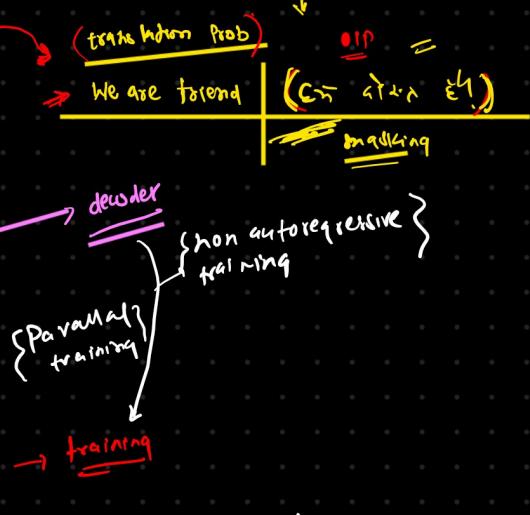
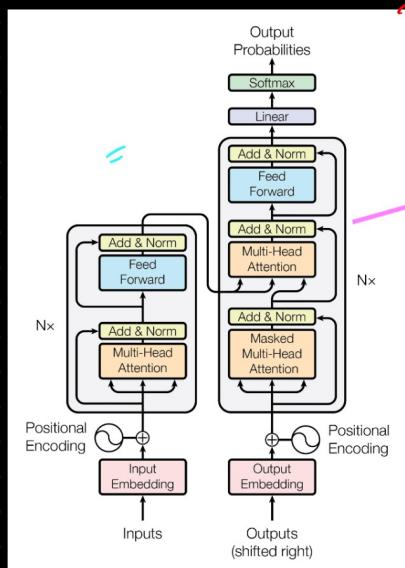


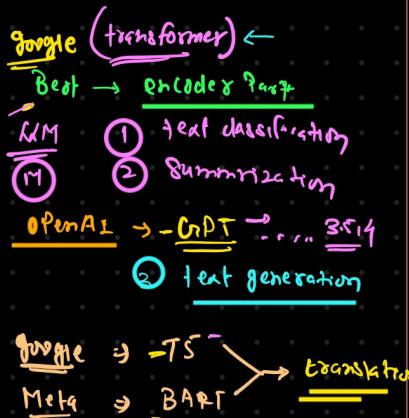
- 1 Encoder → Self attention → multihead self attention
- 2 Entire flow of encoder
- 3 Positional encoding
- 4 Residual connection
- 5 Importance of NW in encoder
- 6 multi head attention } decoder
- 7 cross attention } ==
- 8 Entire flow of decoder → training → }
- 9 Entire flow of inference → Prediction ← }

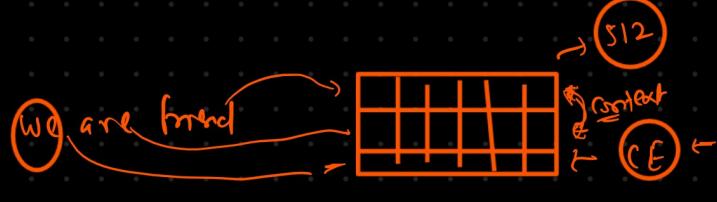
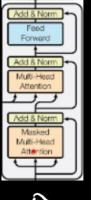
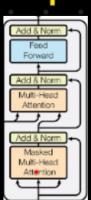
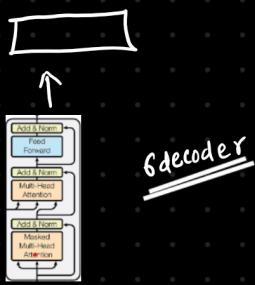
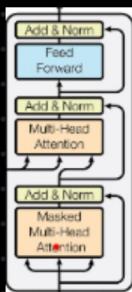
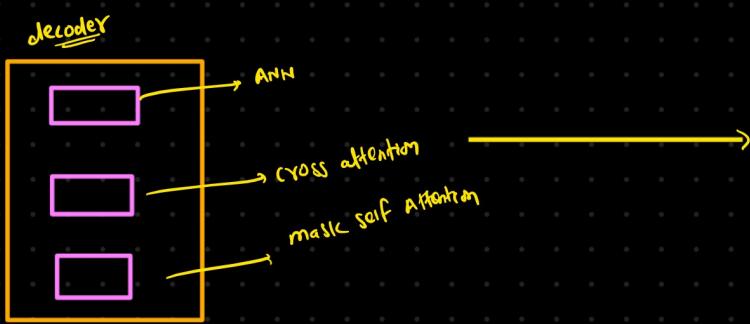
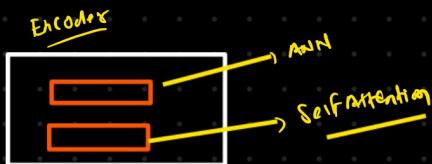
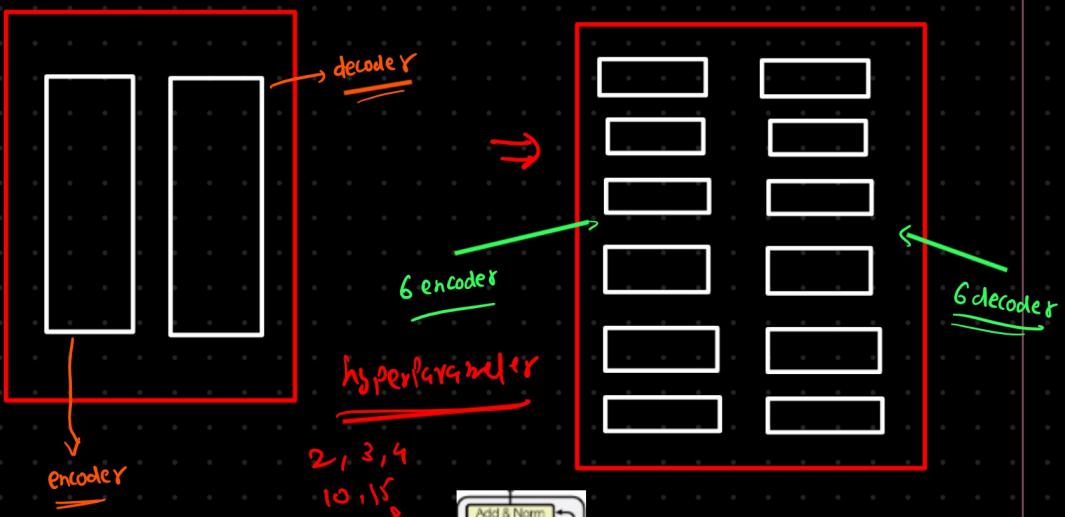
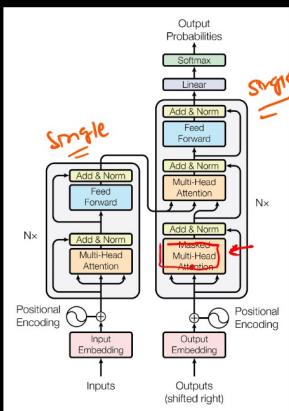


Inference (Prediction)
↳ Autoregressive (sequentially)
↳ sequentially

en. s1 s2 s3
My name is
Sherry

→ transformer → ? → translate





We are friend | Ein g'schickter Teil.
 ||

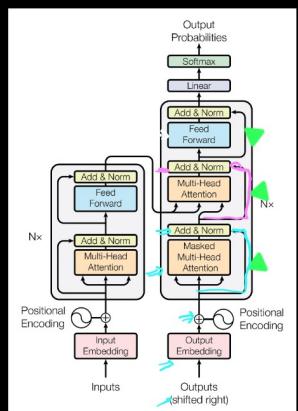
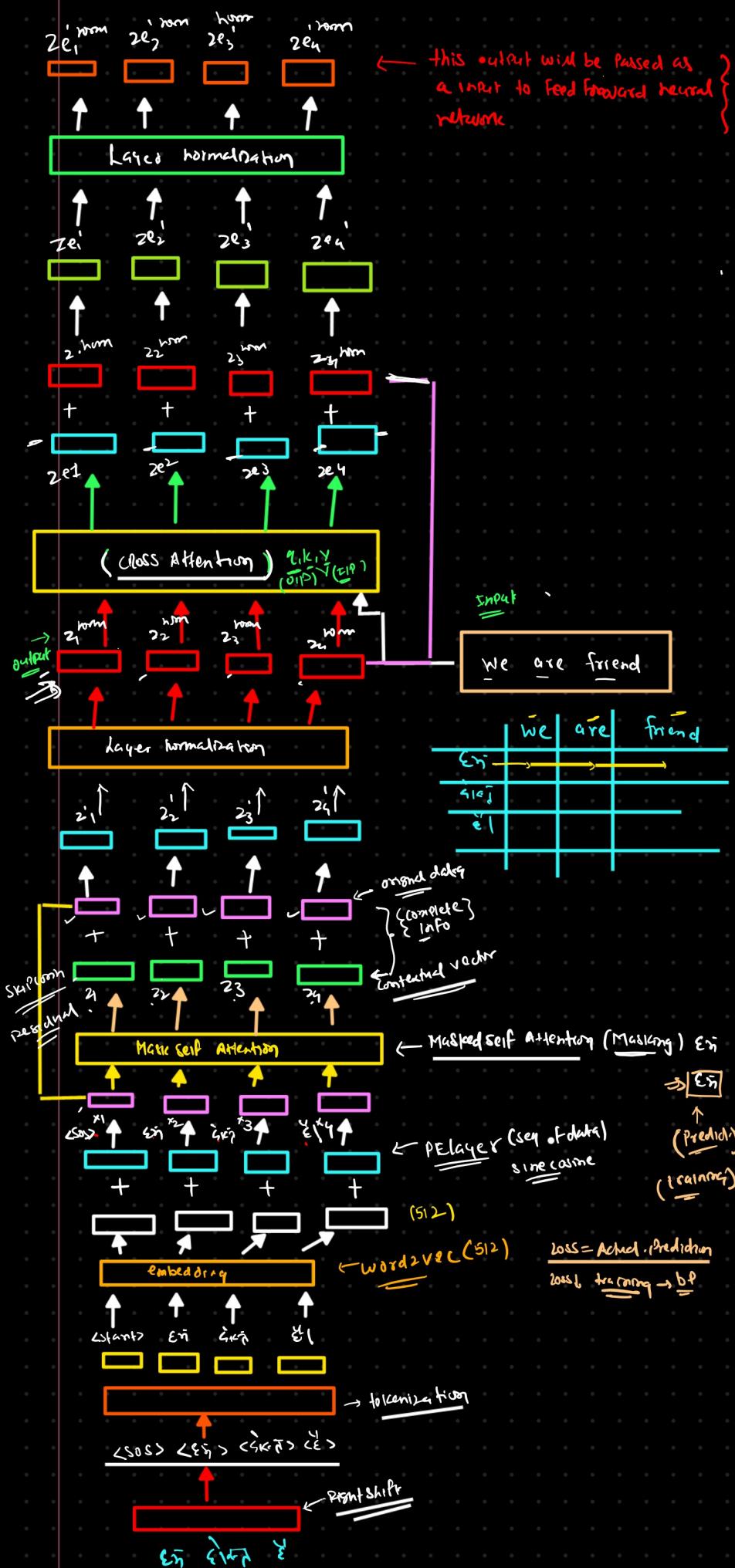


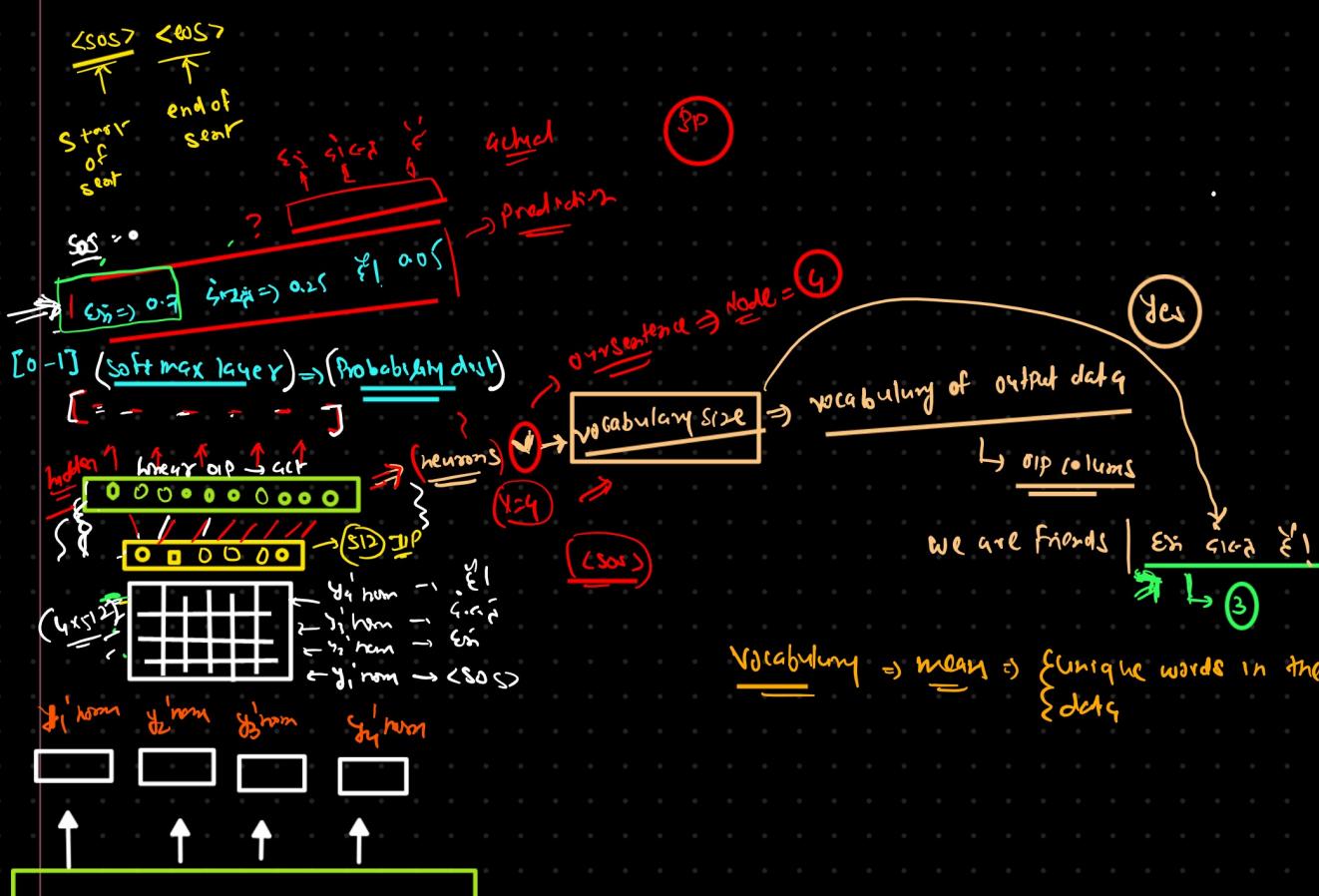
Diagram illustrating the forward pass of a neural network layer:

$$y\hat{=}W_2 \cdot h_1 + b_2$$

$$h_1\hat{=}W_1 \cdot x + b_1$$

The input x is multiplied by W_1 to produce intermediate hidden states h_1 . These states are then multiplied by W_2 to produce the final output $y\hat{=}$.

Pred.
= we are enemy => ?
↳ generate
($\frac{\text{key}}{\text{key}}$)



Vocabulary \Rightarrow mean $\Leftrightarrow \left\{ \begin{array}{l} \text{unique words in the given} \\ \text{data} \end{array} \right\} >$

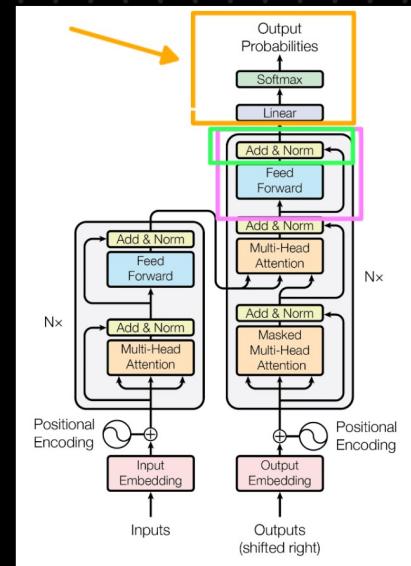
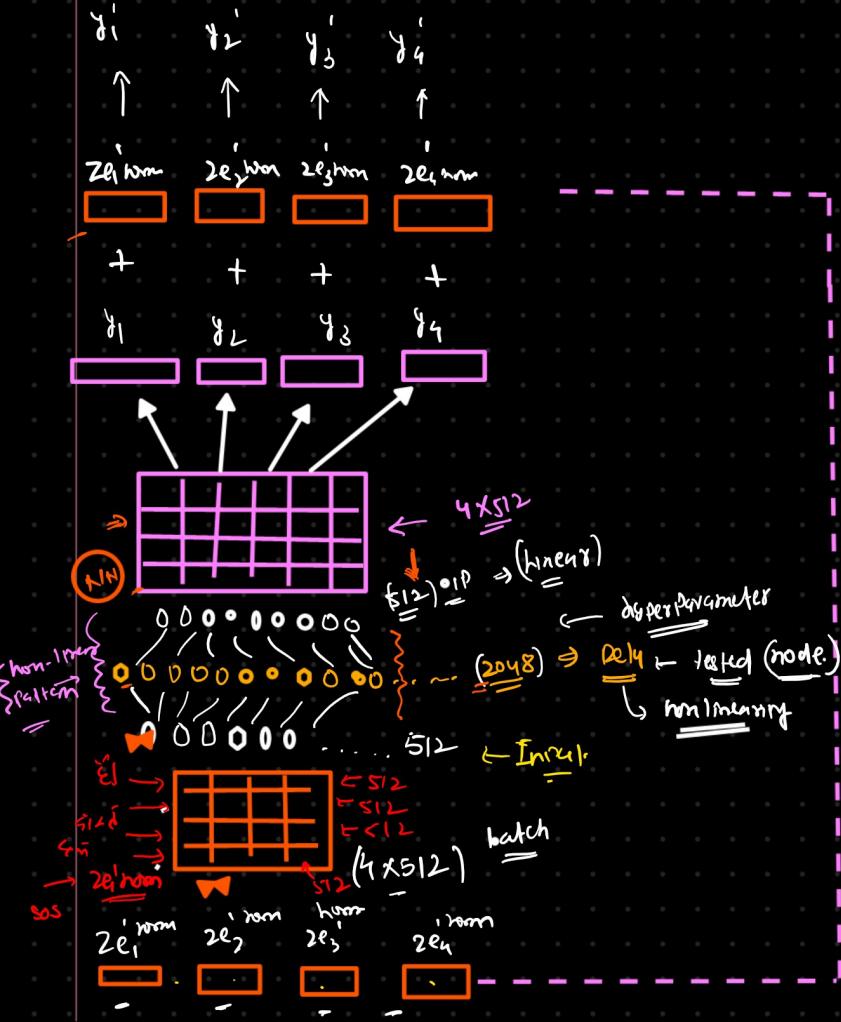
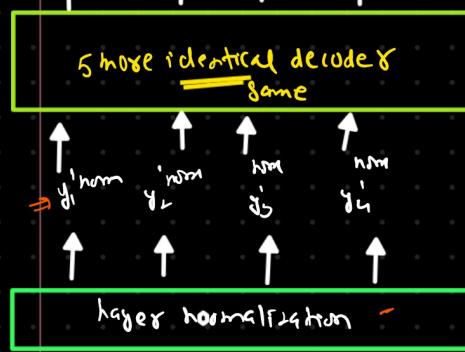
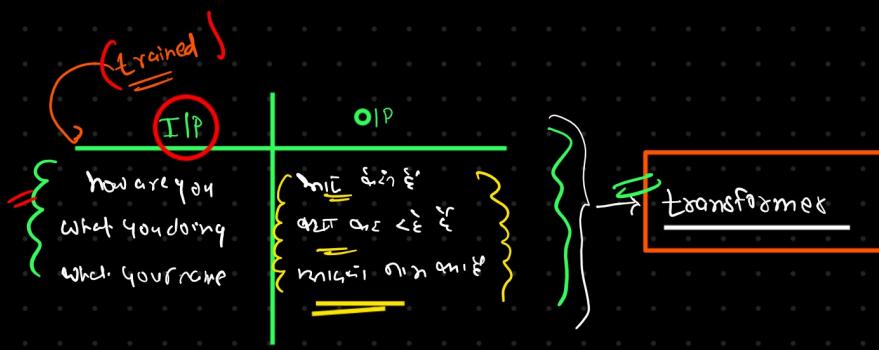
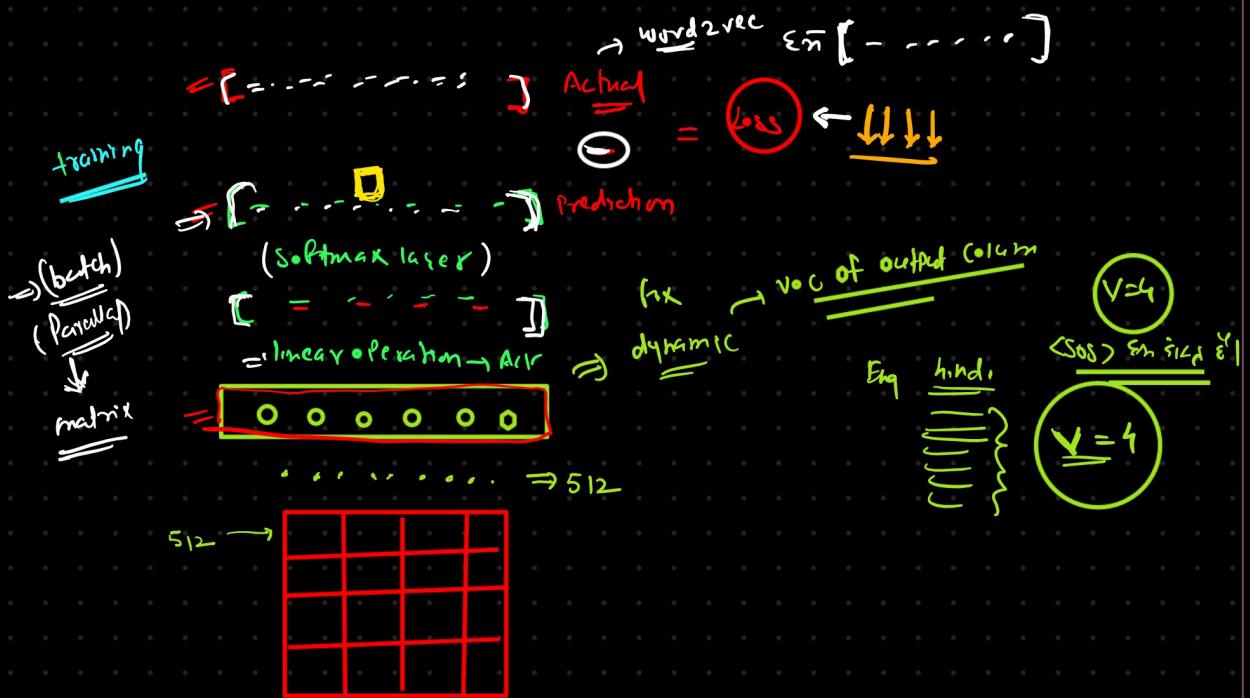
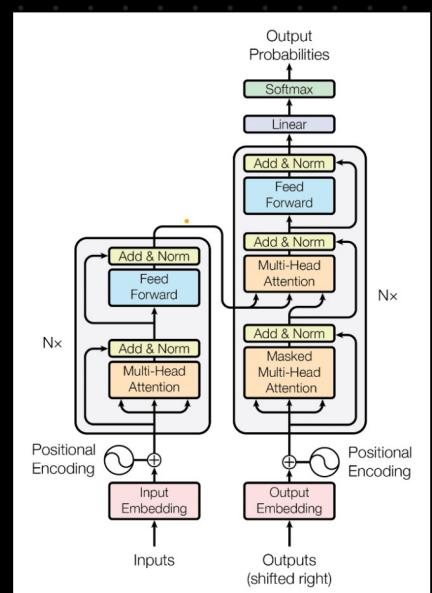
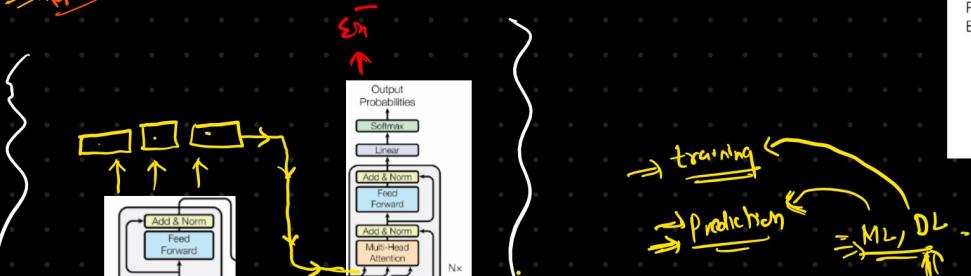


Diagram illustrating a fully connected layer (FC layer) in an Artificial Neural Network (ANN). The input layer (IP) has 4 units and 512 dimensions. It is multiplied by a weight matrix (matrix multiplication) to produce a hidden state (Hidden) with 4 units and 2048 dimensions. The resulting output is labeled as \hat{y} .

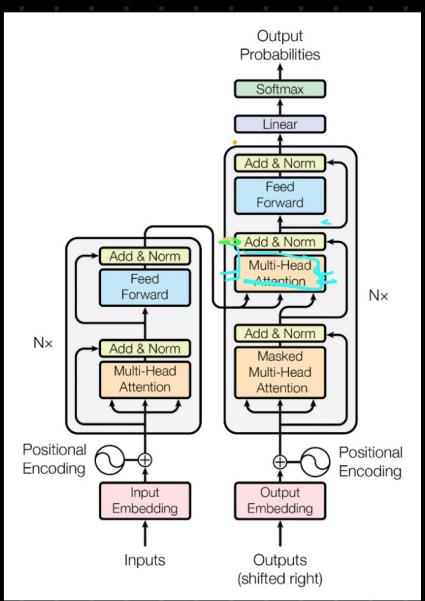
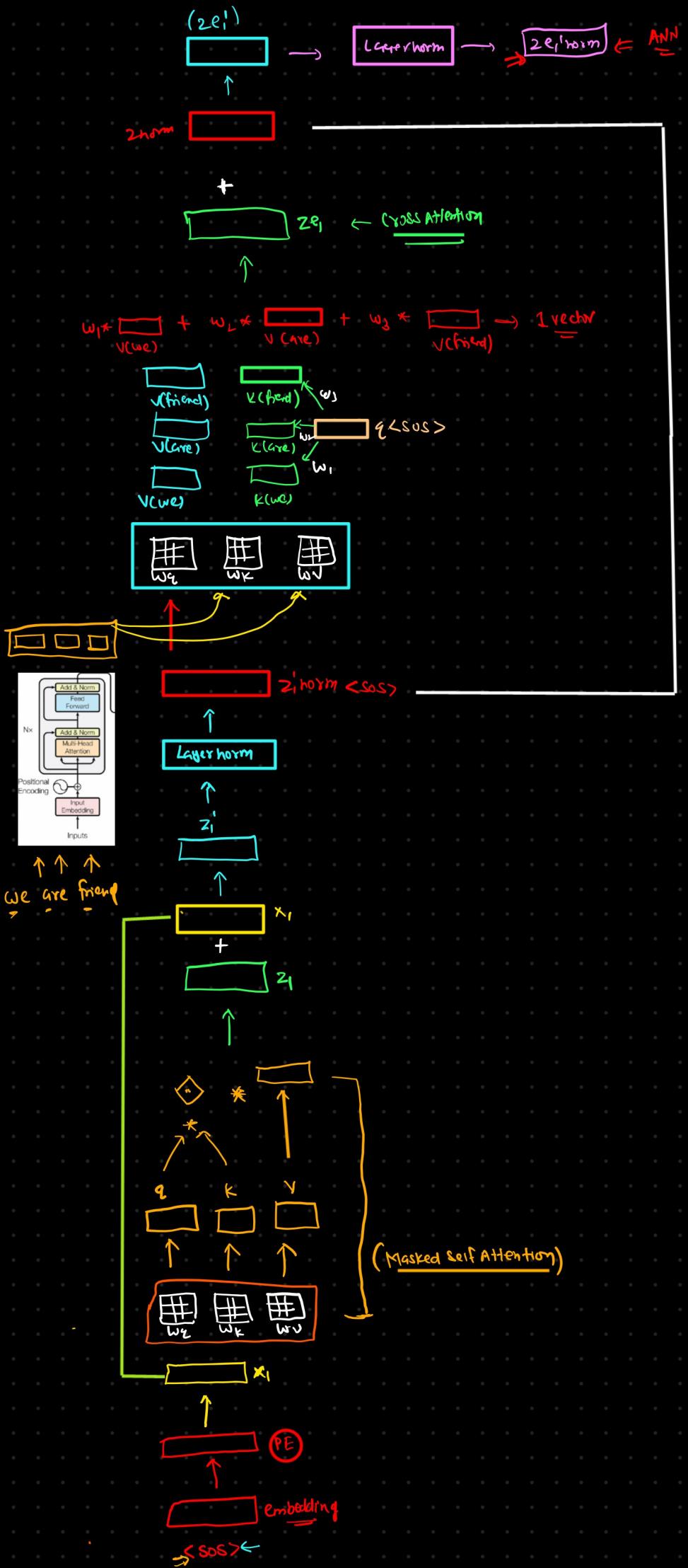


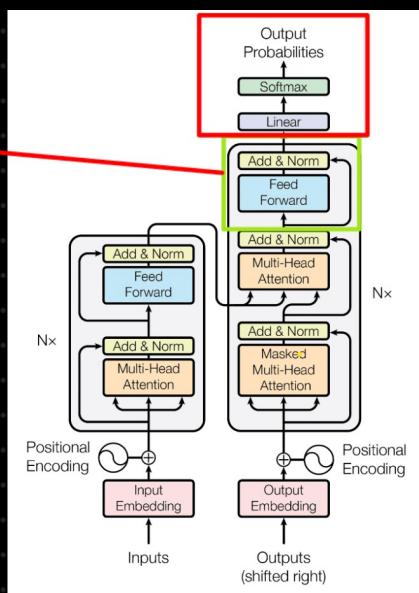
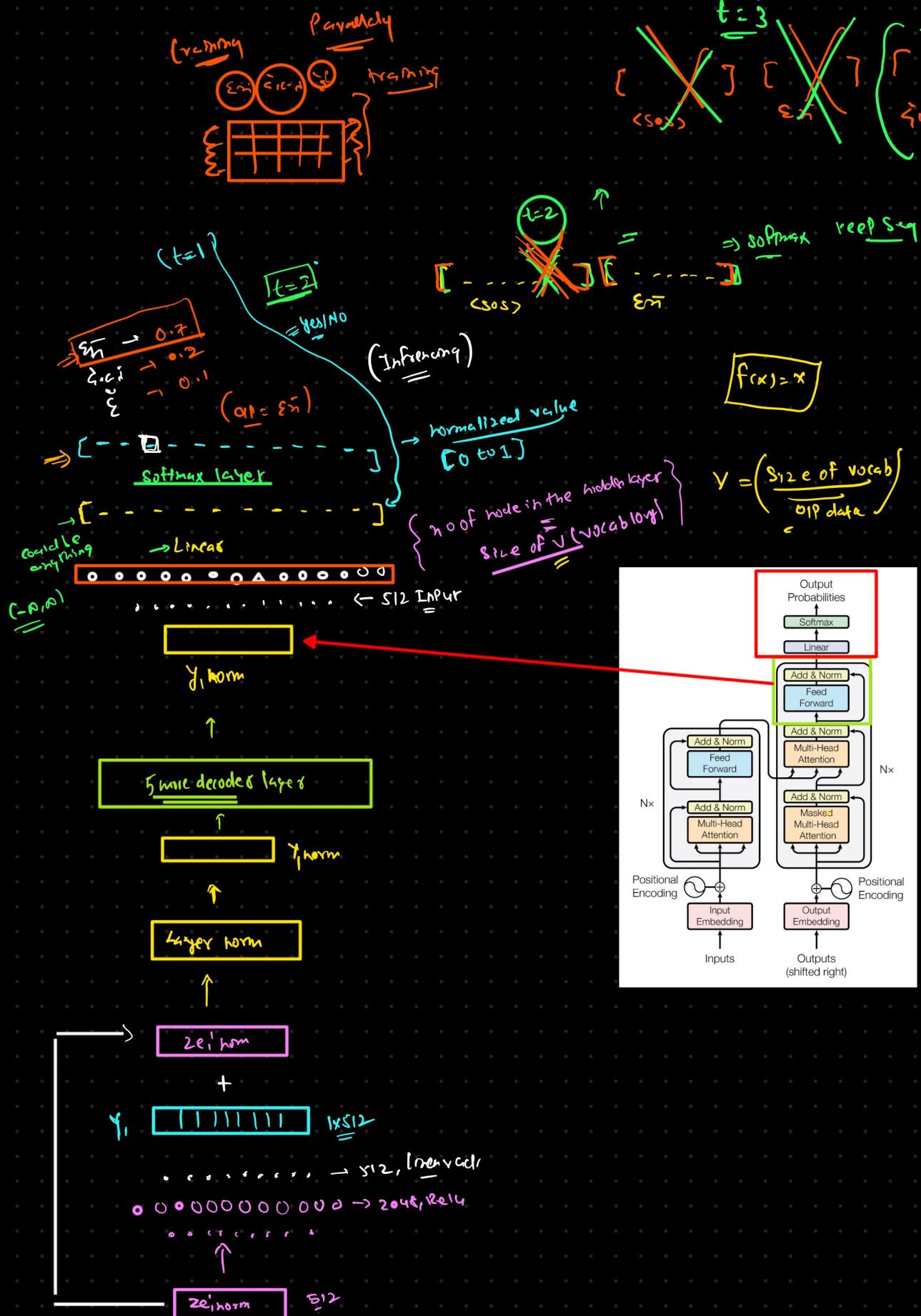
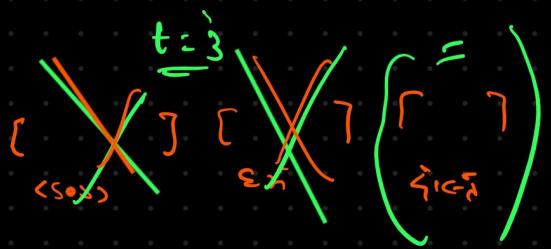
NEW sentence is coming
~~(We are friends) \Rightarrow O/P ?~~
 Prediction



We are friend

<sos>





transformer decoder

(Parallelly)

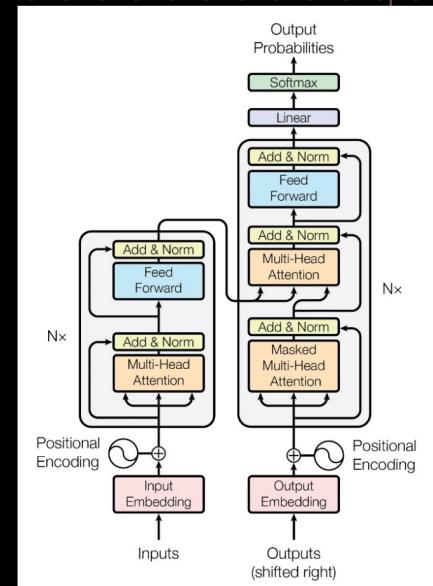
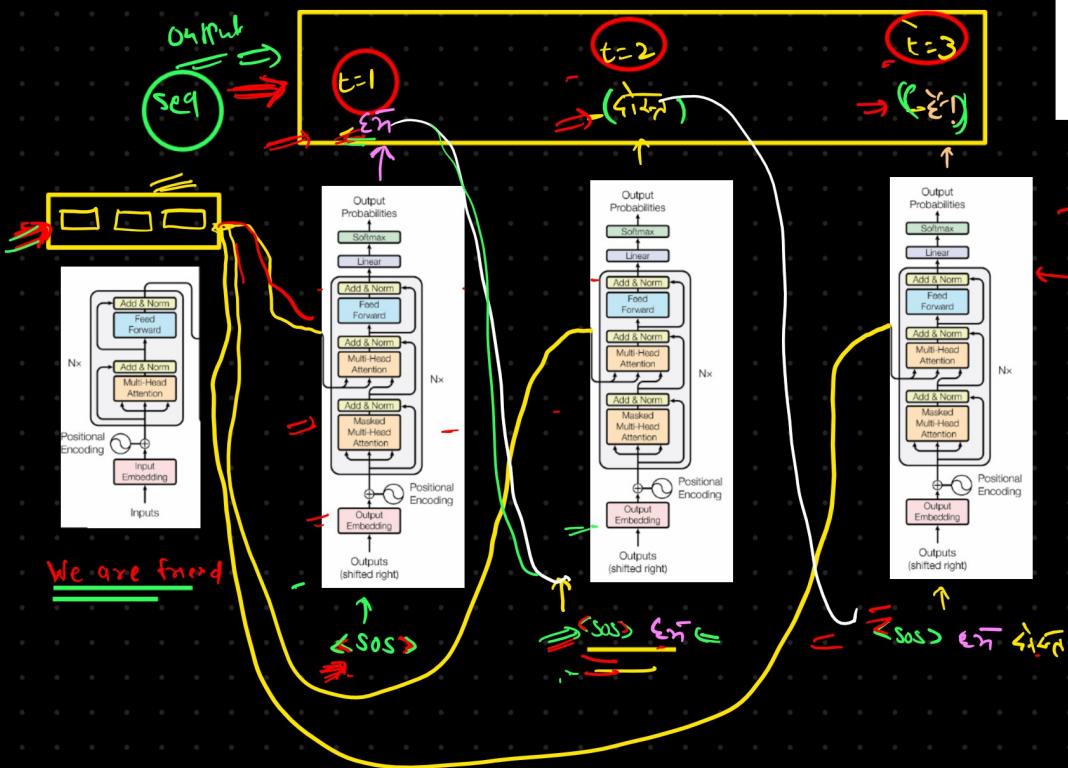
- ↳ it's a non-autoregressive at the time of training
- ↳ It is autoregressive at the time of prediction

Input
We are friends

OIP
($\langle \text{SOS} \rangle \text{ } \tilde{\text{en}}$)

ypred

$\langle \text{EOS} \rangle$

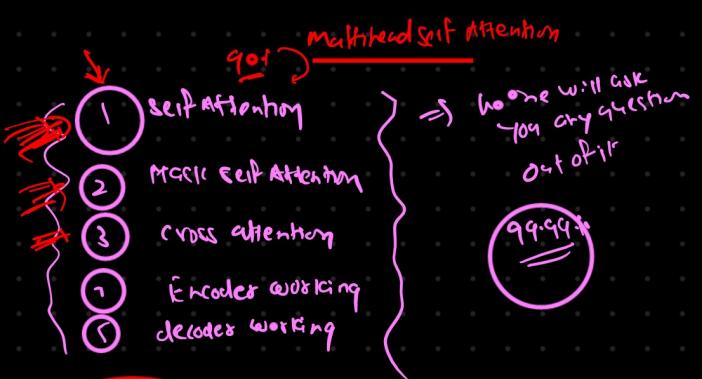


{ 8 →
9 →
10 →

Masking → Data shift

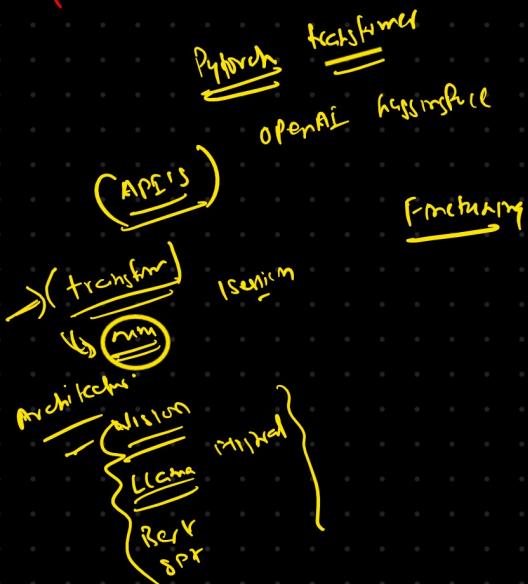
inferring

at the training of masking



Interview

transfer GRU
transfer CRDMR
transfer GPT
transfer BERT
transfer WPE



Isentia PyTorch Isentia → transformer based model
Isentia API's

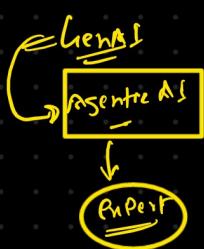
PyTorch Fixing
{ 4-5 }

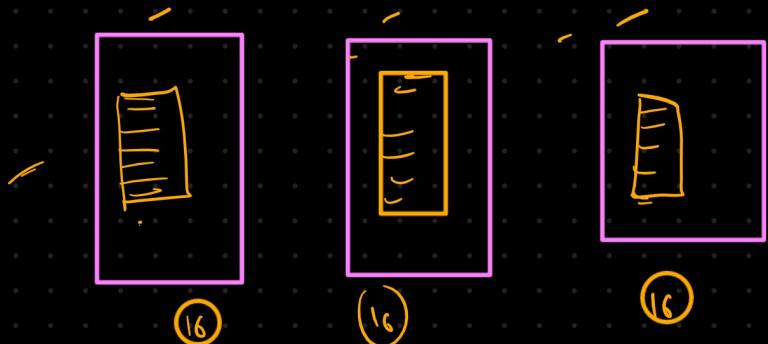
⑧ → ⑨

76.8
6.8

(Project/Code)

{ Permuted
ISOT interviews
question }





Generative Poetrain transformer

