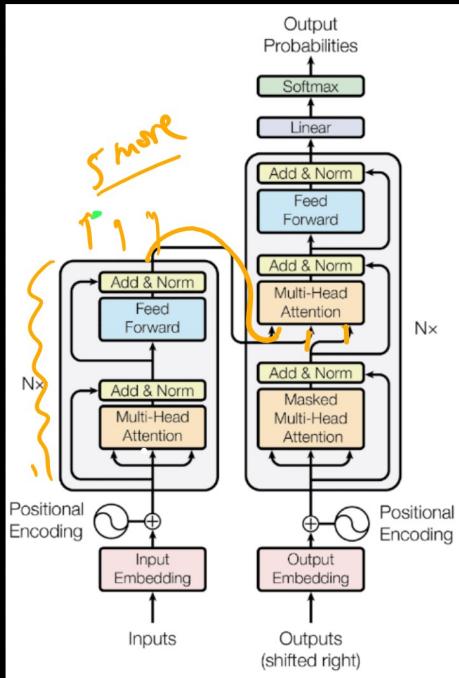


Transformer

Transformer complete summary



Encoder Side:

1. input sequence
2. input embedding
3. positional encoding
4. self-attention, multi-head attention
5. add and layer norm
6. feed forward network(artificial neural network)
7. then again addition(residual connection) and layer normalization
8. Nx (multiple encoder block(in original research paper there was 6 block and all the block were identical means in every block same process happening))

Decoder Side:

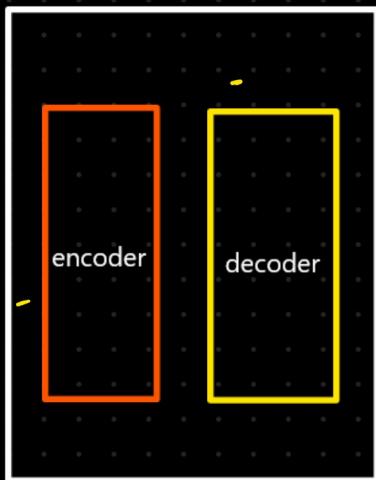
1. output sequence[shifted right means one additional token will add]
2. output sequence embedding[this is word embedding and the dimension of each word is 512]
3. positional encoding for making sure the position of word since we are passing the input parallelly
4. masked multi-head attention
5. addition and layer normalization
6. residual connection
7. multihead attention[it is a cross attention since we are getting input from mask multiheaded attention layer and from encoder as well]
8. then again addition(residual connection) and layer normalization
9. passing output to feed forward network(artificial neural network)
10. addition(residual) and normalization
11. linear activation function
12. then SoftMax activation function
13. final output in the form of probability
14. multiple block of decoder(6 encoder according to the research paper)

this transformer architecture have two sides first at the time of training and second at the time of inferencing

*↳ prediction
(testing)*

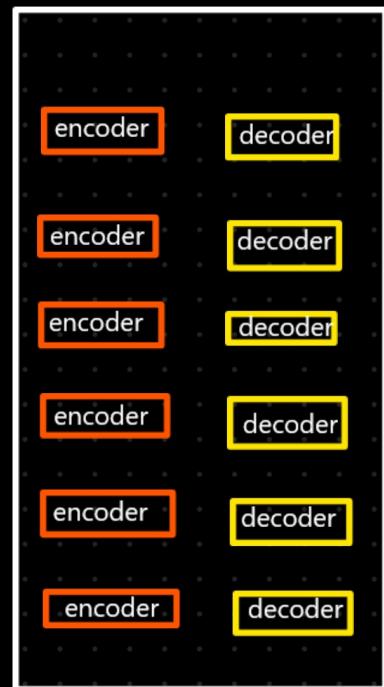


Lets Start with encoder

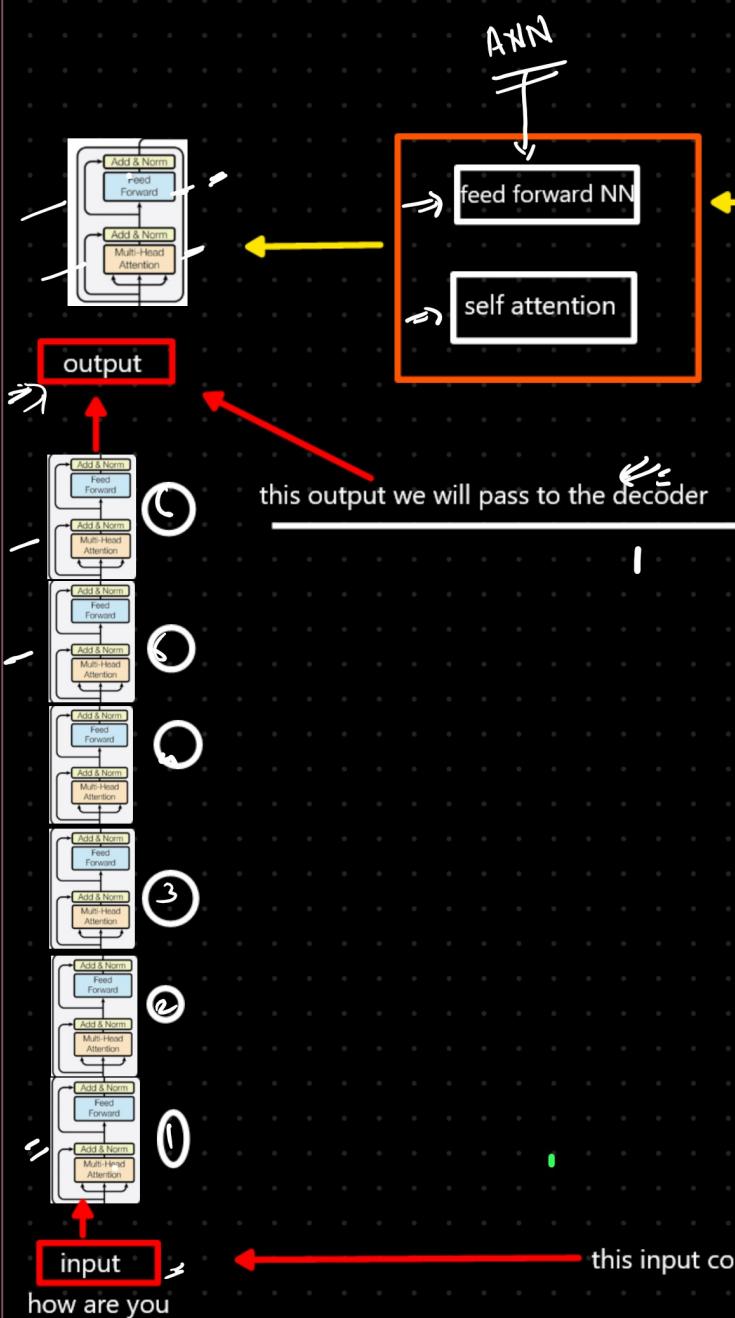


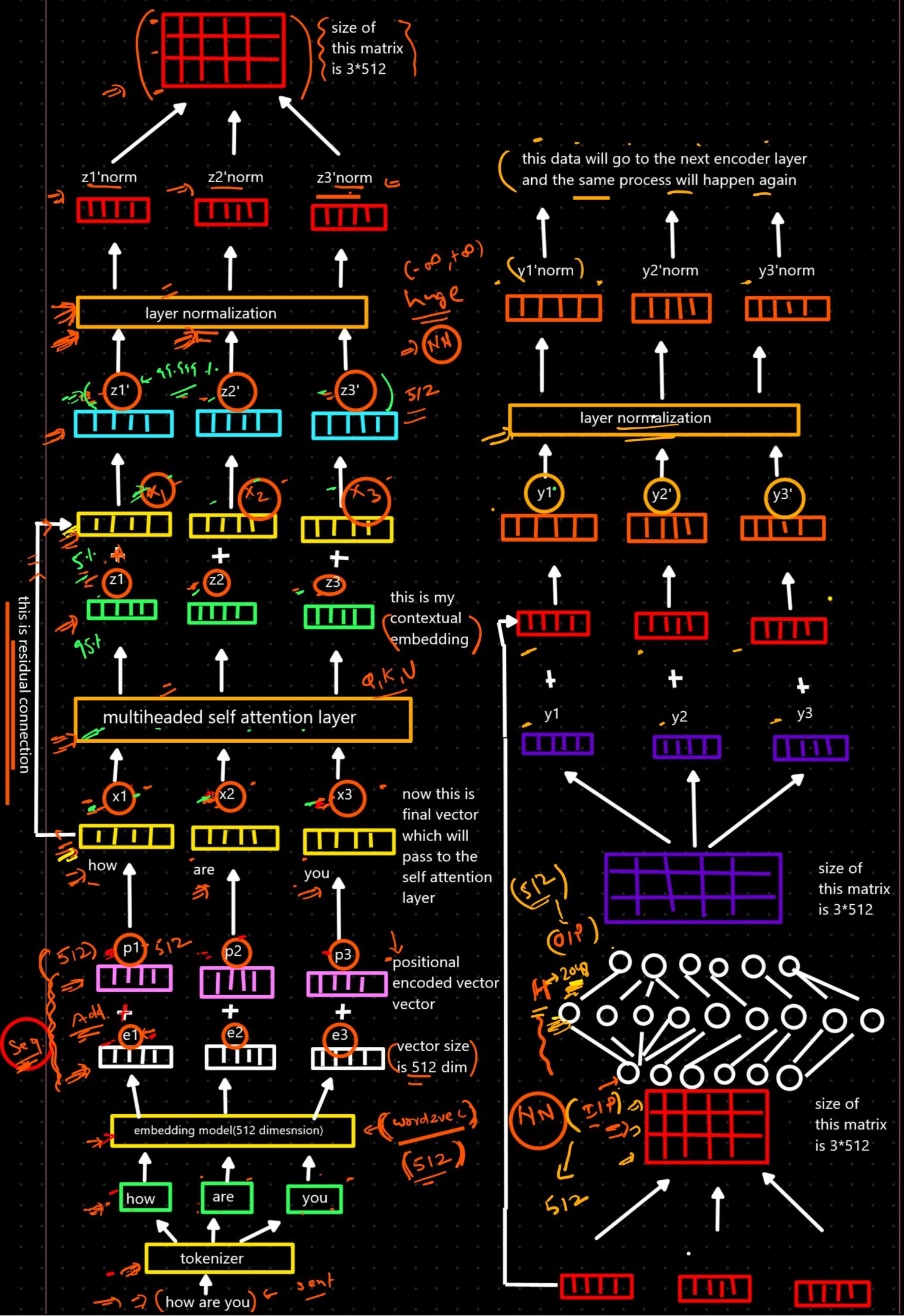
Original Paper

hyperparameter



transformer





here 3 general questions

→ Research Paper

1. why they have used residual connection
2. why used feed forward neural network
3. why they took 6 encoder block

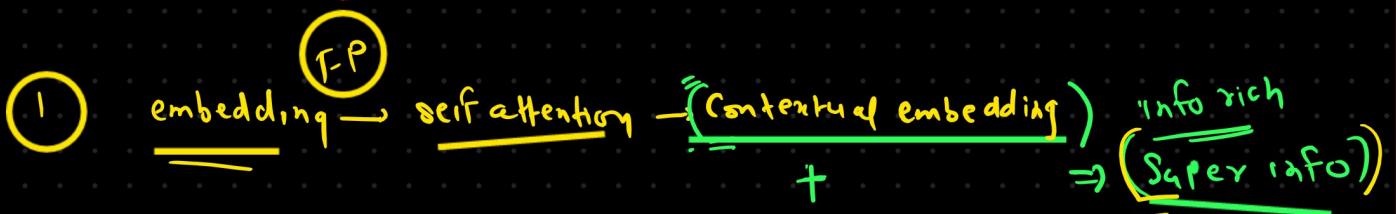
Residual Conn → skip connection

=

= stable training

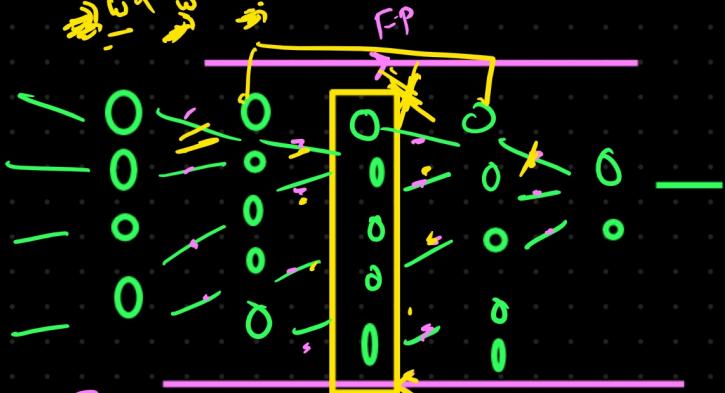
vanishing gradient problem

Resnet ⇒ skip connection



② training → backpropagation → (weight) Self attention

w_Q w_V w_K



neural network (Deep) BP vanishing gradient → (w → small → 0)

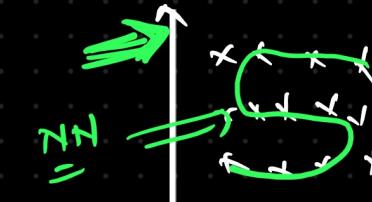
② A NN ⇒ Non-linear information

IIP 512

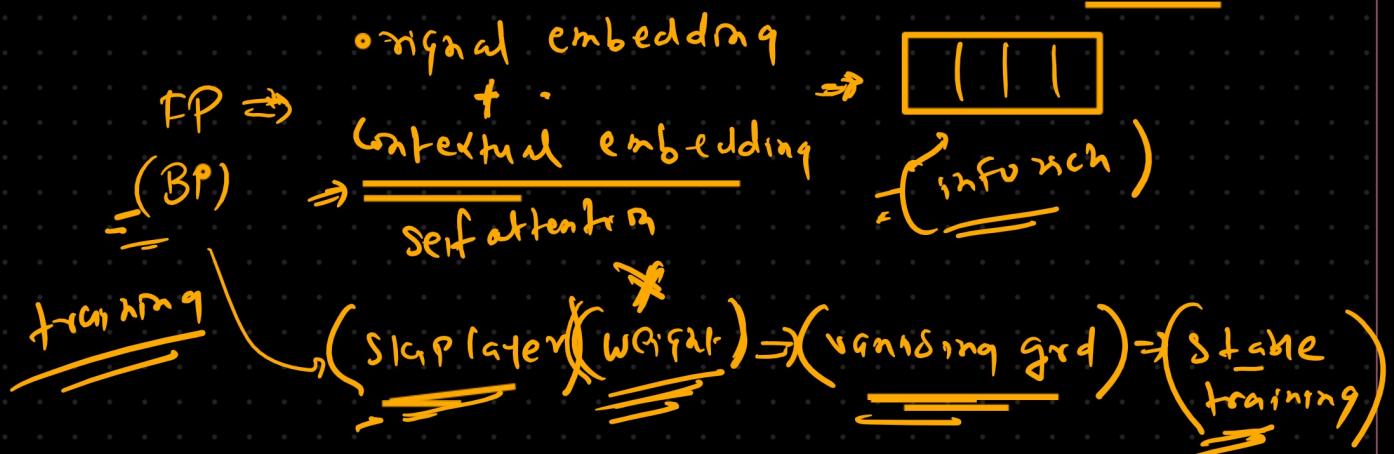


context

linear



Residual Connection → skip connection → identical conn



Add \Rightarrow Non-linearity \Rightarrow hidden layer \Rightarrow Relu \Rightarrow Nonlinearity
out \Rightarrow linear activation fn
(more pattern)

Why 6 encoder & 6 decoder \Rightarrow hyperparameter

experiment

(Concat)

1	1	1	-
---	---	---	---

1	1	1	-
---	---	---	---

dimension ↓
 \Rightarrow increasing
 \rightarrow [1 1 1 1] [1 1 1]

Addition

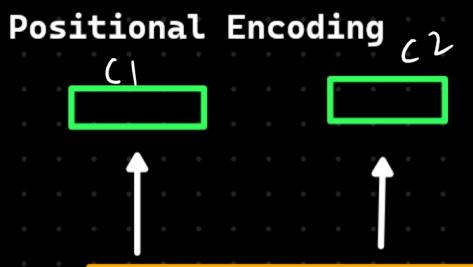
1	2	3	4	-
---	---	---	---	---

+ + + +

5	1	7	8
---	---	---	---

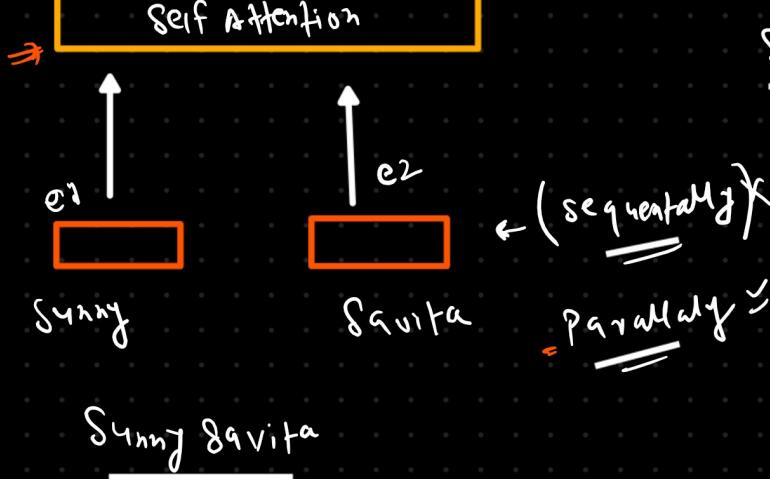
6	8	10	12
---	---	----	----

new vector



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

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



Sunny Savita

$\Rightarrow RNH \rightarrow t=1 \text{ sunny}$
 $t=2 \text{ Savita}$

\downarrow

self attention (seq)
 \Rightarrow parallel.

- ① NLP basics
- ② history (RNN, GRU, LSTM, ELM, ELM Attention)

\Rightarrow ③ Transformer

④ Pytorch

\Rightarrow ⑤ transformer Code in Pytorch

= text
 ↳ multimodal
 ↳ audio based
 ↳ video based

⑥ API (openAI, Gemini, huggingface)

\Rightarrow ⑦ Fracturing (Bert, LLaMa, mstral, CPT ...)

↳ Quantization

↳ LoRA, GLORA

↳ Preframe optimization → RLHF

↳ Framework

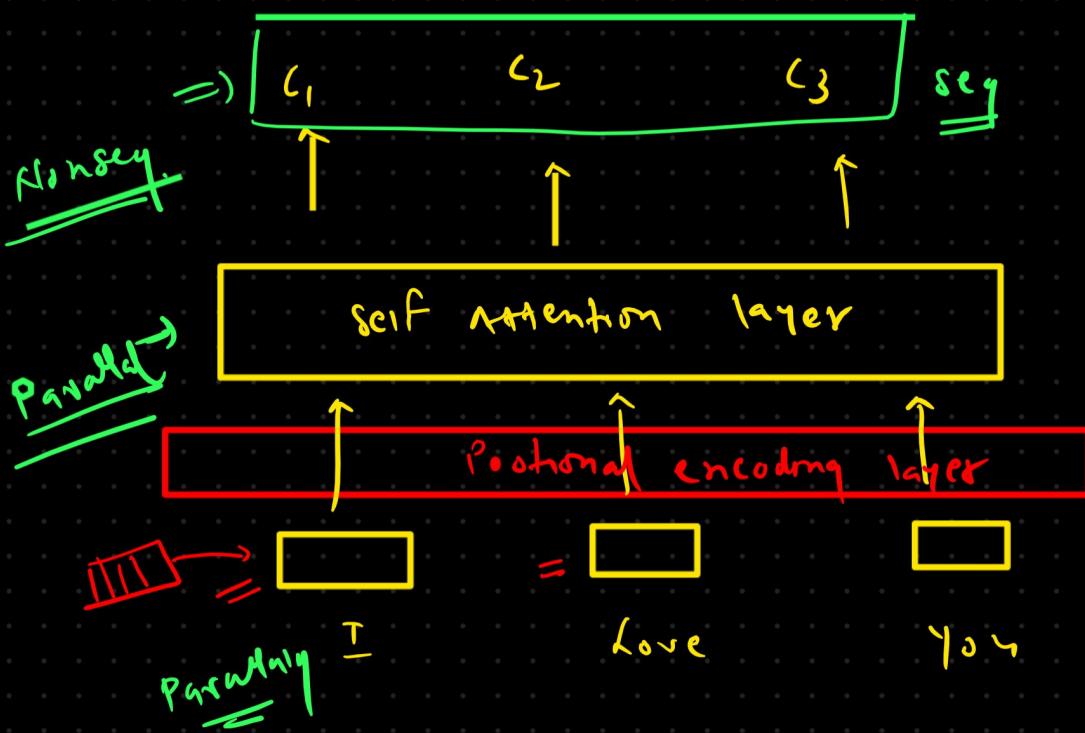
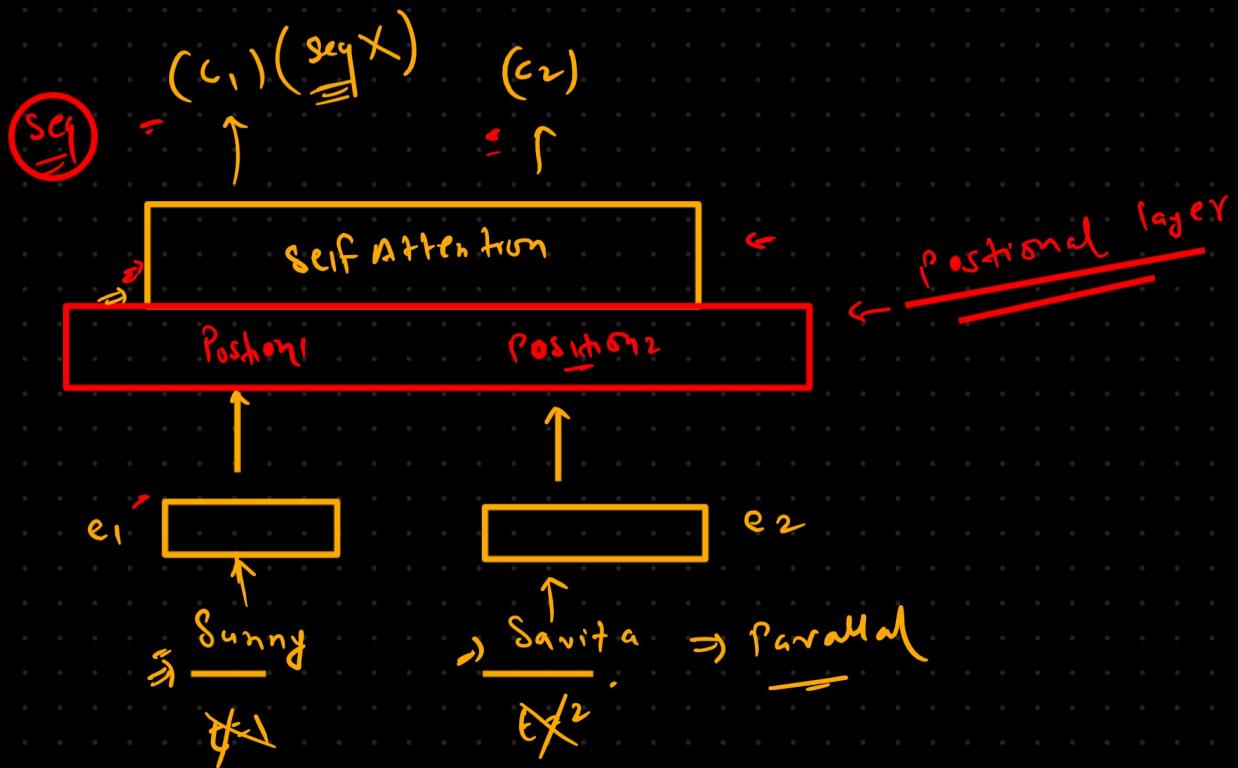
↳ Instruction

↳ SFT (TFL)

Seq → Positional - vector

~~Sunny~~ is taking generate tension & teaching transform

~~Sunny~~ is generating & teaching transform.



(Self Attention) → (masked Self attention) → (Cross Attention) ⑧

\Rightarrow weight constraint

little difference

L. Stxs

long triple weight
matrix

→ (decoder)

The diagram illustrates the architecture of a self-attention layer and its role in generating a final output. At the top, three input tokens c_1 , c_2 , and c_3 are shown with arrows pointing up to a yellow box labeled "self attention layer". Below this box, three vectors e_1 , e_2 , and e_3 are shown, each with a red box around it. Above each vector is a red bracket labeled $p_1 \uparrow$, $p_2 \uparrow$, and $p_3 \uparrow$ respectively. An arrow points from the bottom of the e_1 vector to the word "I". The bottom row consists of the words "I", "LOVE", and "you" followed by a right-pointing arrow. To the right of the "you" word, there is a handwritten note: "You love I". Above the "you" word, there is a handwritten note: "Basic idea". A curved arrow points from the "Basic idea" note to the "you" word. Above the "Basic idea" note, there is a handwritten note: "researched".

Passing \rightarrow book \Rightarrow 3 word $\underline{\underline{3L}} \underline{\underline{302}}$

$\underline{\underline{(NN)}}$ (SelfAttention)

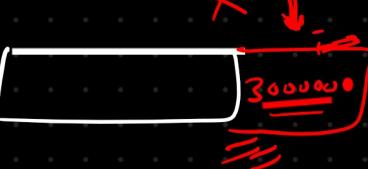
X

desperate

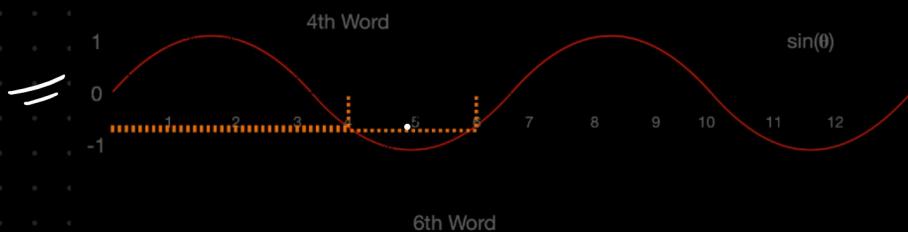
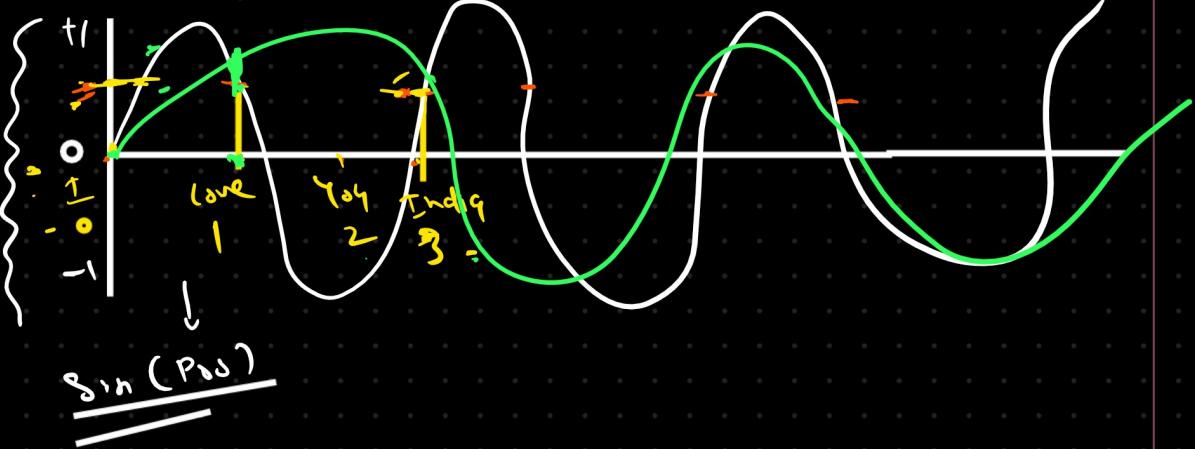


\Rightarrow

(periodic)



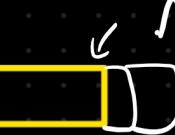
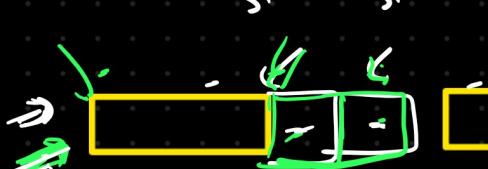
\sin



$\sin(0)$ $\sin(\frac{\pi}{2})$
 $\sin(\pi)$ $\sin(\frac{3\pi}{2})$

$\sin(\theta)$ $\sin(\frac{\pi}{2})$

$\sin(\theta)$ $\sin(\frac{3\pi}{2})$
 No
 Seq of data \approx



I Scalar love

you

longer

$\frac{\text{Pos0}}{I}$

1

2

3

Pos1 love

4

5

6

Pos2 you

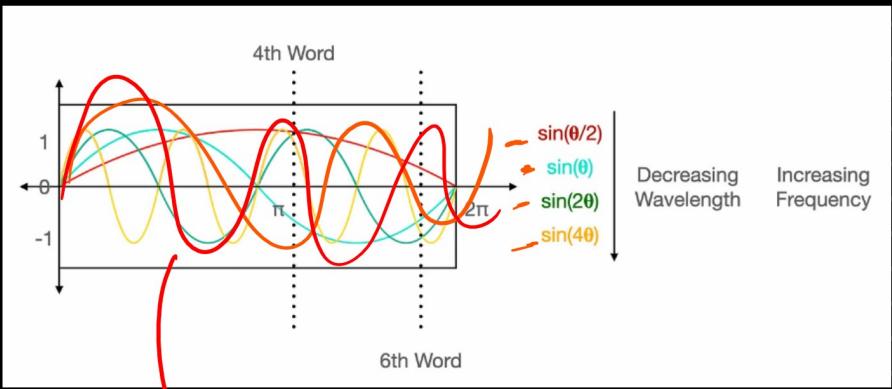
7

8

9

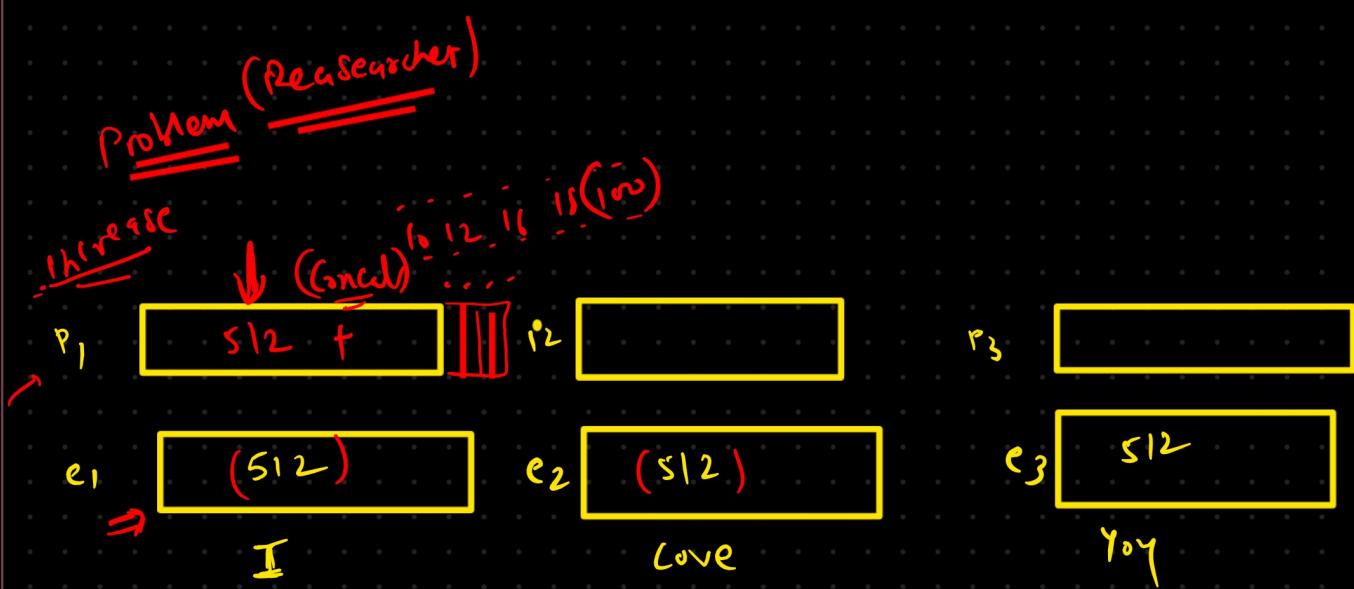
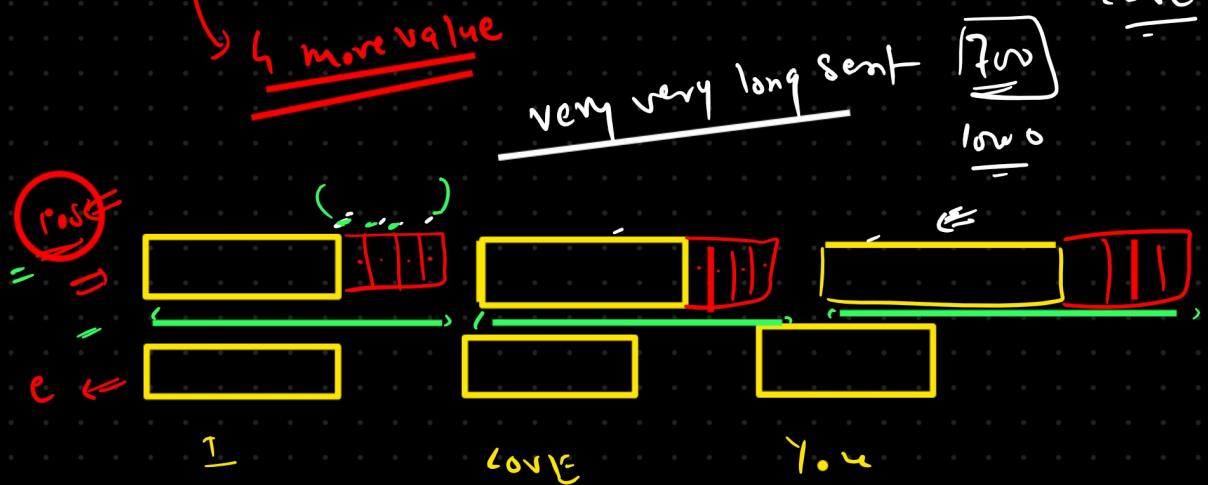
Pos3 India

Pos4 very much



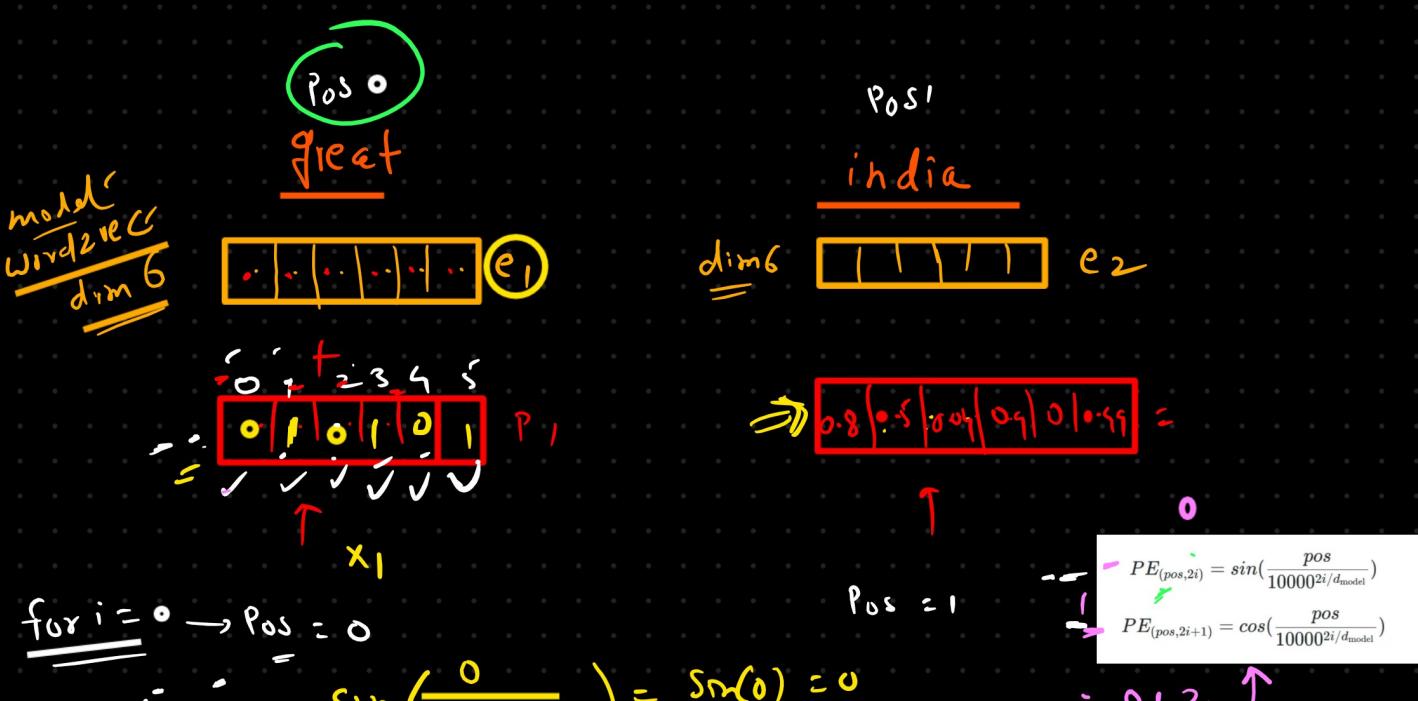
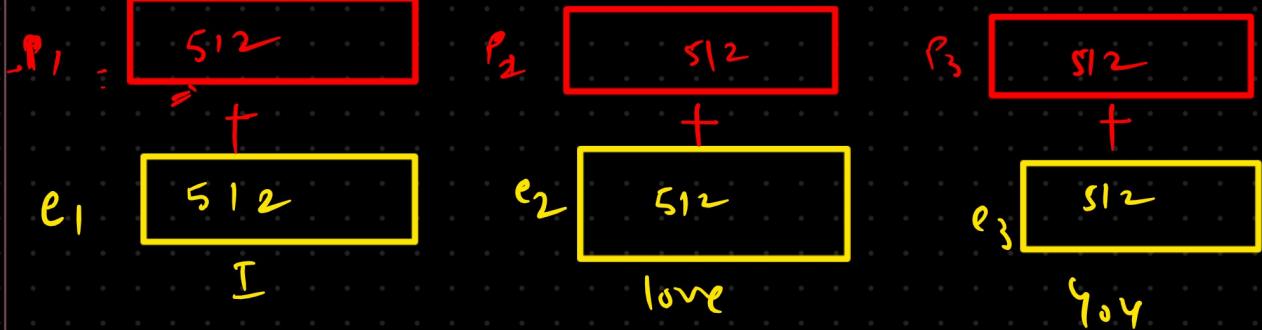
$$\cos(\theta/2)$$

$$\cos(\theta)$$



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

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



$$\underline{\underline{\text{for } i = 0 \rightarrow pos = 0}}$$

$$\Rightarrow PE(0,0) = \sin\left(\frac{0}{10000^{2 \times 0}/6}\right) = \sin(0) = 0$$

$$\Rightarrow PE(0,1) = \cos\left(\frac{0}{10000^{2 \times 0}/6}\right) = \cos(0) = 1$$

$$\underline{\underline{\text{for } i = 1}}$$

$$PE(0,2) = \sin\left(\frac{0}{10000^{2 \times 1}/6}\right) = \sin\left(\frac{0}{10000^{2 \times 1}/3}\right) = 0$$

$$PE(0,3) = \cos\left(\frac{0}{10000^{2 \times 1}/6}\right) = \cos\left(\frac{0}{10000^{2 \times 1}/3}\right) = 1$$

$$\underline{\underline{\text{for } i = 2}}$$

$$PE(0,4) = \sin\left(\frac{0}{10000^{2 \times 2}/6}\right) = \sin\left(\frac{0}{10000^{2 \times 2}/3}\right) = 0$$

$$PE(0,5) = \cos\left(\frac{0}{10000^{2 \times 2}/6}\right) = \cos\left(\frac{0}{10000^{2 \times 2}/3}\right) = 1$$

