

28/07/24 - Friday

Transformer Decoder Architecture

Output probabilities

Softmax

Linear

Add & Norm
Feed Forward

Add & Norm

Multi-head
attention

from encoder

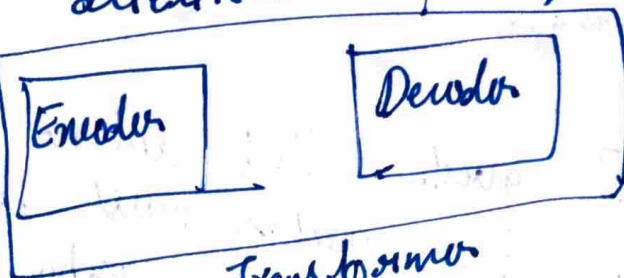
Add & Norm

Masked Multi-head
attention

Positional Encoding

Output Embedding

Self attention and cross
attention → processing



Encoder → encodes input
Decoder → decode the encoded
input

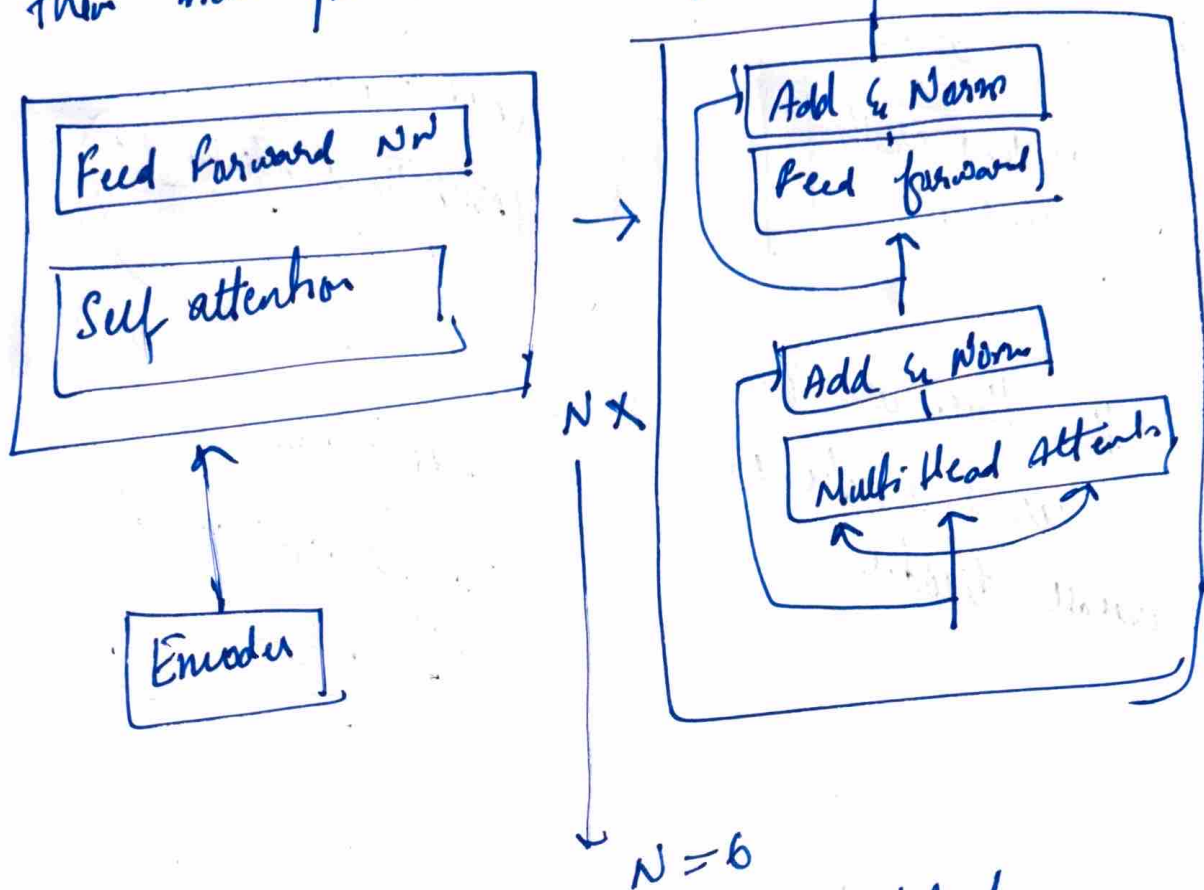


Encoder and Decoder are
made up of 6 blocks each



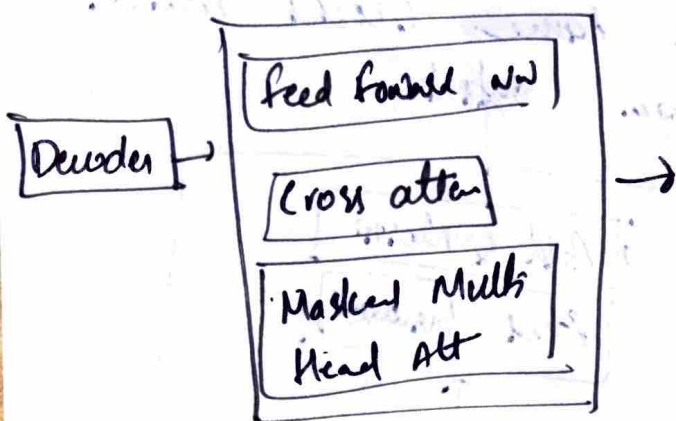
in each encoder, its
made up of 2 blocks

→ Each encoder block has same blocks inside them but parameters change. ^{output}

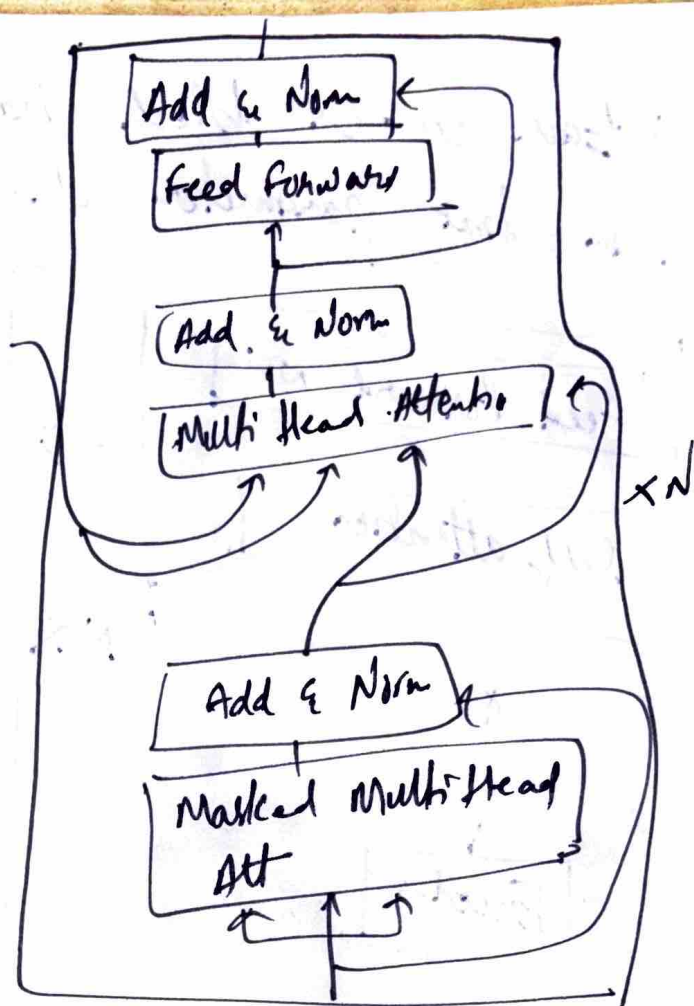


So there are 6 such blocks
The output from this Encoder will be given to decoder.



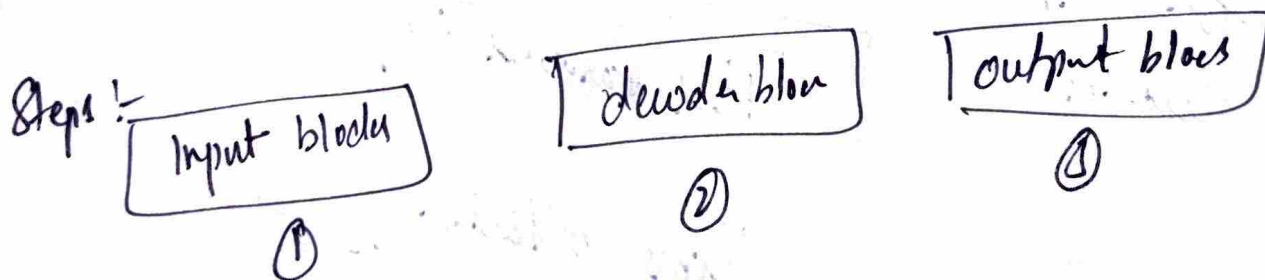


like this there are 6 decoder blocks in the overall decoder



Q) what happens in a single decoder?

A) As all the decoder blocks are same, so if you understand one decoder block, rest all are same.



Example: Machine Translation

English to Hindi

English = we are friends

Hindi = Hum dost Hai

English Sentence goes into Encoder

↓
generate contextual Embedding
for decoder to
work

① input block

there are 4 operations in this:-

1. Shifting

2. Tokenization

3. Embedding

4. Positional Embeddings.

→ the work of input block is to take output
sentence (Hindi sentence) and process it and gets
in the form to send it to decoder's first block.

Hum dost Hai

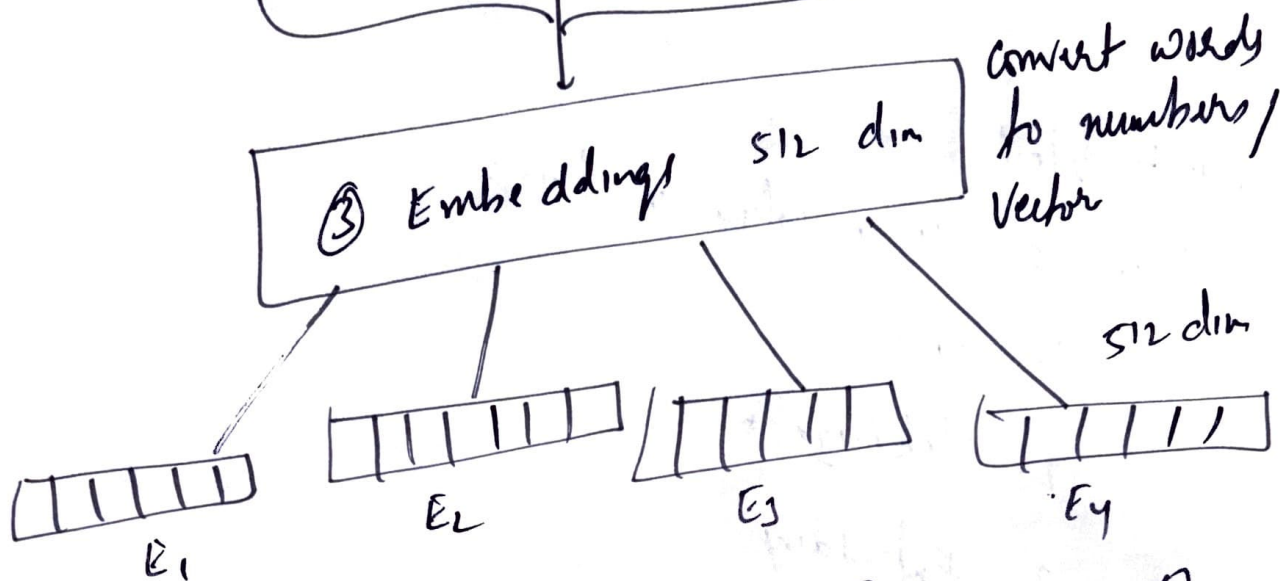
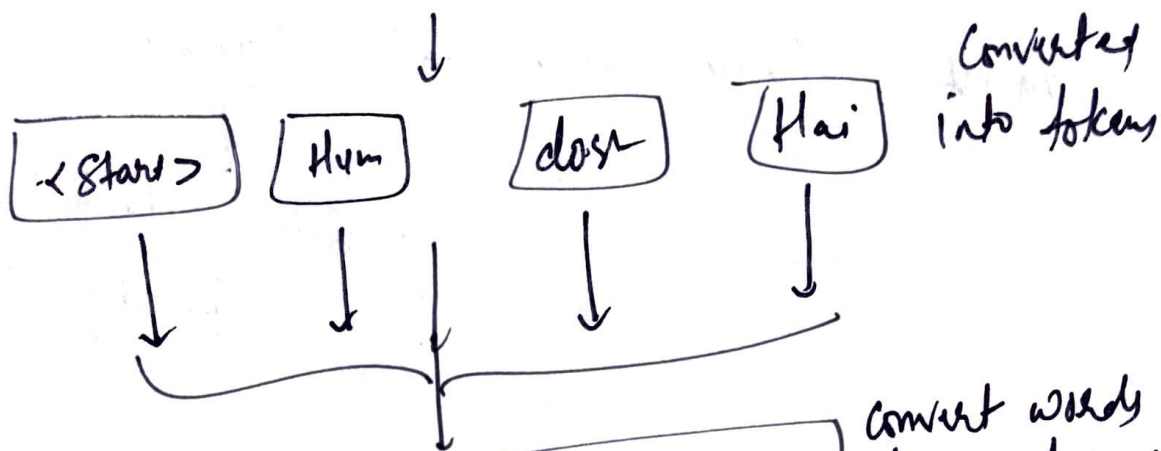
↓
① Right Shift

→ add a token start
before Hum

||
<start> Hum dost Hai

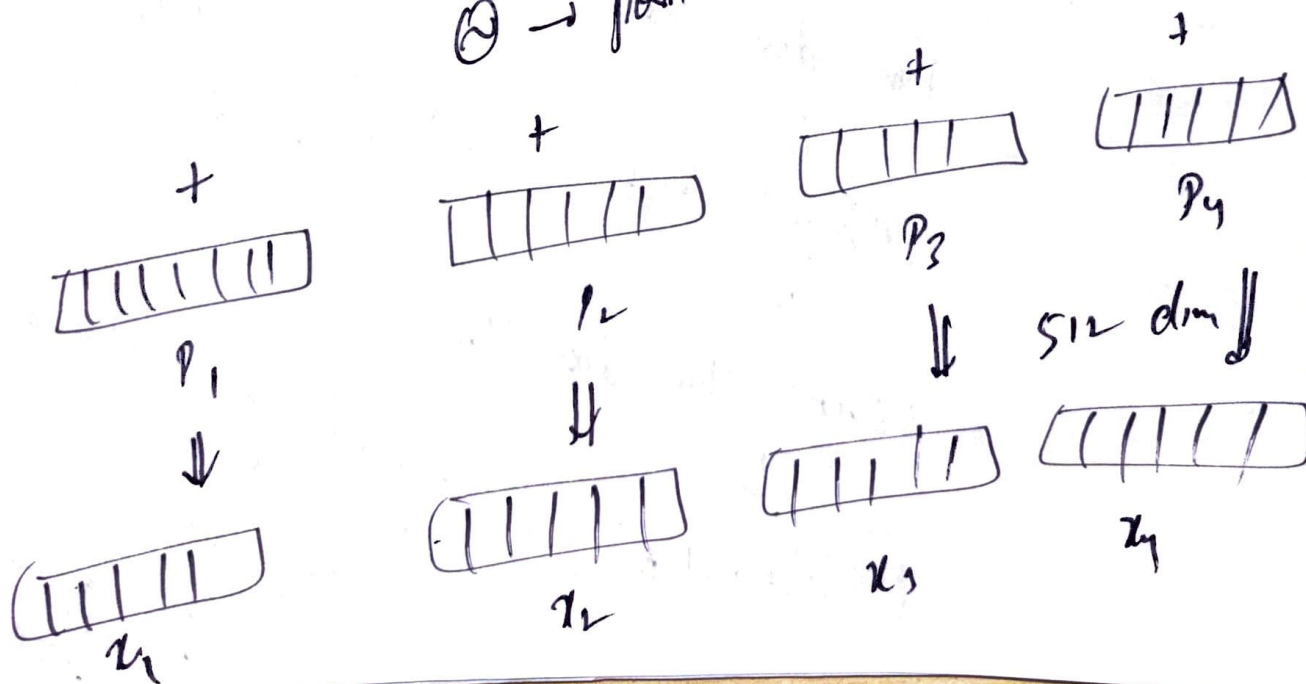
↓
② Tokenization

→ Break down sentences
into tokens
any gram
eg: 1 gram
token



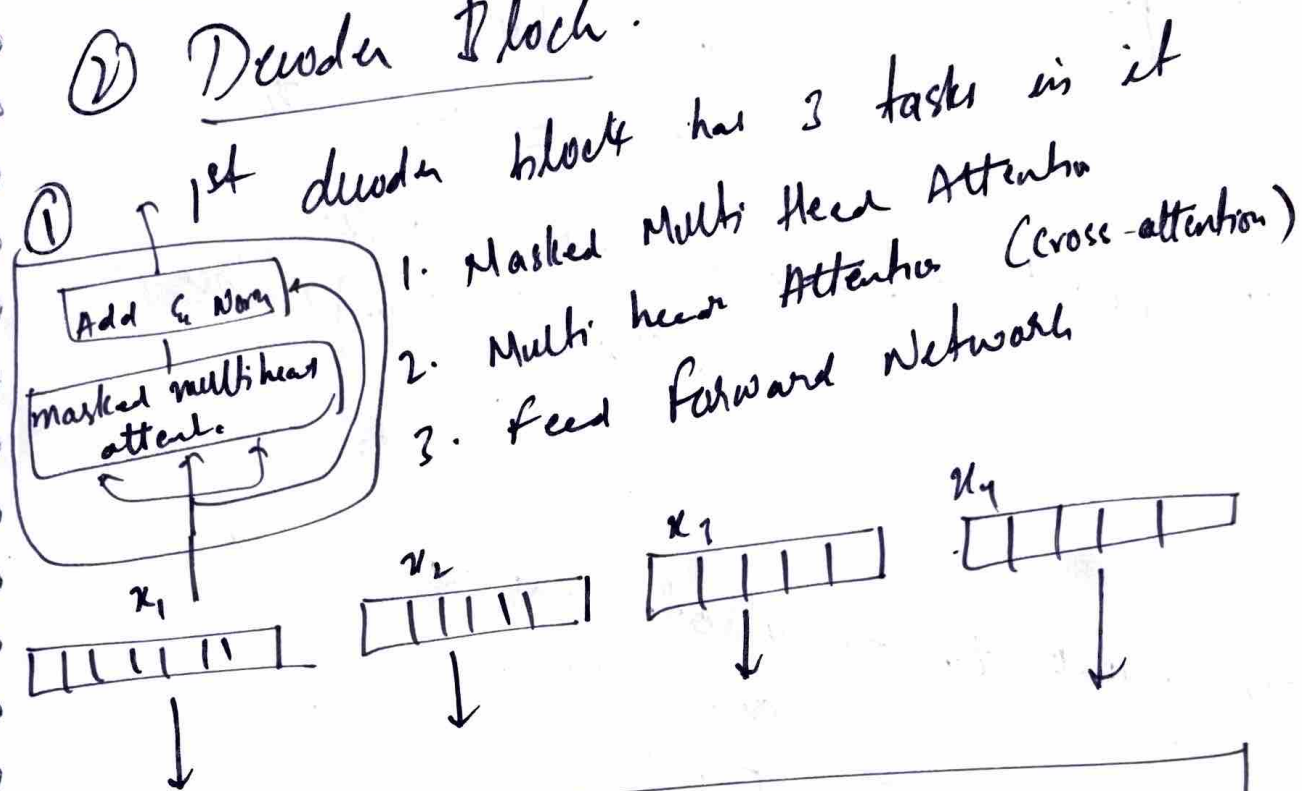
here we don't know when word is before like, the position of the words so I tackle that issue we see this

② → positional Encodings.

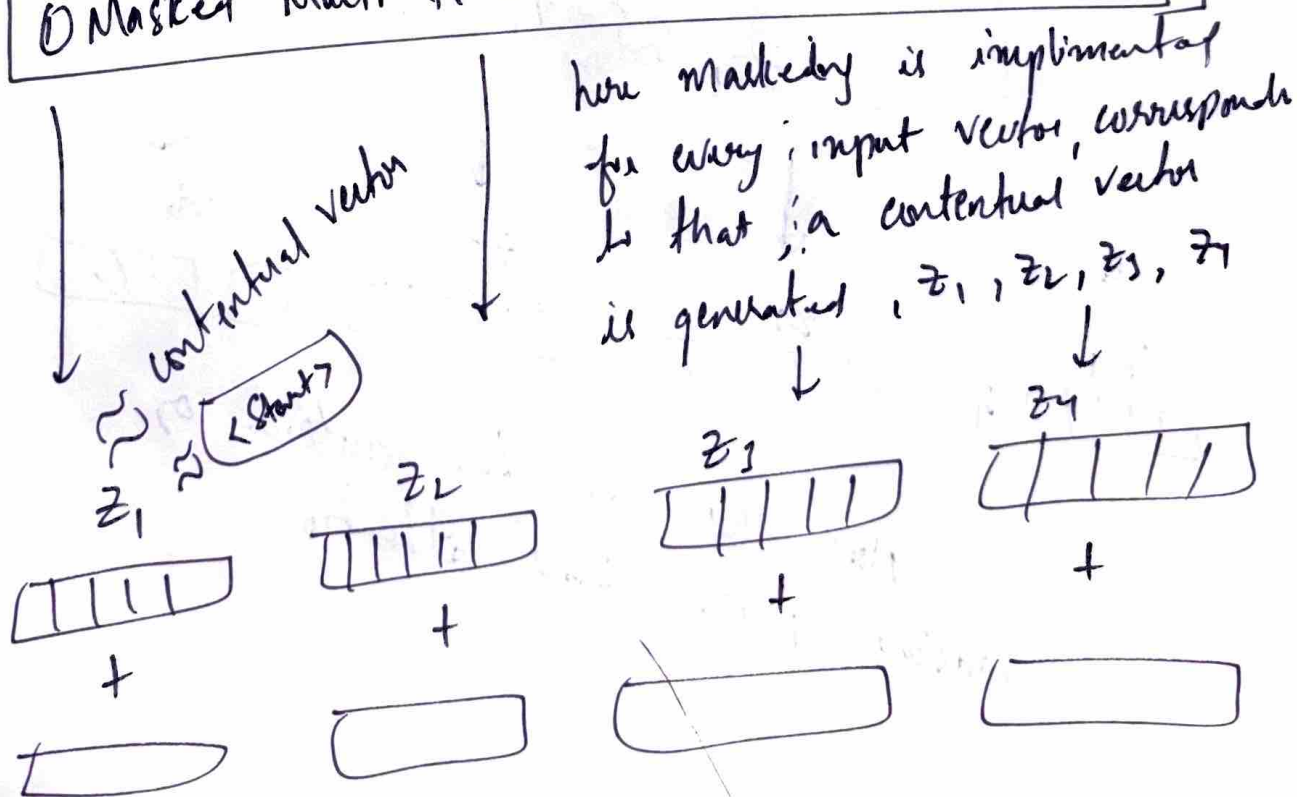


Now we are done with input block.
 the final x_1, x_2, x_3, x_4 vectors are
 the inputs for decoder

② Decoder Block.



③ Masked Multi Head Attention



To create this masked contextual vectors, we only consider z_1 vector, i.e., <start> and we do not consider rest. $z_1 \approx x_1$

To generate $z_2 \rightarrow$ we only consider <start> and Hum

$$z_2 \approx x_1 \quad x_2$$

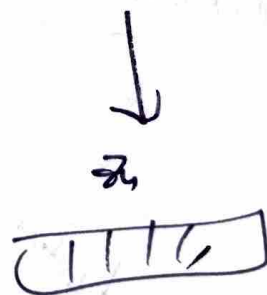
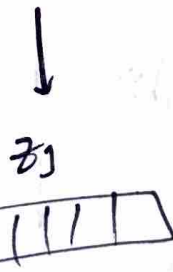
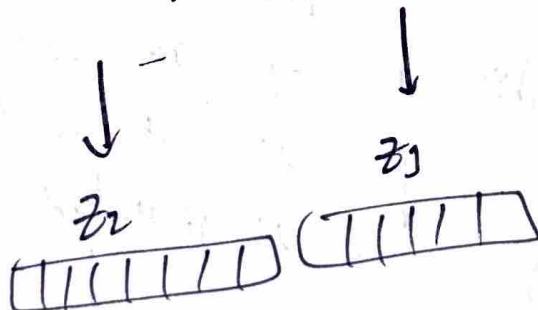
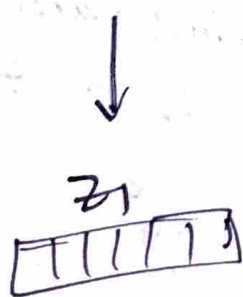
To generate $z_3 \rightarrow$ we only consider <start> Hum close

$$z_3 \approx x_1 \quad x_2 \quad x_3$$

To generate $z_4 \rightarrow$ we only consider previous and current

Start Hum close Hai

$$z_4 \approx x_1 \quad x_2 \quad x_3 \quad x_4$$



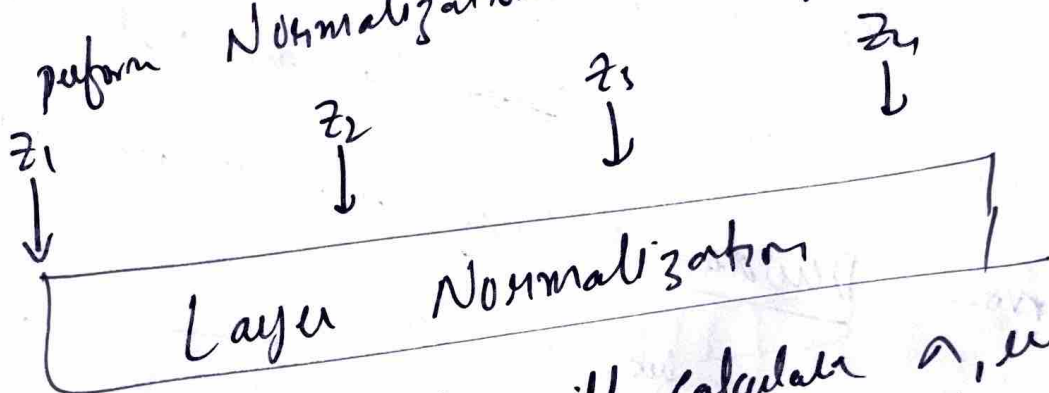
Now these are the output of masked Multi Head attention

Now perform Add & Norm

$$\text{Add} = \boxed{\text{Outputs of masked multi head attn}} + \boxed{\text{original inputs to the decoder}}$$

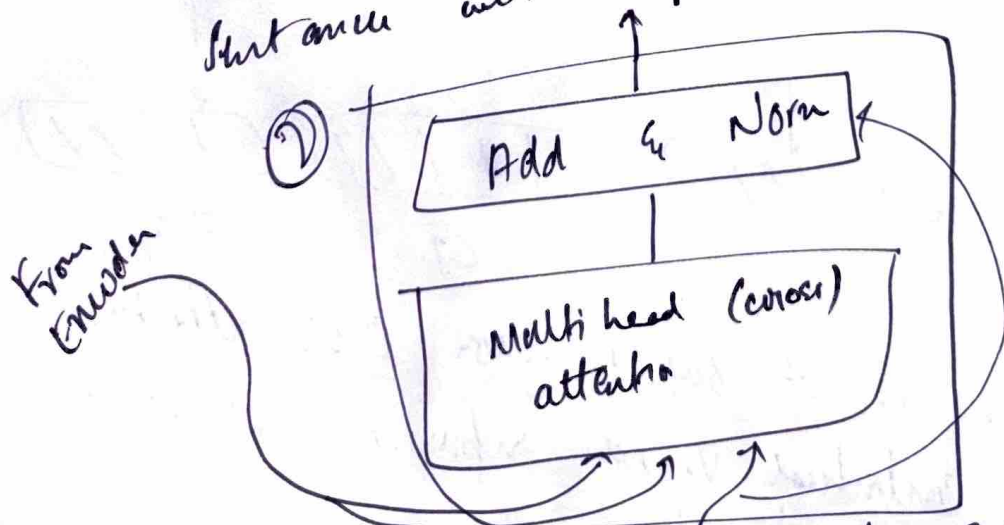
$z_1 \ z_2 \ z_3 \ z_4$ $x_1 \ x_2 \ x_3 \ x_4$

Now perform Normalization \approx layer Norm



for each vector, this will calculate μ, σ and with that help it normalize, this is done to make training process stable

Now cross attention comes into field here, the connection between Hindi and English sentences will happen.



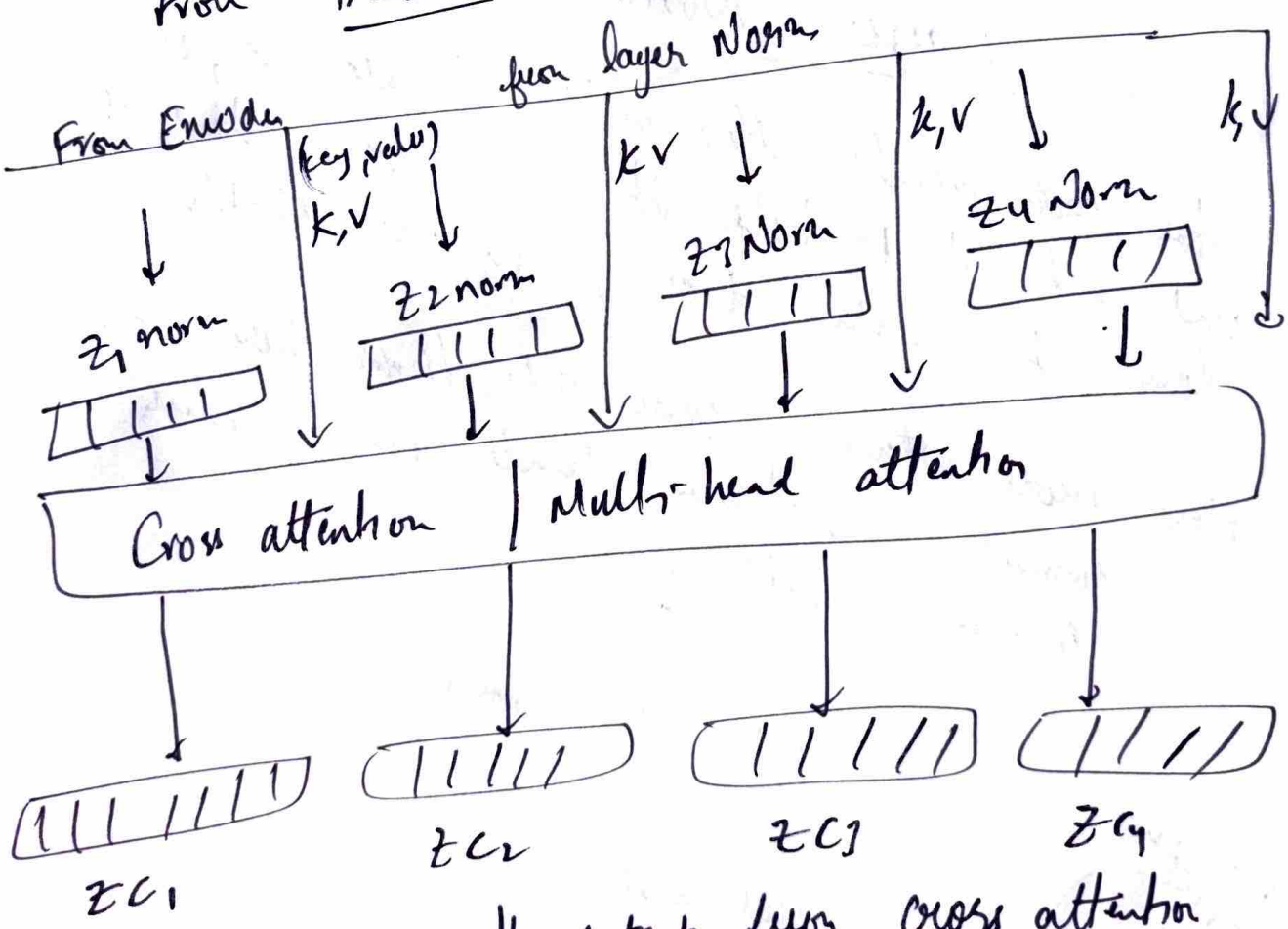
$z_{1Norm}, z_{2Norm}, z_{3Norm}, z_{4Norm}$ inputs from part 1 of decoder

Cross attention

1 Seq (English) from encoder
 2 Seq (Hindi) from previous step of decoder

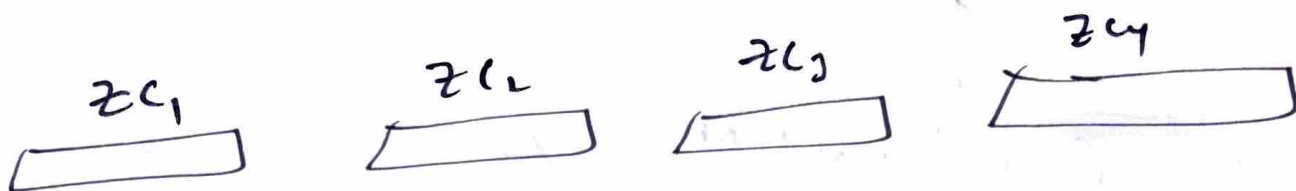
→ this is done as we need 3 vectors
 key, value and Query

From Encoder → key and value is given
 from Decoder → Query vector is given and



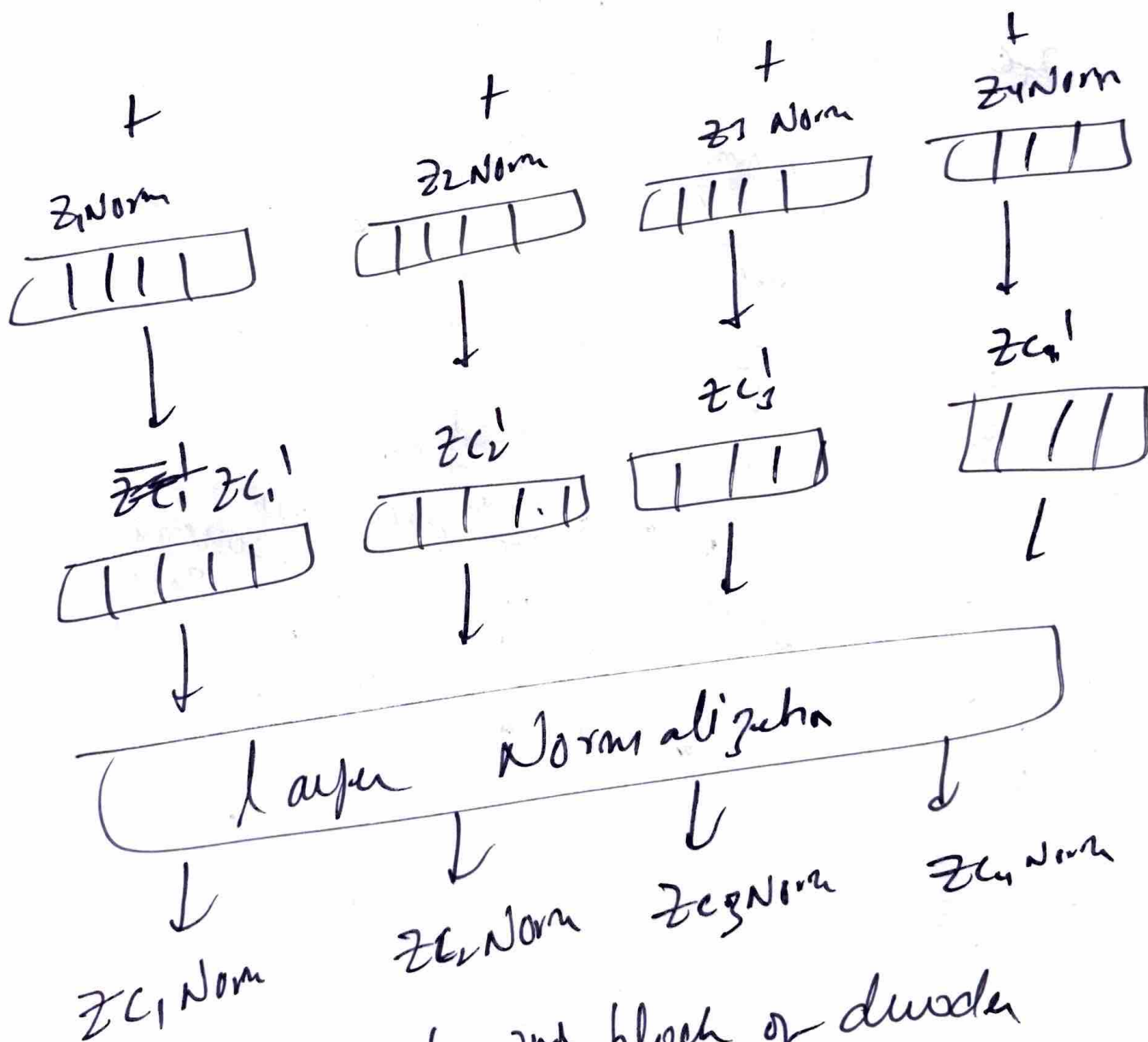
These are the outputs from cross attention
 where the contextual vector information is attached

for each token of the output vector (z_1, z_2, z_3, z_4) there will be a contextual embedding vector



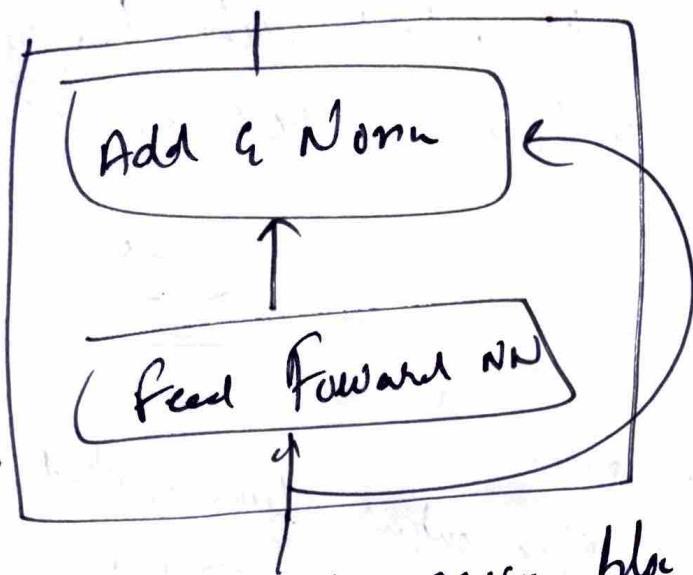
Add and Norm

Previous step output from masked multi-head att (z_1^{Norm} , z_2^{Norm} , z_3^{Norm} , z_4^{Norm})



Done with 2nd block of decoder

②



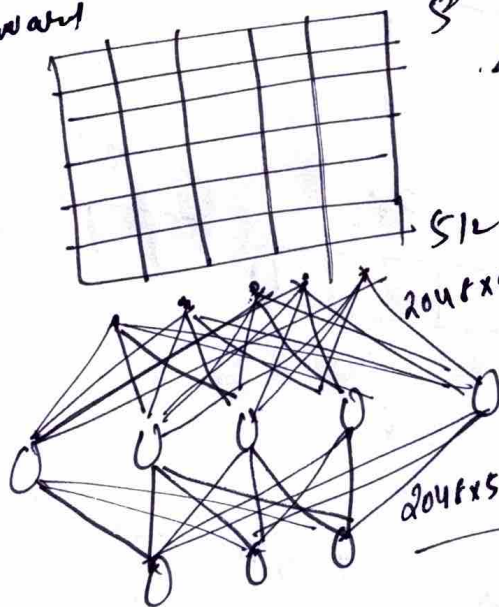
input from previous block

$z_{C1} Norm, z_{C2} Norm, z_{C3} Norm, z_{C4} Norm$

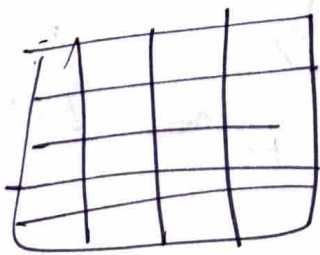
Diagram



Same architecture as
encoder
has 2 layers, 2048 neurons



ω_1 2048 neurons
ReLU
 β 2048 bias
 ω_2 512 neurons
Linear Activation
 β 512 bias



output
512

$$(w_1 \cdot z + b_1)$$

$$[4 \times 512 \quad 512 \times 2048 \approx 4 \times 2048] \text{ dimension check}$$

$$\text{ReLU}(w_1 \cdot z + b_1)$$

$$\text{Shape} = 4 \times 2048$$

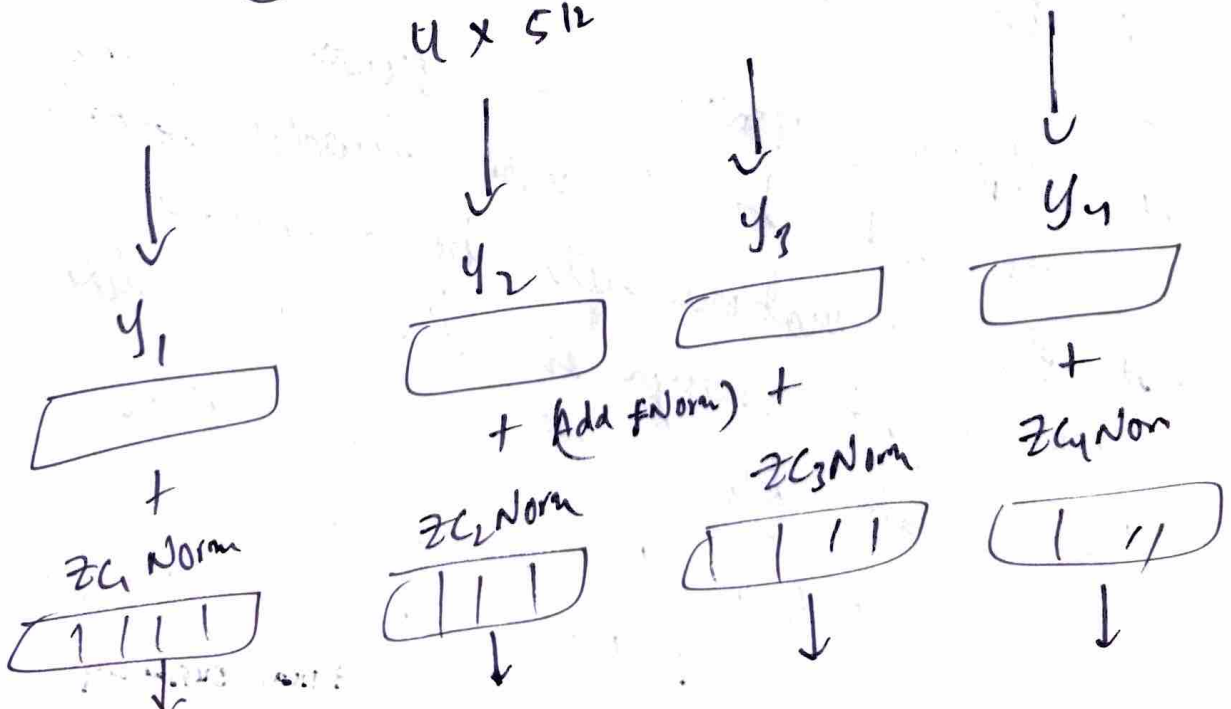
in next layer

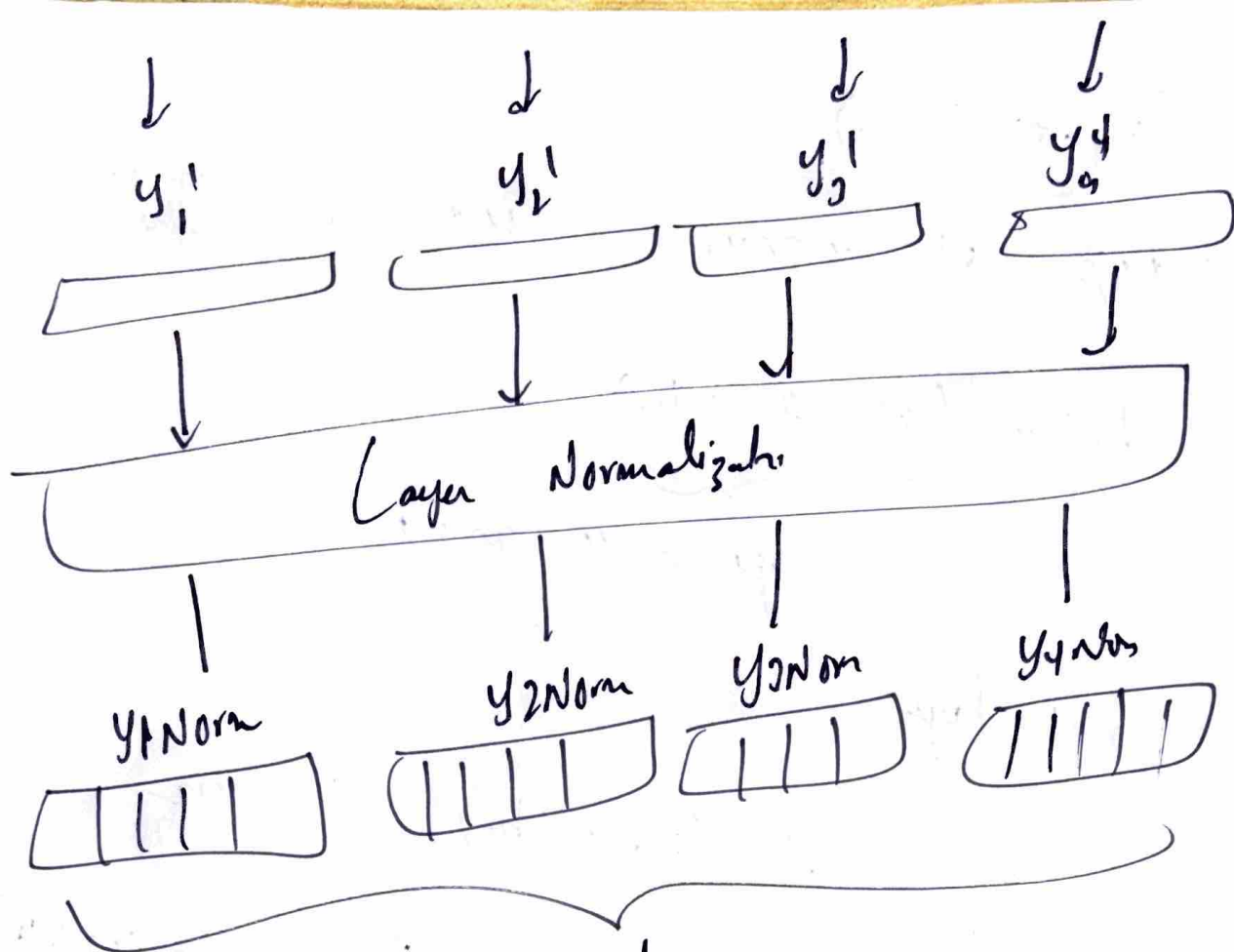
$$w_2 [\text{ReLU}(w_1 \cdot z + b_1)] + b_2$$

dimension check $(4 \times 2048) \cdot 2048 \cdot 512 = 4 \times 512$

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

$$4 \times 512$$

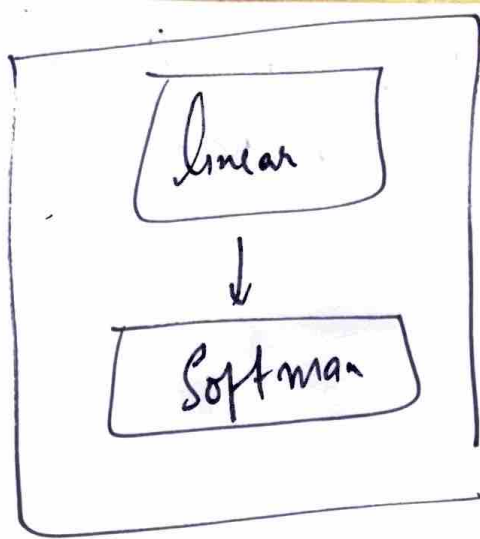




these 4 vectors are the ~~first~~ output of decoder 1

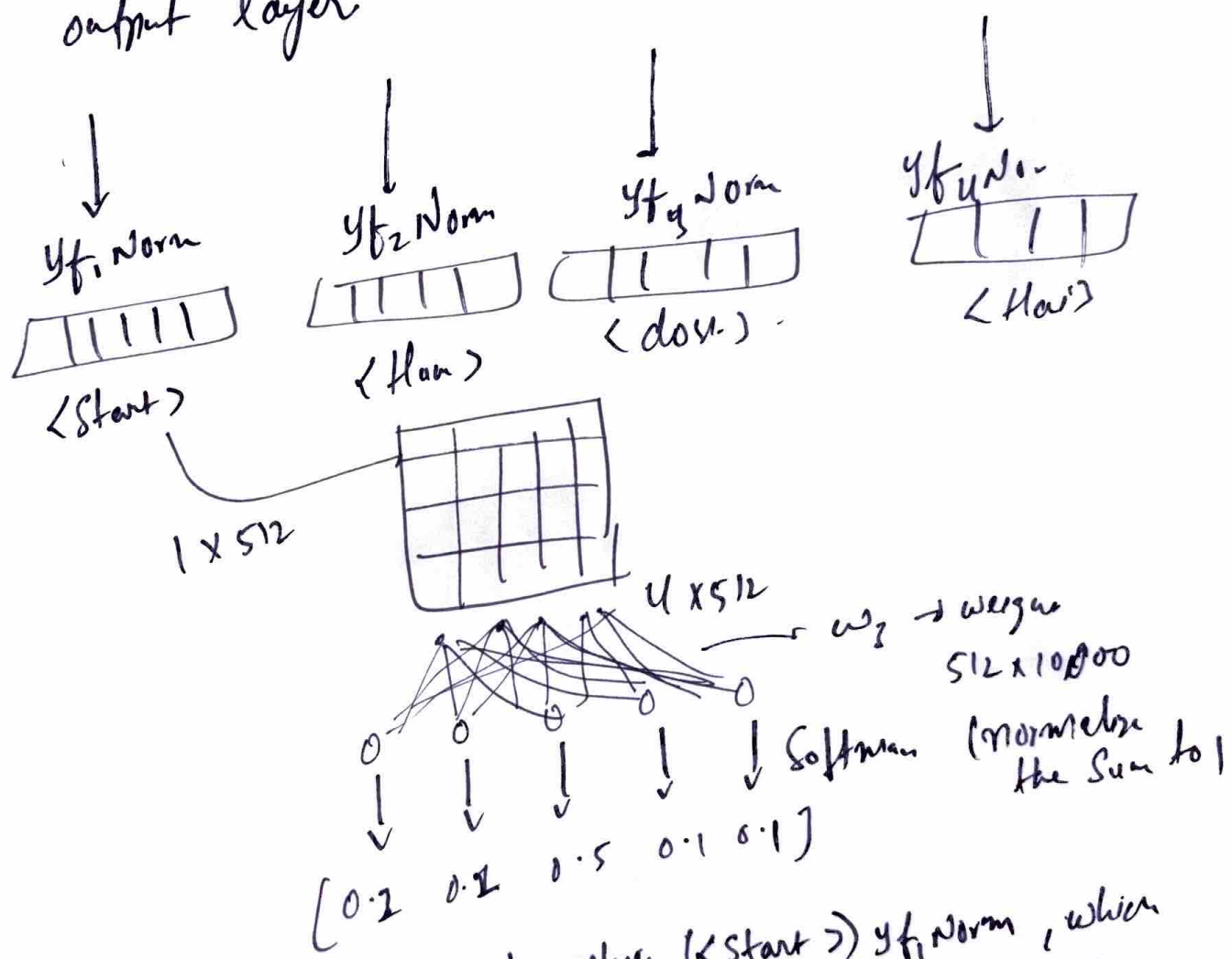
③ Output block

These y_1^{Norm} , y_2^{Norm} ... y_n^{Norm} layers are sent to 5 more decoder blocks and some operations will be executed with slight change in the parameters the outputs from 6th decoder layer the outputs are called as y_{f1}^{Norm} , y_{f2}^{Norm} ... y_{fn}^{Norm} final output



Output block

made of 2 layers, linear and softmax.
Can be compared to feed forward NN
output layer



Now is corresponding to where (<Start>) y_{f1_norm} , which word has high probability from the neural network
now send (y_{f2_norm}) to the neural network and see which word

has high
translation
reachability, that is the
of that english word.
kind: