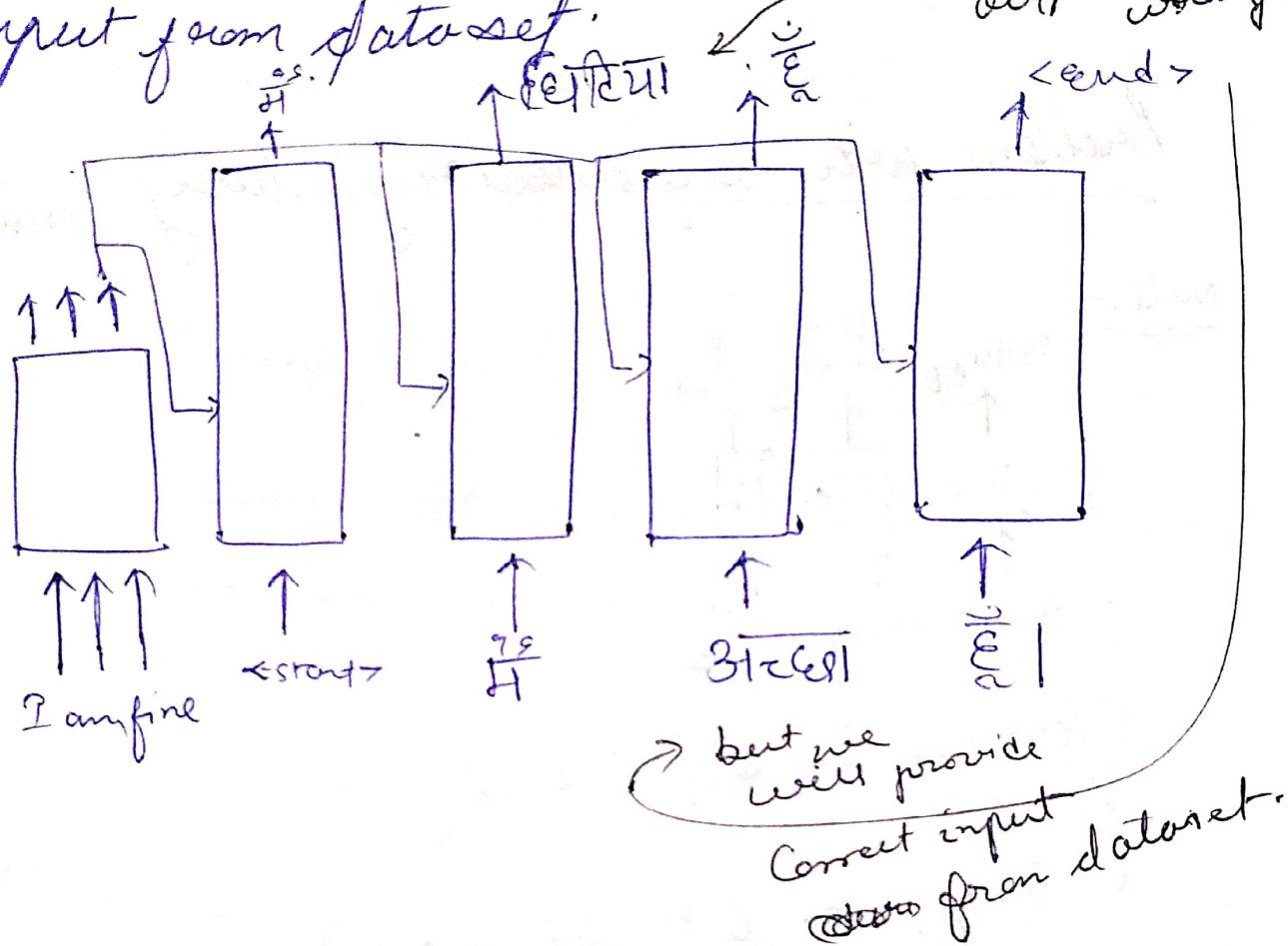
Autoregressive models are the class of models that generate data points in a sequence by conditioning each new point on a previously generated point.

- * Inference = prediction.

- * Training :-

- oo During training we have the data set with us, we need not to send the output of every step to next as it may cause false training, because if output become wrong during a step & we provide the wrong output next step, it will be trained on wrong outputs.

that is why during training whatever be the output we always provide correct input from dataset.



* Due to above reasons during training we need to be non-autoregressive.

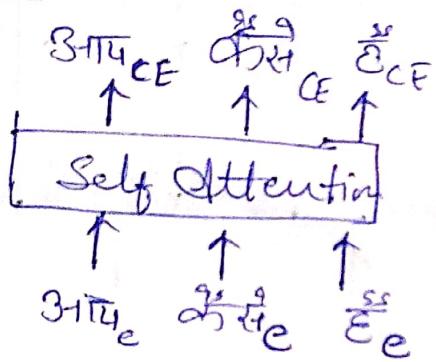
Prediction:-

∴ During prediction, ~~we don't~~ we don't have dataset with us & so that is why we need to be depended upon ~~not~~ previous values to predict next word.

* that if we need to be autoregressive during prediction.

Problem with self embedding during Prediction

Note:-



According to self Embedding:-

E.g:-

$$\begin{aligned} \overline{3T4}_{CE} &= 0.8 \overline{3T4}_C + 0.1 \overline{3T4}_F + 0.1 \overline{3T4}_E \\ \overline{3T4}_{FT} &= 0.25 \overline{3T4}_C + 0.9 \overline{3T4}_F + 0.25 \overline{3T4}_E \\ \overline{3T4}_E &= 0.5 \overline{3T4}_C + 0.5 \overline{3T4}_F + 0.9 \overline{3T4}_E \end{aligned}$$

Now During Prediction, at first time stamp we will only have $\overline{3T4}_{CE}$ (first word) generated, but self attention to make

Contextual Embedding of the word
will demand future words Embedding
as well & which during Prediction
we don't even know.

i.e.

During prediction also,

Current token value is dependent on
future token values,

which is not
possible during prediction,

During training

it is possible as we have the dataset
available with us

"That is why simple self
Attention Arch's will not
work here"

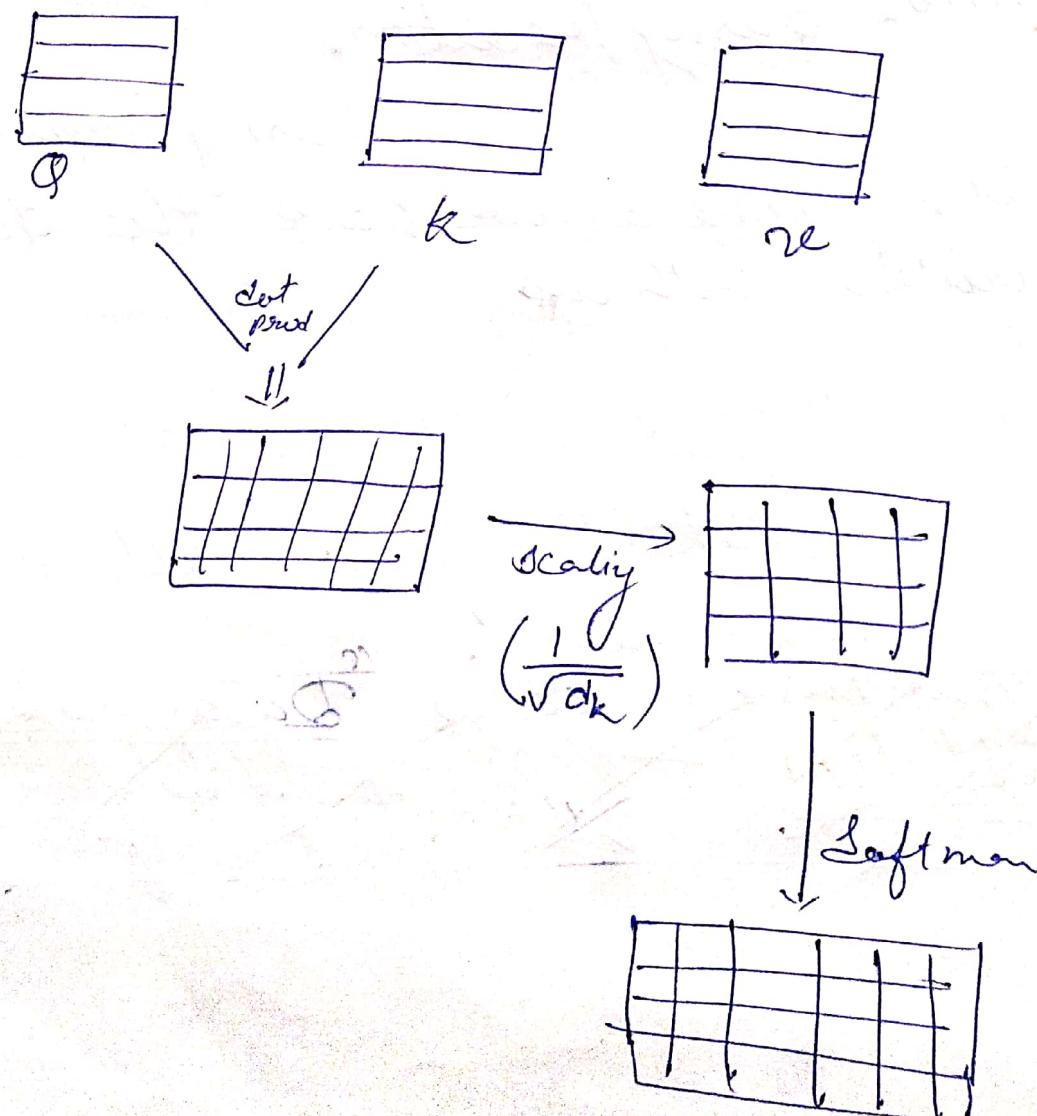
~~This issue is called "Data leakage"
as we used future data during
training but didn't have it.
Masked Self Attention.~~

* Data leakage:-

During training, the Decoder model has future tokens available & that is why it was able to perform training but during prediction it is not available, this issue is called data leakage.

Solution => " Masked Self Attention"

Let:-



Now,

v_1
v_2
v_3
v_4
v_{52}

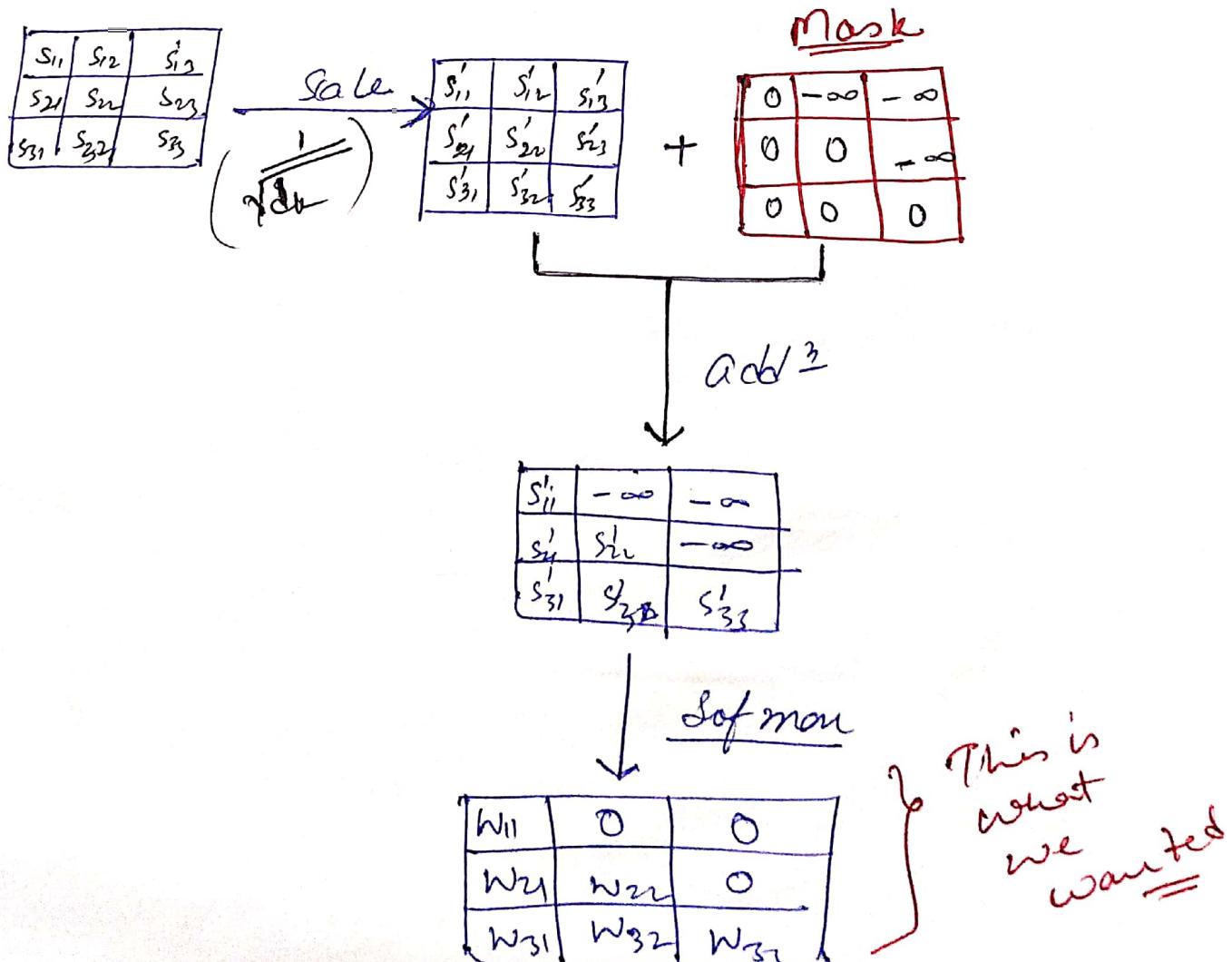
w_{11}	w_{12}	w_{13}

" v ".

$$\text{Output}_{CE} = w_{11} * (v_{3114}) + w_{12} * (v_{2510}) + w_{13} * (v_{82})$$

∴ During prediction these two will not be available.

Masked self Attention to avoid future tokens.



Now, in this matrix we can see the future weights are "0", this will avoid future tokens.



Cross Attention.

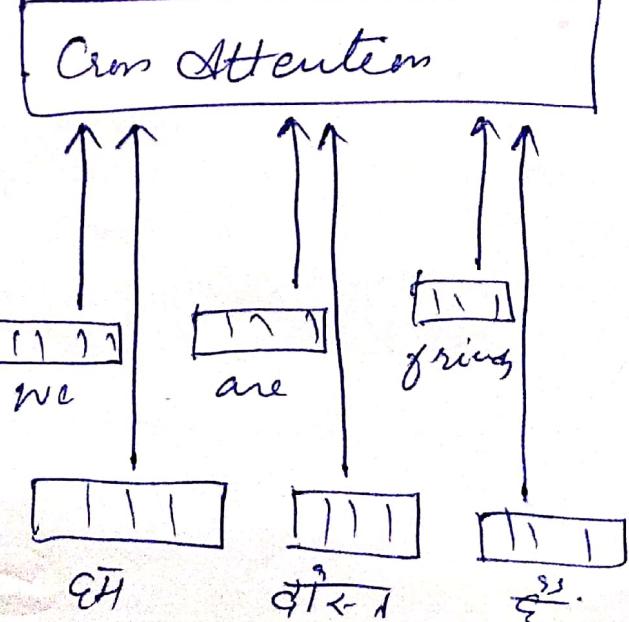
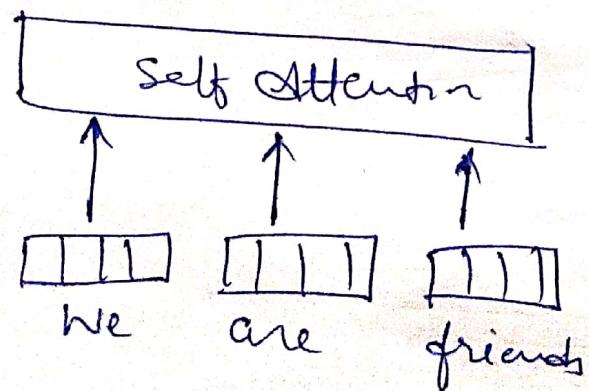
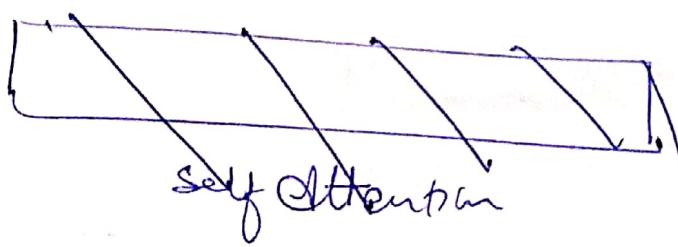
→ Cross-attention is a mechanism used in transformer architecture, particularly in tasks involving sequence to sequence data like translation or summarization. It allows a model to focus on different parts of an input sequence when generating an output sequence.

Self Attention

vs

Cross Attention

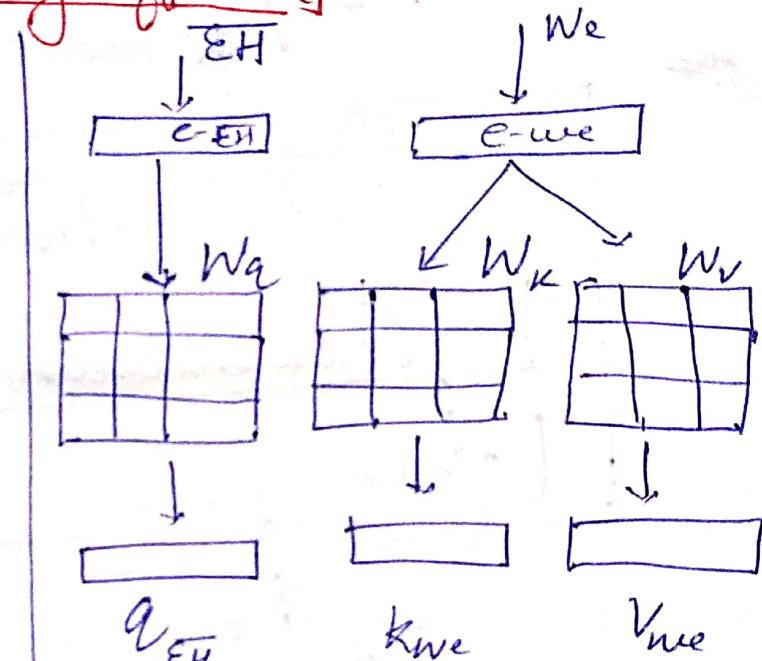
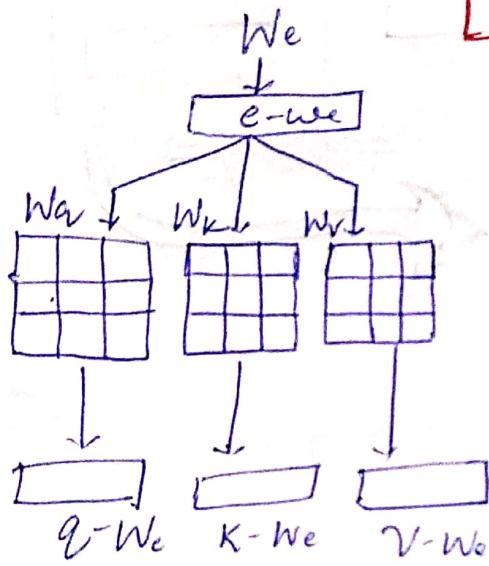
[Input Difference]



So,

Cross attention need both input & output sequences embedding while self attention need only input sequence.

Processing Difference



If further processing are same.

Further processing are same.

* In Cross Attention, the

Output

↳ the query vector

→ Comes from

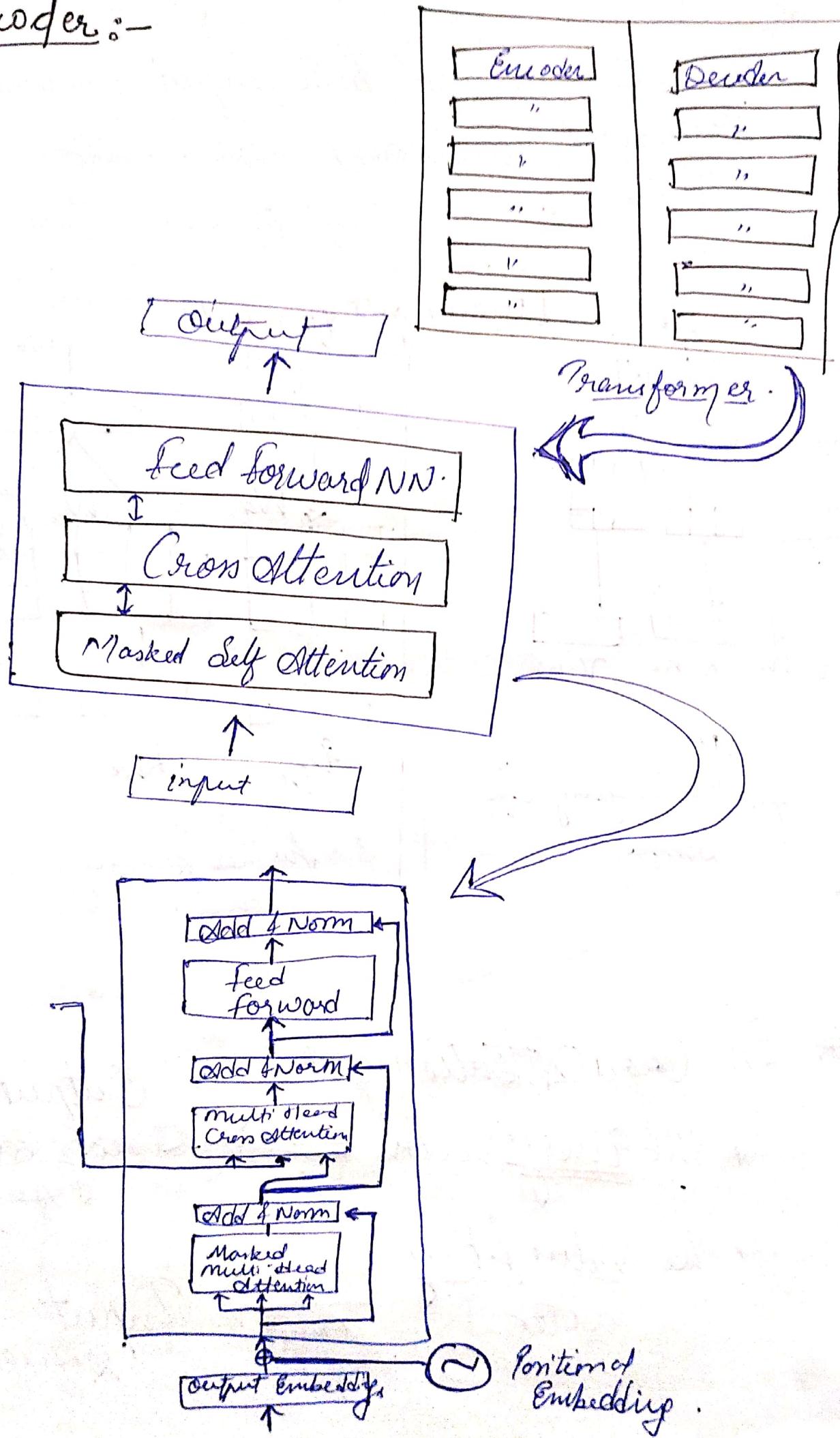
Decoded Sequence

↳ the value & key vector

→ Comes from

Input Sequence

Decoder :-



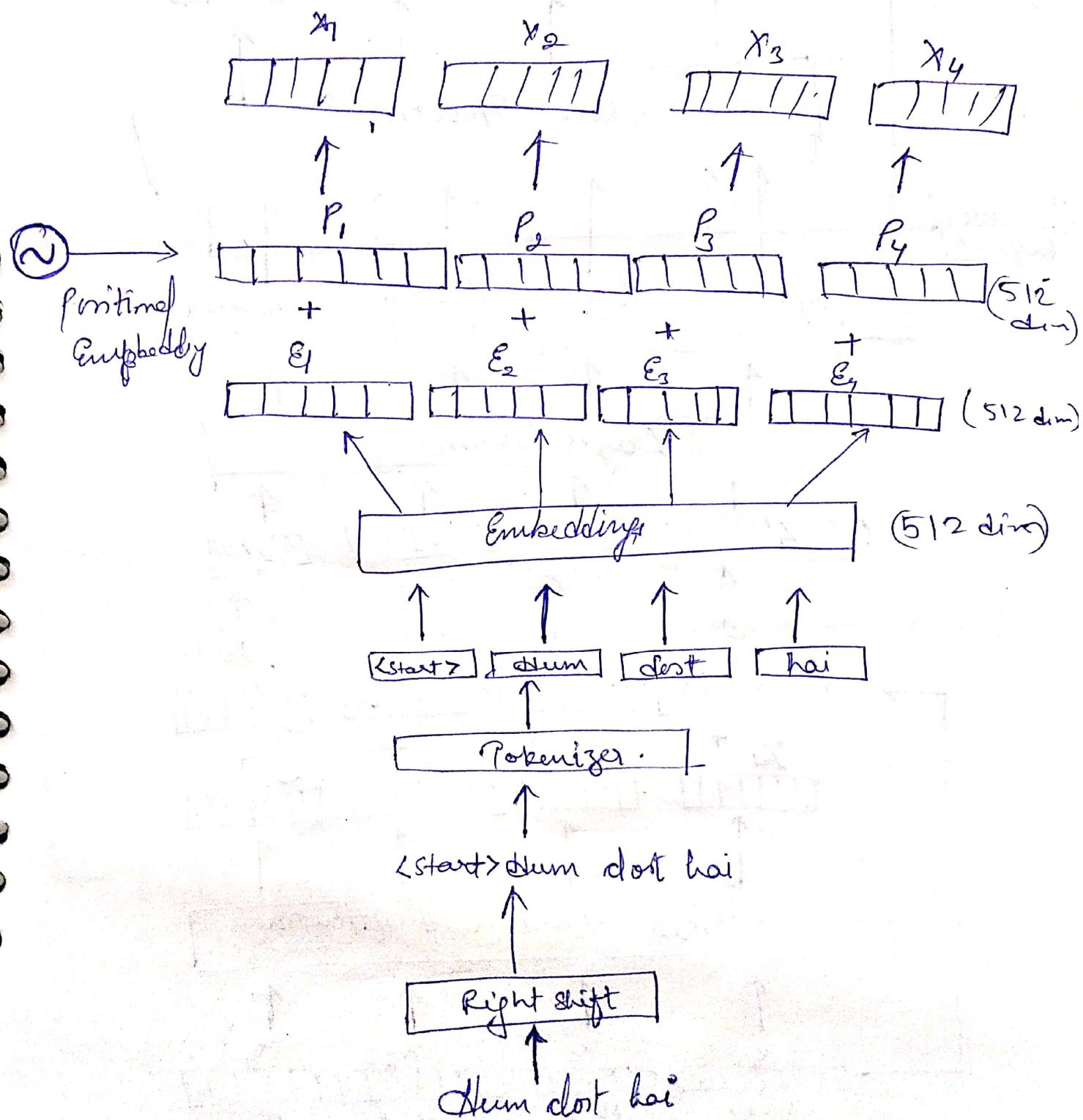
Input Block :-

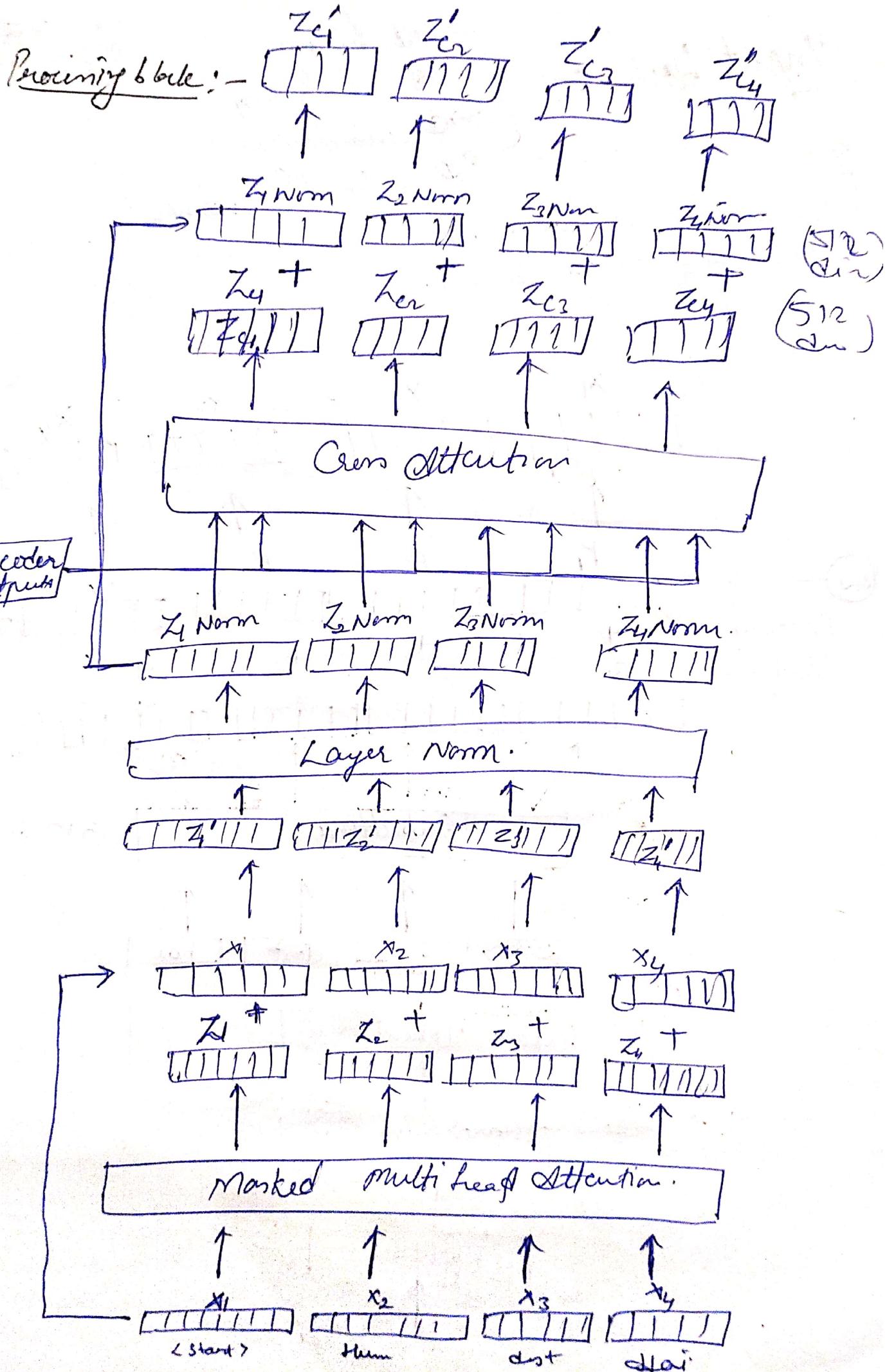
- ① Shifting
- ② Tokenization
- ③ Embedding
- ④ Position of Encoding

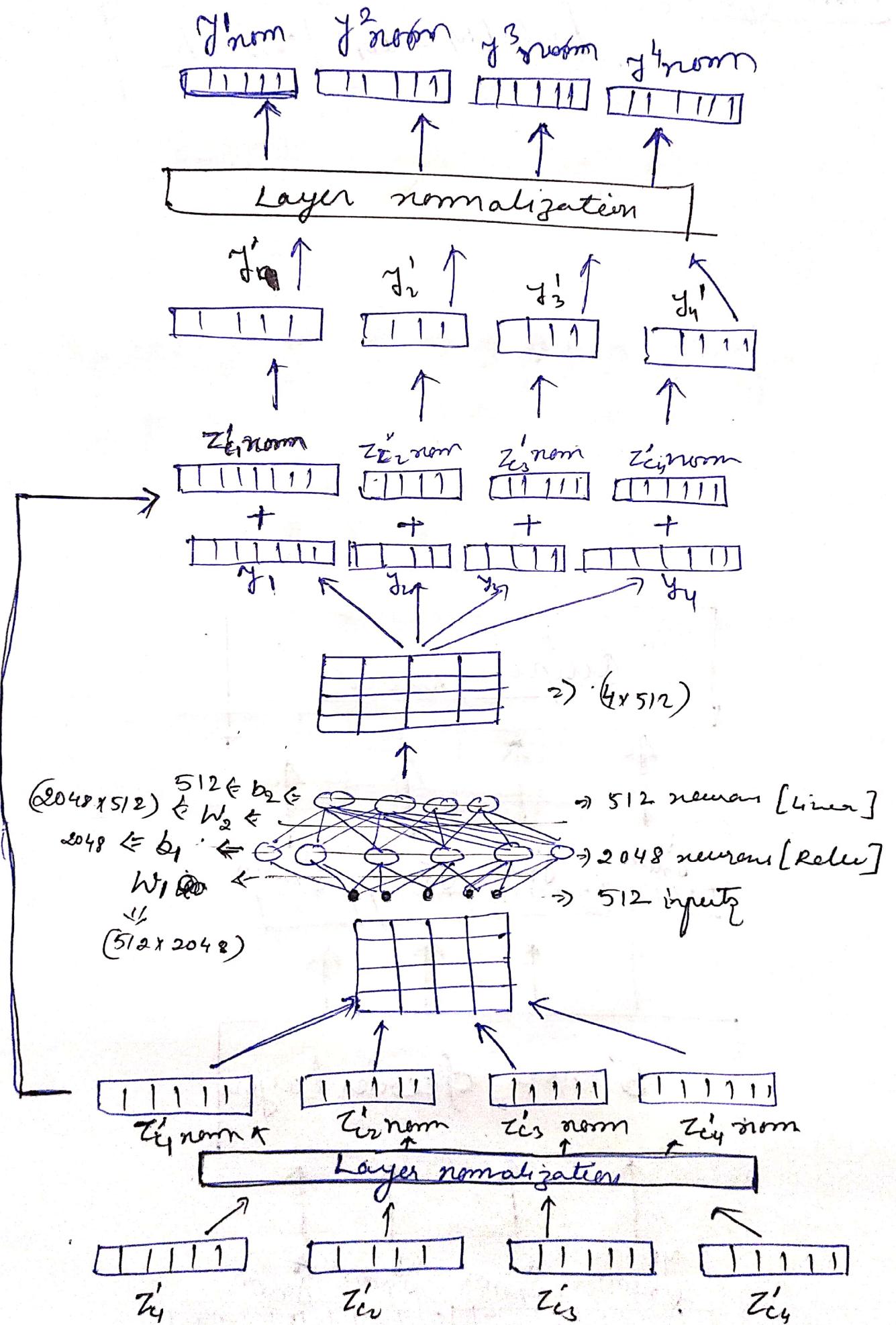
Eg:-

English sentence :- We are friends.

Hindi sentence :- शुभ दूष हाई।



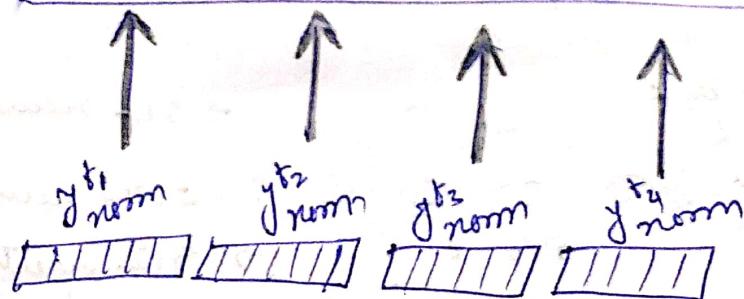




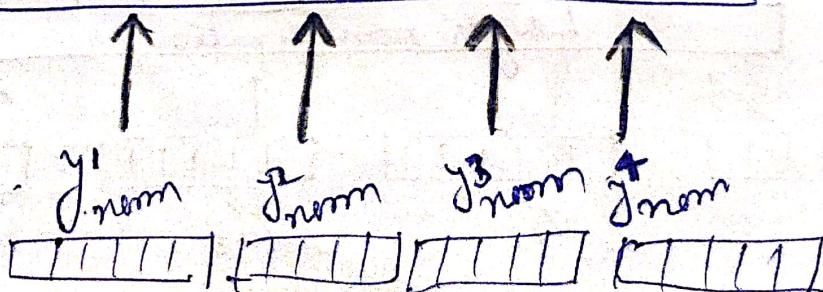
FFNN Equaⁿ⁻⁰

$$\boxed{[\text{relu}(w_1 \cdot z + b_1)] \cdot w_2 + b_2}$$

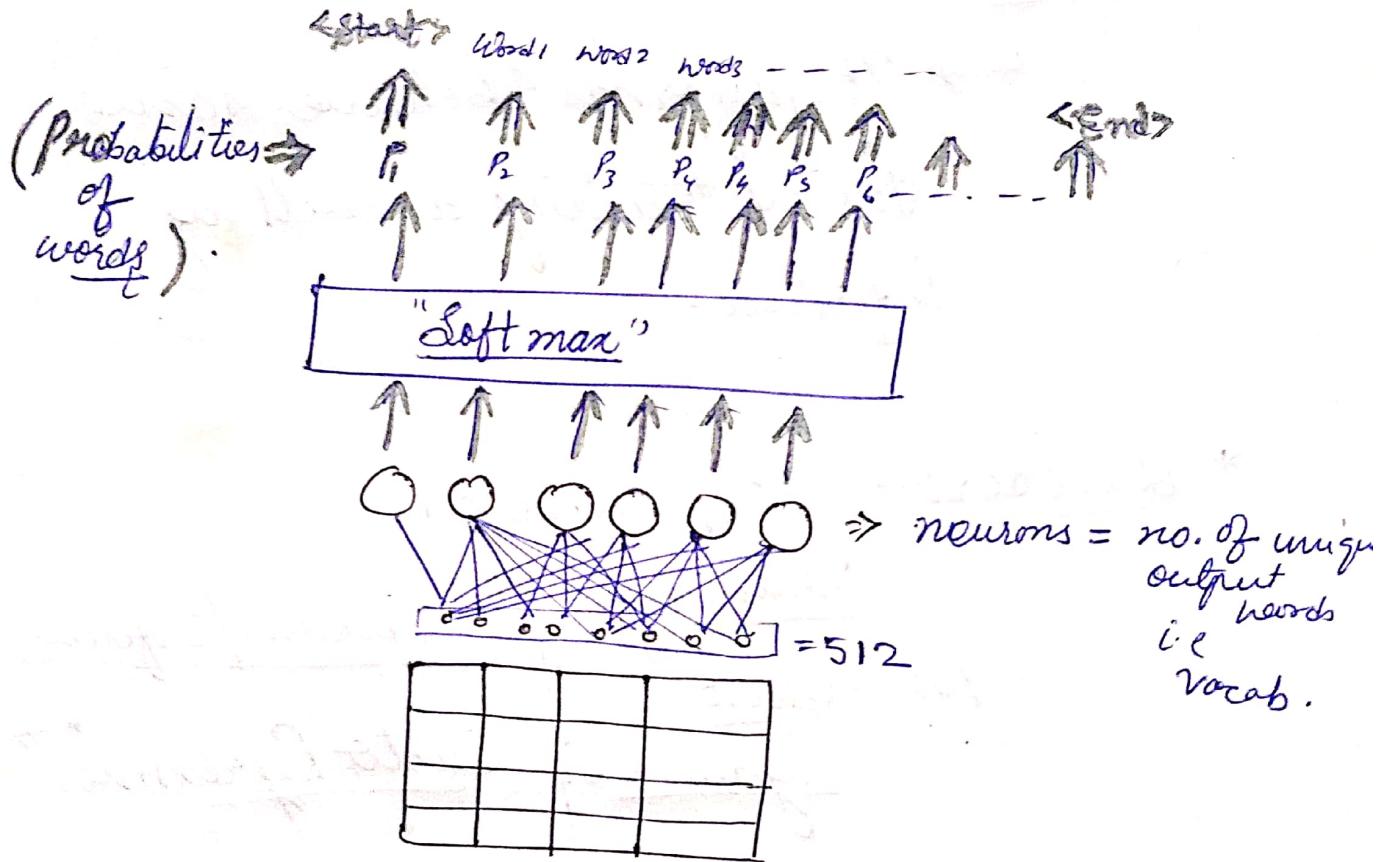
"Output Block"



5 more decoder layers



Output block :-



Transformer

During

Inference [Prediction]

* Encoder :-

→ It performs & behave same during training as well as inference.

* Decoder :- It behave as :-

<u>Training</u>	<u>Prediction / Inference</u>
" <u>Non-Auto Regressive</u> "	" <u>Auto Regressive</u> ."