

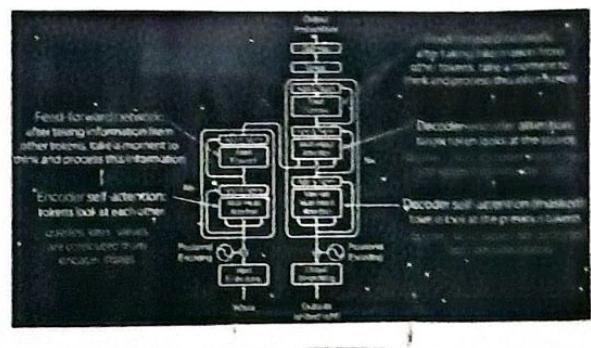
TRANSFORMERS

Introduction To Transformers

- 1) RNN/LSTM/GRU RNN
- 2) Encoder Decoder Architecture
- 3) ATTENTION MECHANISM
- 4) TRANSFORMERS

- 1) Why Transformers?
- 2) Architecture of Transformers?
- 3) SELF ATTENTION $\rightarrow Q, K, V$
- 4) Positional Encoding
- 5) Multi Head ATTENTION
- 6) Combining the Working of Transformers

Architecture



Generative AI \rightarrow UN, Multimodel

BERT, GPT



Open AI \rightarrow ChatGPT

↑
GPT-4

0) What And Why \rightarrow Transformers

Transformers in natural language processing (NLP) are a type of deep learning model that use self-attention mechanisms to analyze and process natural language data. They are encoder-decoder models that can be used for many applications, including machine translation. \Rightarrow Seq2Seq Task

Eg: Language Translation \rightarrow Google Translation

English \rightarrow French

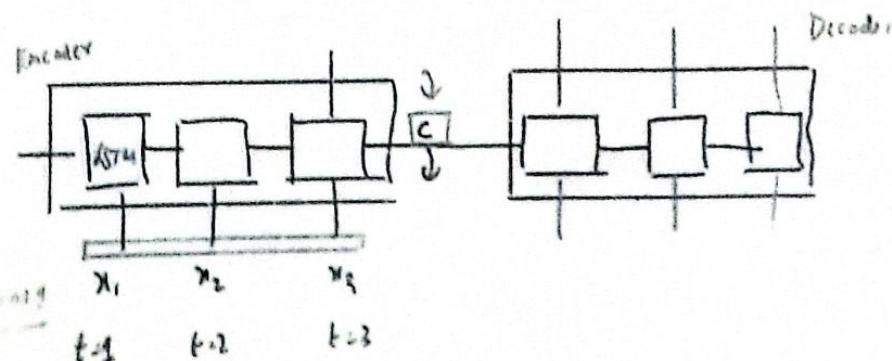
ip \Rightarrow Many \rightarrow o/p: Many. {length of the sentence}

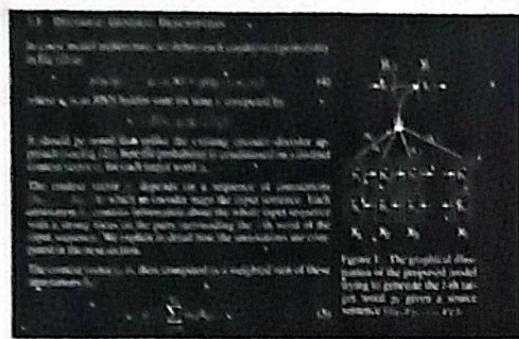
Encoder - Decoder

Sentence length $\uparrow\uparrow$

Beam Score $\downarrow\downarrow$

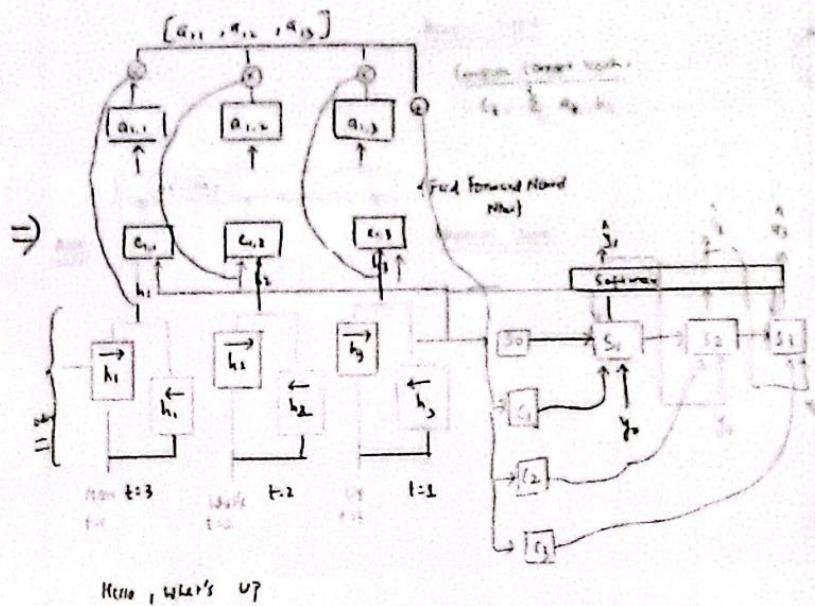
Length Sentence $\uparrow\uparrow$





Additional Context \rightarrow Decoder

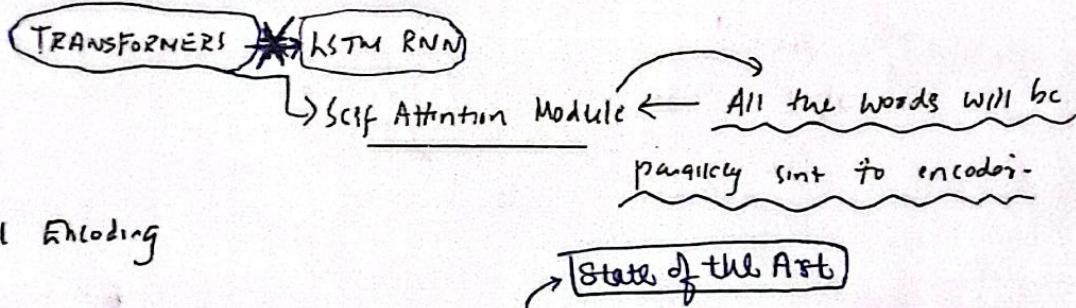
long sentence Accuracy ↑↑



① Attention Mechanism \rightarrow (Problem)

① Parallelly, we cannot send all the words in a sentence \rightarrow **Scalable**

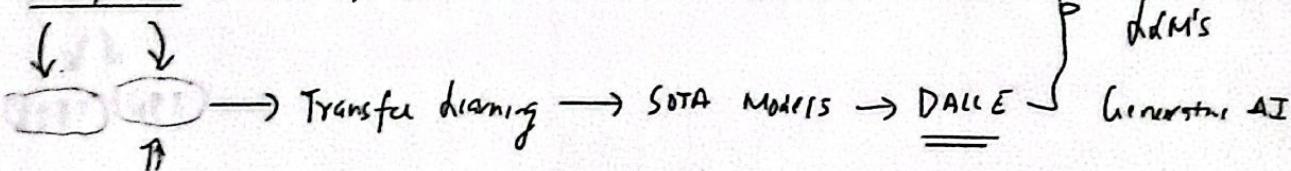
DATASET \rightarrow Huge \rightarrow Scalable with Respect to Training.



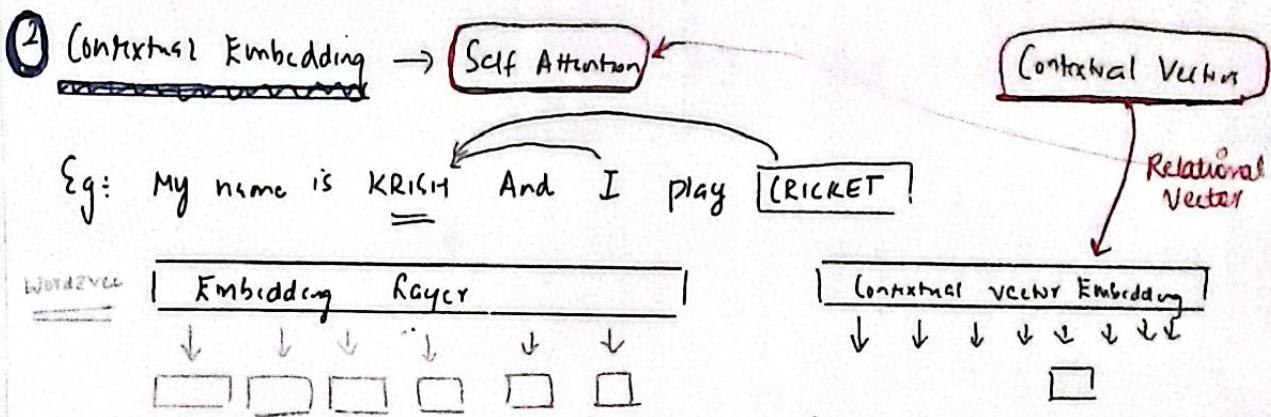
✓ Transformer \uparrow DATASET \rightarrow Amazing SOTA \leftarrow NLP

✓ Transfer learning \rightarrow MultiModel Task \rightarrow NLP + Image \leftarrow

✓ Transformers \div AI Space \rightarrow SOTA Model \rightarrow

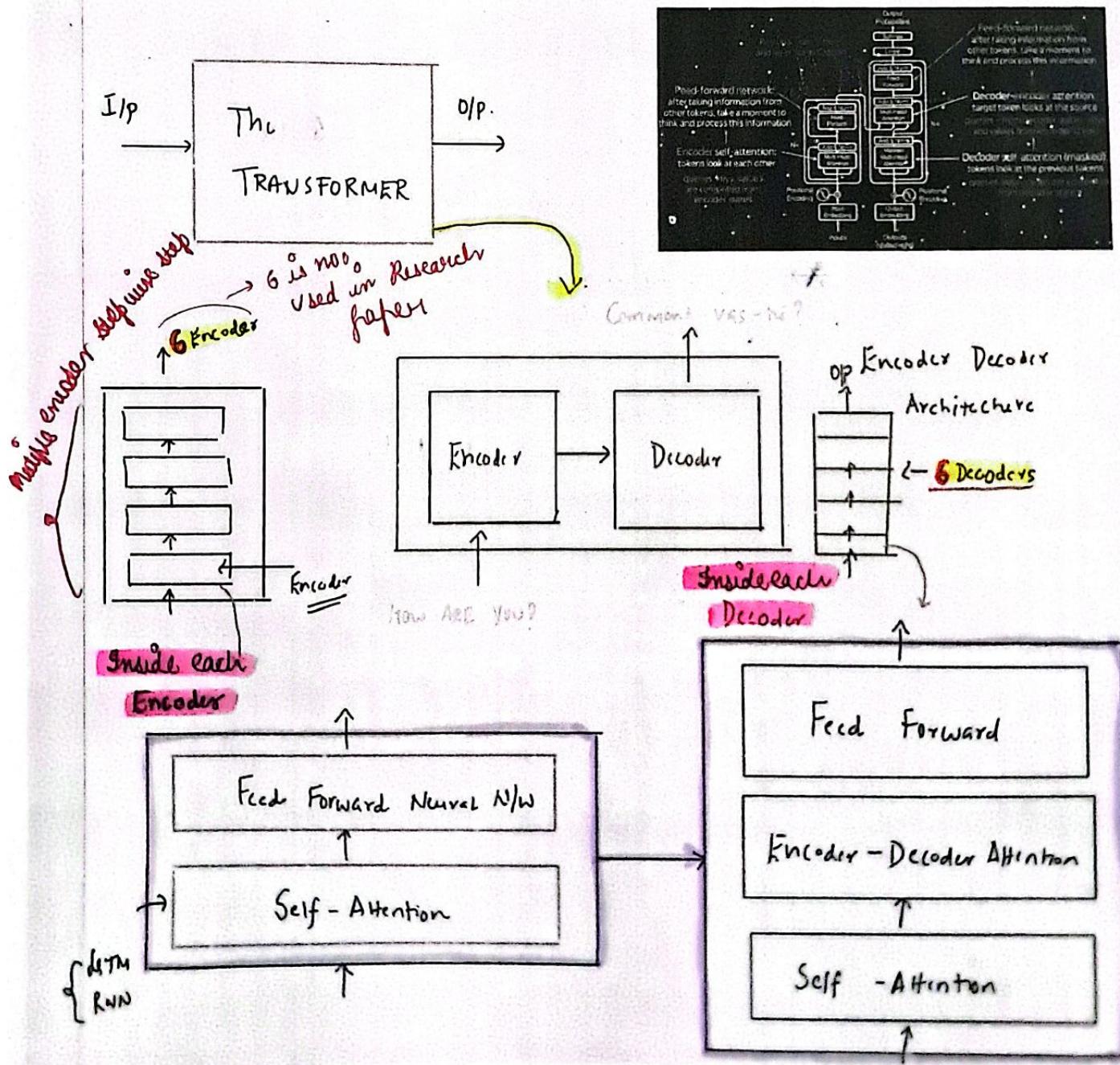


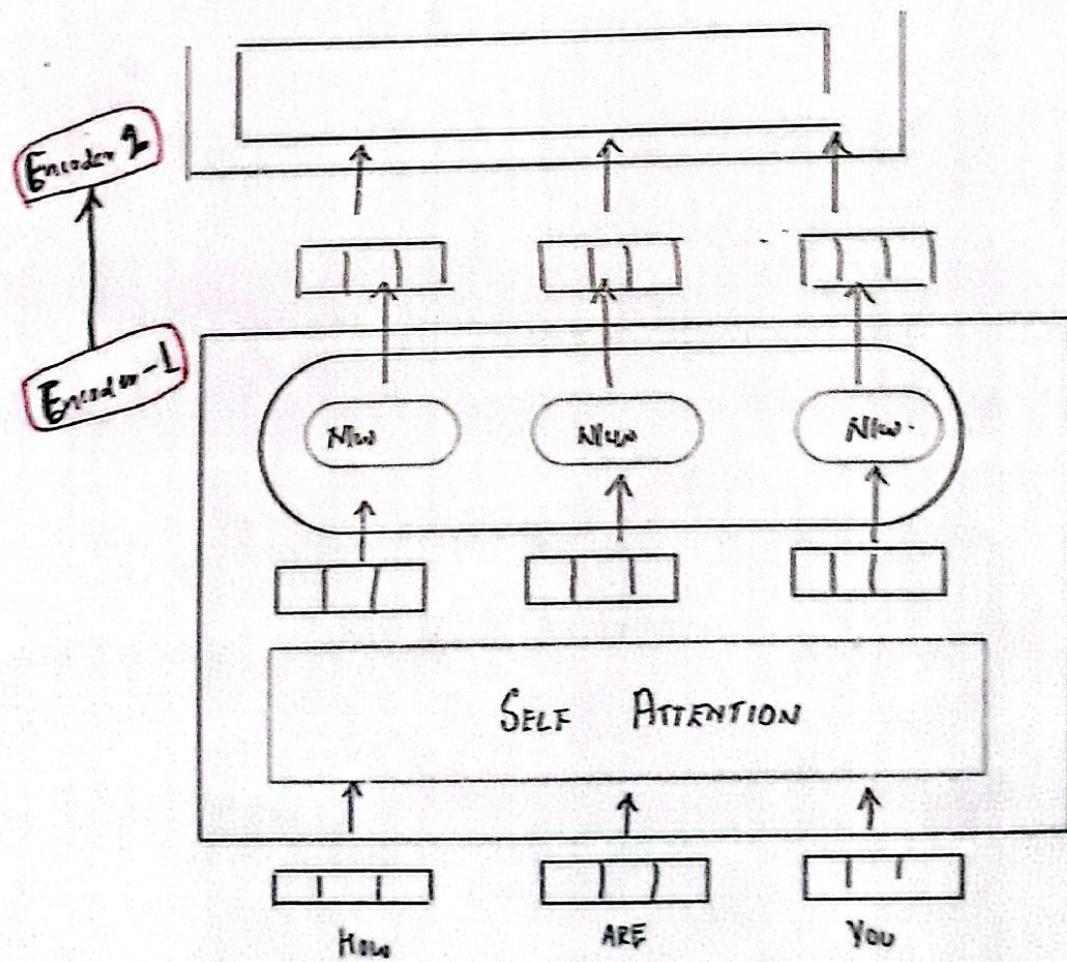
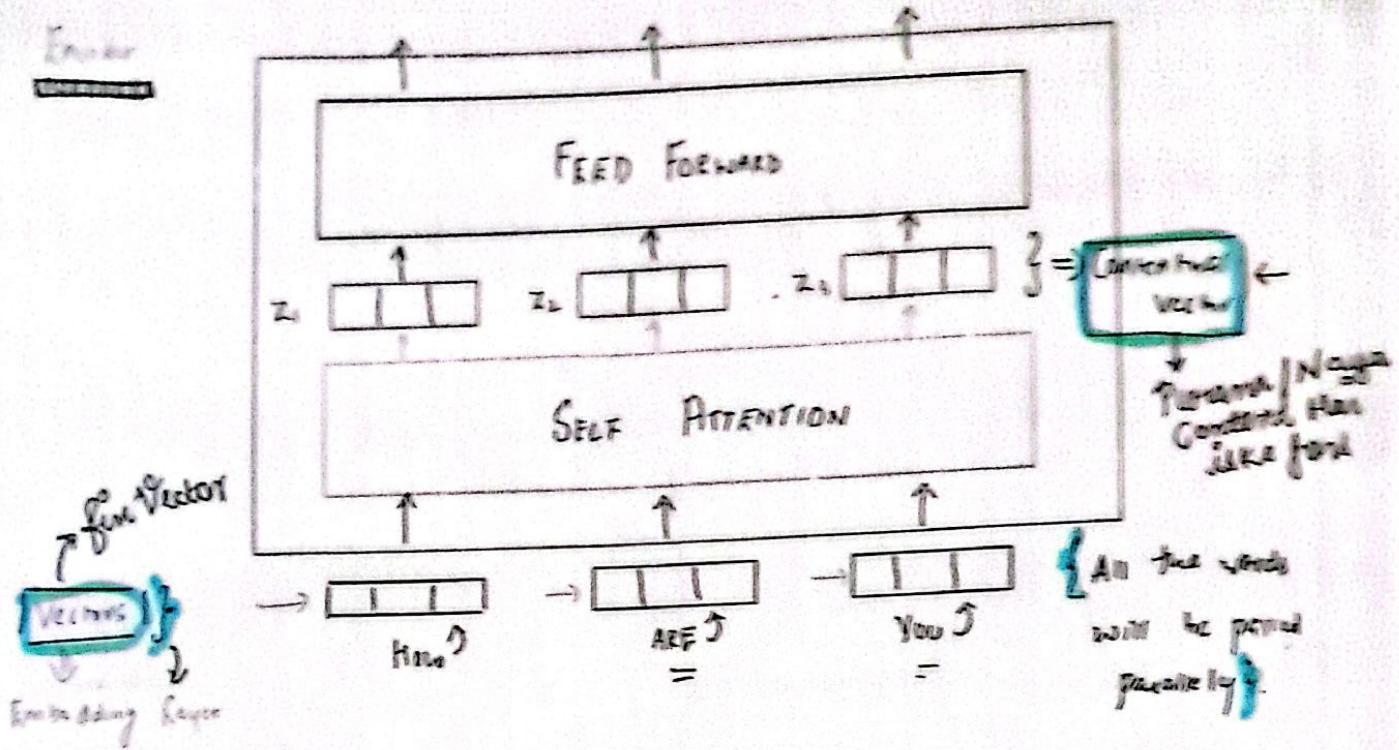
Train Huge Data



② Basic Transformer Architecture

{Seq2Seq Task} → Language Translation {Eng → French}





Self Attention At a Higher Level

Word Embedding

4 4 4 8 4 4 4 4 4 4 4 4 4 4

SELF ATTENTION

4 4 4 8 4 4 4 4 4 4 4 4 4 4



4 4 4 8 4 4 4 4 4 4 4 4 4 4

The

The

→ [Cat]

→ [the]

→ [the]

1 → [Sat]

Sat

{ Contextual Embedding }

Row 2 → 0n

0n

3 → [the]

the

mat

mat

mat → [the]



Self Attention In Detail

① To create 3 vectors from each of the encoder (IP) using key, query, value

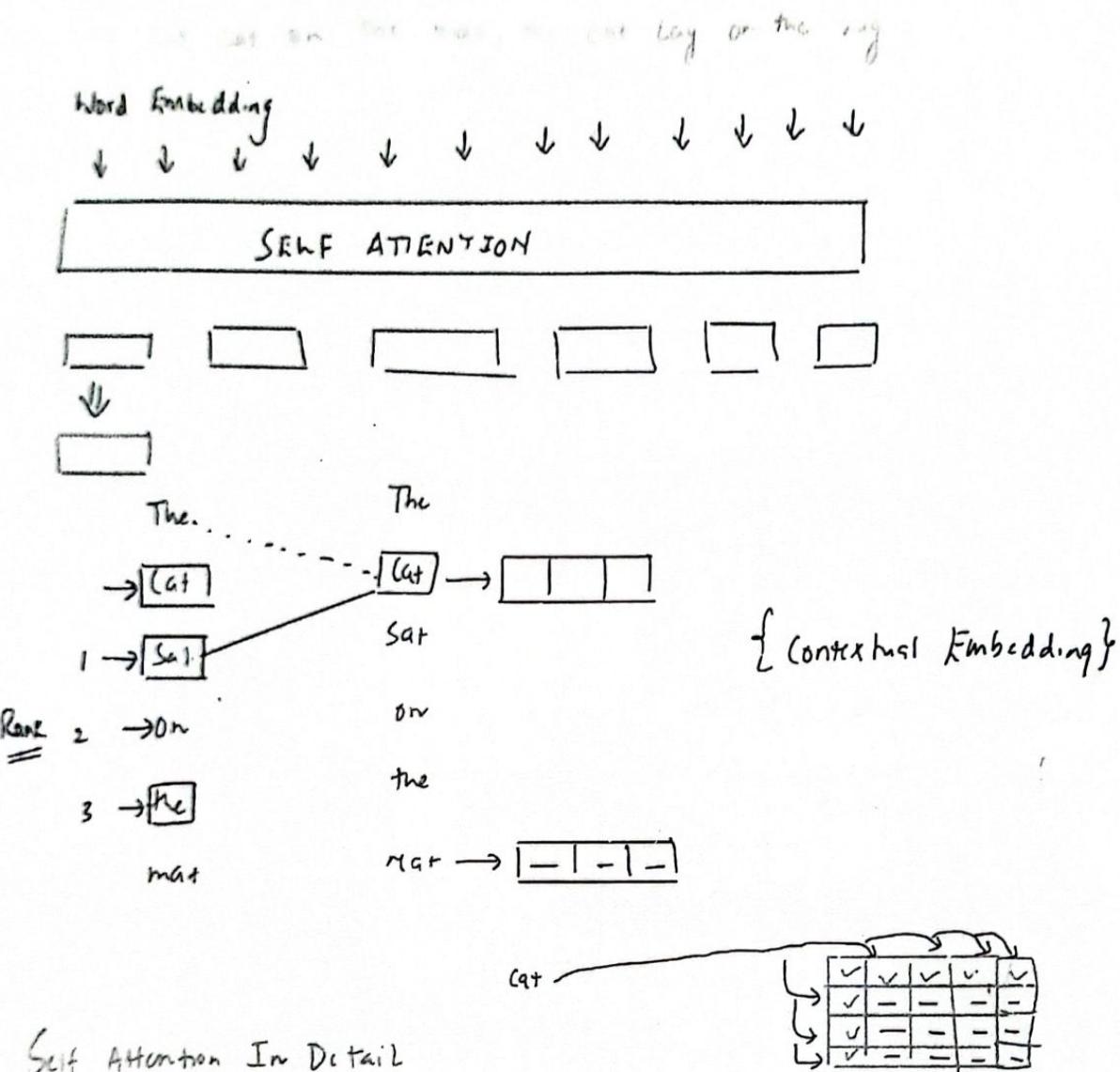
key, query, value \rightarrow Contextual Embedding

↳ \rightarrow Search keywords \rightarrow {Query}

Query \rightarrow {Key} $\xrightarrow{\text{Tags}}$ Tags
Description $\xrightarrow{\text{Tags}}$ Tags
 $\xrightarrow{\text{Title}}$

② Second step in calculating self-attention is to calculate the dot

Self Attention At a Higher Level



① To create 3 vectors from each of the encoder i/p. Query vector,

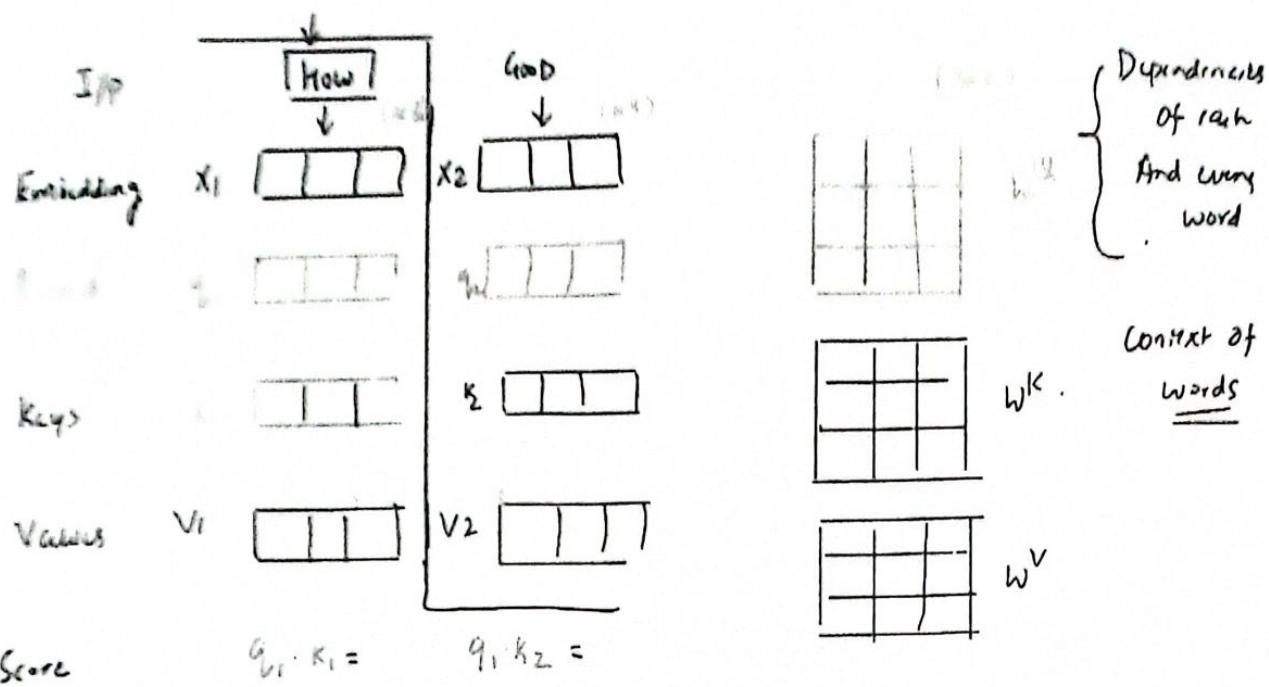
Key vector, Value vector \Rightarrow Contextual Embedding

Eg. $YT \Rightarrow$ Search keywords $\xrightarrow{\text{Topic}}$ {Query}.

Query \rightarrow {Key} $\xrightarrow{\text{Tags}}$ Tags
 $\xrightarrow{\text{Description}}$ Description \Rightarrow Value $\xrightarrow{\text{Title}}$ Value \Rightarrow O/P Video

② Second step in calculating self attention is to calculate the Score.

The score determines how much focus to place on the other part of the sequence.



- ② Focus to place on other part of the I/P Sentence as we encode a word at certain position

Divide \sqrt{dk}

More stable

More gradients

Dimension = 3 Dimension $\frac{1}{\sqrt{dk}}$

Output $\uparrow \uparrow \uparrow$

$\frac{\text{Variance}}{\sqrt{dk}} \approx \frac{\text{Variance } \uparrow \uparrow \uparrow}{\sqrt{dk}}$

Softmax

$\frac{[0-1]}{0.88} + \frac{[0-1]}{0.12} = 1$

Softmax \Rightarrow The Softmax score determines how much each word will be experienced at this position

$\times \quad \times \quad \times$

Value Vectors $v_1 [1|1|1]$ $v_2 [1|0.6|0.12]$

$0.88 [1|1|1]$

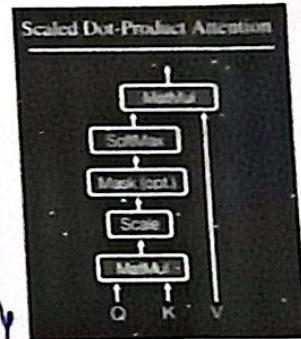
$\frac{0.88 \ 0.88 \ 0.88}{0.88 \ 0.88 \ 0.88} + \frac{0.12 \ 0.6 \ 0.12}{0.12 \ 0.6 \ 0.12}$

$\downarrow \quad \leftarrow$

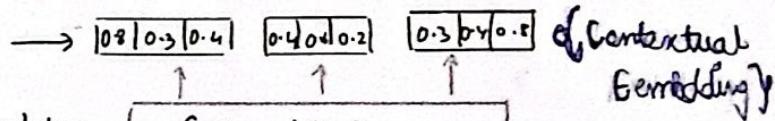
$\text{softmax} \rightarrow [1|1|1] \quad \rightarrow [1|1|1|1]$

Self Attention At Higher And Detailed level

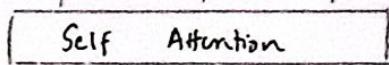
Self-attention, also known as scaled dot-product attention, is a crucial mechanism in the transformer architecture that allows the model to weigh the importance of different tokens in the input sequence relative to each other



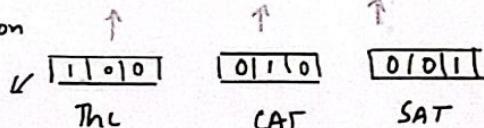
Idea :



✓ Language Translation



✓ Text Summarization



{ Sentence, Dataset }

Embedding

layer → Fixed Vector

1) Inputs: Queries, keys, And Values

Model → Queries, Keys And Values

1. Query Vectors (Q):

Role: Query vectors represent the token for which we are calculating the attention. They help determine the importance of other tokens in the context of the current token.

Importance:

Focus Determination: Queries help the model decide which parts of the sequence to focus on for each specific token. By calculating the dot product between a query vector and all key vectors, the model assesses how much attention to give to each token relative to the current token.

Contextual Understanding: Queries contribute to understanding the relationship between the current token and the rest of the sequence, which is essential for capturing dependencies and context.

2. Key Vectors (K):

Role: Key vectors represent all the tokens in the sequence and are used to compare with the query vectors to calculate attention scores.

Importance:

Relevance Measurement: Keys are compared with queries to measure the relevance or compatibility of each token with the current token. This comparison helps in determining how much attention each token should receive.

Information Retrieval: Keys play a critical role in retrieving the most relevant information from the sequence by providing a basis for the attention mechanism to compute similarity scores.

3. Value Vectors (V):

Role: Value vectors hold the actual information that will be aggregated to form the output of the attention mechanism.

Importance:

Information Aggregation: Values contain the data that will be weighted by the attention scores. The weighted sum of values forms the output of the self-attention mechanism, which is then passed on to the next layers in the network.

Context Preservation: By weighting the values according to the attention scores, the model preserves and aggregates relevant context from the entire sequence, which is crucial for tasks like translation, summarization, and more.

Input Sequence = $["\text{The}", "\text{CAT}", "\text{SAT}"]$

• Let's assume

Embedding Size = 4

$$Q, K, V \Rightarrow 4$$

$$\xrightarrow{b} [1 \ 0 \ 1 \ 0] \rightarrow$$

$[1 \ 0 \ 1 \ 0]$ ← Content vector

$[1 \ 0 \ 1 \ 0]$ $\square \ \square$

Context, Dataset

$$E_{\text{The}} = [1 \ 0 \ 1 \ 0]$$

$$E_{\text{SAT}} = [1, 1, 1, 1]$$

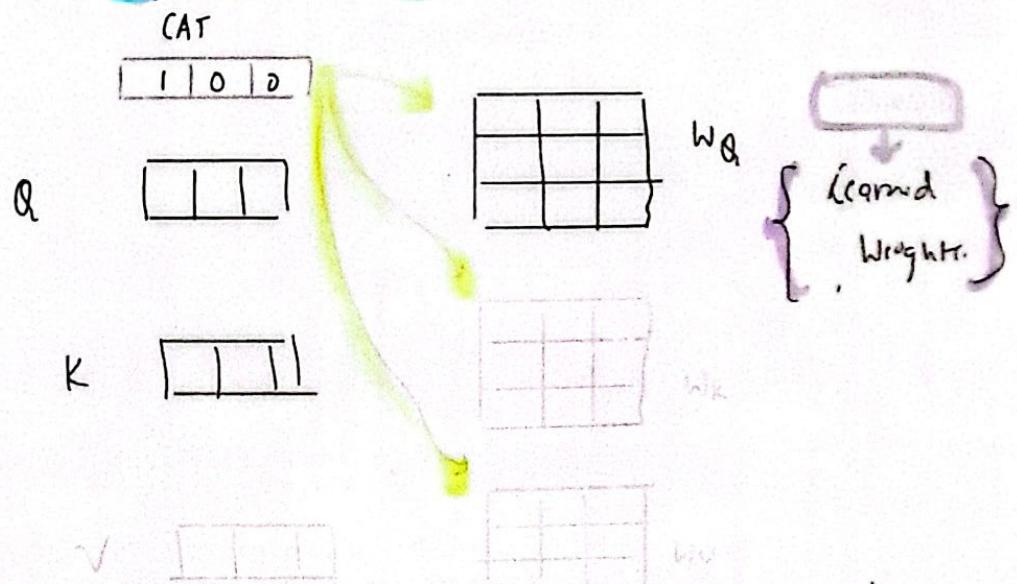
$$E_{\text{CAT}} = [0 \ 1 \ 0 \ 1]$$

Listen, We have vectors like Content Vector, Context Vector, for that with different datasets & sentences the Content Vector will change due to their mutual connections. So to solve this we have Q, K, V to Averag the Content Vector.

up with Self Attention

② Linear Transformation

We create Q, K, V by multiplying the embeddings by learned weights matrices W_Q , W_K and W_V



Let's consider

$$W_Q = W_K = W_V = I$$

$$\begin{bmatrix} 1 & 0 & 0 \\ 0 & 1 & 0 \\ 0 & 0 & 1 \end{bmatrix}$$

$$Q_{\text{The}} = \begin{bmatrix} 1 & 0 & 0 \end{bmatrix} \cdot \begin{bmatrix} 1 & 0 & 0 \\ 0 & 1 & 0 \\ 0 & 0 & 1 \end{bmatrix} = \begin{bmatrix} 1 & 0 & 0 \end{bmatrix}$$

$$K_{\text{The}} = \begin{bmatrix} 1 & 0 & 0 \end{bmatrix} \cdot \begin{bmatrix} 1 & 0 & 0 \\ 0 & 1 & 0 \\ 0 & 0 & 1 \end{bmatrix}$$

$$V_{\text{The}} = \begin{bmatrix} 1 & 0 & 0 \end{bmatrix} \cdot \begin{bmatrix} 1 & 0 & 0 \\ 0 & 1 & 0 \\ 0 & 0 & 1 \end{bmatrix}$$

$$Q_{\text{The}} = K_{\text{The}} = V_{\text{The}} = \begin{bmatrix} 1 & 0 & 0 \end{bmatrix}$$

$$② Q_{\text{CAT}} = K_{\text{CAT}} = V_{\text{CAT}} = \begin{bmatrix} 0 & 1 & 0 \end{bmatrix}$$

$$③ Q_{\text{SAT}} = K_{\text{SAT}} = V_{\text{SAT}} = \begin{bmatrix} 1 & 1 & 1 & 1 \end{bmatrix}$$

Why we do

→ To know how much attention to give to each token relative to the current token

$$\begin{bmatrix} 1 & 0 & 1 & 0 \end{bmatrix} \begin{bmatrix} ? \\ ? \\ ? \\ ? \end{bmatrix}$$

③ Compute Attention Scores

The

$$\text{Score}(Q_{\text{The}}, K_{\text{The}}) = [1 \ 0 \ 1 \ 0] \cdot [1 \ 0 \ 1 \ 0]^T = 2$$

$$\text{Score}(Q_{\text{The}}, K_{\text{CAT}}) = [1 \ 0 \ 1 \ 0] \cdot [0 \ 1 \ 0 \ 1]^T = 0$$

$$\text{Score}(Q_{\text{The}}, K_{\text{SAT}}) = [1 \ 0 \ 1 \ 0] \cdot [1 \ 1 \ 1 \ 1]^T = 2$$

For the token CAT

$$\text{Score}(Q_{\text{CAT}}, K_{\text{The}}) = [0 \ 1 \ 0 \ 1] \cdot [1 \ 0 \ 1 \ 0]^T = 0$$

$$\text{Score}(Q_{\text{CAT}}, K_{\text{CAT}}) = [0 \ 1 \ 0 \ 1] \cdot [0 \ 1 \ 0 \ 1]^T = 2$$

$$\text{Score}(Q_{\text{CAT}}, K_{\text{SAT}}) = [0 \ 1 \ 0 \ 1] \cdot [1 \ 1 \ 1 \ 1]^T = 2.$$

For the Token SAT

$$\text{Score}(Q_{\text{SAT}}, K_{\text{The}}) = [1 \ 1 \ 1 \ 1] \cdot [1 \ 0 \ 1 \ 0]^T = 2$$

$$\text{Score}(Q_{\text{SAT}}, K_{\text{CAT}}) = 2$$

$$\text{Score}(Q_{\text{SAT}}, K_{\text{SAT}}) = 4$$

④ Scaling =

We take up the scores and scale down by dividing the scores by the

$$\sqrt{d_K} \Rightarrow d_K = 4 \quad \sqrt{d_K} = 2.$$

↓ dimension of key vector

Scaling in the attention mechanism is crucial to prevent the dot product from growing too large. \Rightarrow Ensure stable gradients during Training
↓ Here dot product is b/w query & key vector.

If d_K is large \rightarrow

① Gradient Exploding

② Softmax Saturates

other token \Rightarrow weight ≈ 0

$$Q = \begin{bmatrix} 2 & 3 & 4 & 1 \end{bmatrix} \quad K_1 = \begin{bmatrix} 1 & 0 & 1 & 0 \end{bmatrix} \quad K_2 = \begin{bmatrix} 0 & 1 & 0 & 1 \end{bmatrix}$$

Without Scaling

$$Q \cdot K_1^T = 2 \times 1 + 3 \times 0 + 4 \times 1 + 1 \times 0 = 2 + 4 = 6.$$

$$Q \cdot K_2^T = 2 \times 0 + 3 \times 1 + 4 \times 0 + 1 \times 1 = 0 + 3 + 0 + 1 = 4$$

∴ Score $[6, 4] \Rightarrow$ Scaling Not Applied

$$\text{Softmax} \left(\begin{bmatrix} 6, 4 \end{bmatrix} \right) = \left[\frac{e^6}{e^6 + e^4}, \frac{e^4}{e^6 + e^4} \right] = \left[\frac{e^6}{e^6(1 + e^{-2})}, \frac{e^4}{e^4(e^2 + 1)} \right]$$

② Property of Softmax

$$\left(\begin{bmatrix} 10, 1 \end{bmatrix} \right) = \left[\frac{1}{(1 + e^{-2})}, \frac{1}{(e^2 + 1)} \right]$$

Dot product \uparrow = Large values
 \downarrow = Small values

Job description
 value K_1 we bring
 down. Be prepared
 the weight will not
 update. But of
 small value $\approx [0.88, 0.12]$

Most of the attention weight is assigned to the first key vector,
 very little to the second vector,

With Scaling

① Compute Scaled Dot Product

$$\left[\begin{bmatrix} 6, 4 \end{bmatrix} \right] \Rightarrow \text{Scale} \Rightarrow \left[\begin{bmatrix} 6/2, 4/2 \end{bmatrix} \right] = \left[\begin{bmatrix} 3, 2 \end{bmatrix} \right]$$

$$\sqrt{d_K} = \sqrt{4} = \frac{2}{\sqrt{2}}$$

See Research
 Variance of
 dot product
 should be same
 propagation of
 difference
 $2, 3$ Dimension
 $\sqrt{d_K}$ same

Dimension \uparrow Variance \uparrow
 $\sqrt{d_K} \uparrow$ same

$$\text{Softmax} \left(\begin{bmatrix} 3, 2 \end{bmatrix} \right) = \left[\frac{e^3}{e^3 + e^2}, \frac{e^2}{e^3 + e^2} \right] = \left[\frac{e^3}{e^3(1 + e^{-1})}, \frac{e^2}{e^2(e^1 + 1)} \right] = \left[0.73, 0.27 \right] \Rightarrow \text{Attention Weights}$$

Here, the attention weights are more balanced compared to the unscaled case.

Summary of Importance

Stabilizing Training: Scaling prevents extremely large dot products, which helps in stabilizing the gradients during backpropagation, making the training process more stable and efficient.

Preventing Saturation: By scaling the dot products, the softmax function produces more balanced attention weights, preventing the model from focusing too heavily on a single token and ignoring others.

Improved Learning: Balanced attention weights enable the model to learn better representations by considering multiple relevant tokens in the sequence, leading to better performance on tasks that require context understanding.

Scaling ensures that the dot products are kept within a range that allows the softmax function to operate effectively, providing a more balanced distribution of attention weights and improving the overall learning process of the model.

0 Scaling = $\sqrt{dk} = \sqrt{4} \Rightarrow 2$

Scaling will be done for all other terms

Scaled-Score $(Q_{The}, K_{The}) = \frac{1}{2} = 1$

Scaled-Score $(Q_{The}, K_{CAT}) = 0/2 = 0$

Scaled-Score $(Q_{The}, K_{SAT}) = 2/2 = 1$

1 Apply Softmax

Contextual Embedding Vector

ATTENTION WEIGHTS $_{The} = \text{softmax}([1, 0, 1]) = [0.4223, 0.1574, 0.4223]$

ATTENTION WEIGHTS $_{CAT} = \text{softmax}([0, 1, 2]) = [0.1574, 0.4223, 0.4223]$

ATTENTION WEIGHTS $_{SAT} = \text{softmax}([2, 2, 1]) = [0.2119, 0.2119, 0.5762]$

2 Weight Sum of Values

We multiply the attention weights by corresponding value vector

For the Token The =

$$\text{Output}_{(\text{The})} = 0.4223 * V_{\text{The}} + 0.1554 * V_{\text{CAT}} + 0.4223 * V_{\text{sat}}$$

$$= 0.4223 [1 0 1 0] + 0.1554 [0 1 0 1] + 0.4223 [1 1 1]$$

$$= [0.4223, 0, 0.4223, 0] + [0, 0.1554, 0, 0.1554] + [0.4223, 0.4223, 0.4223, 0.4223]$$

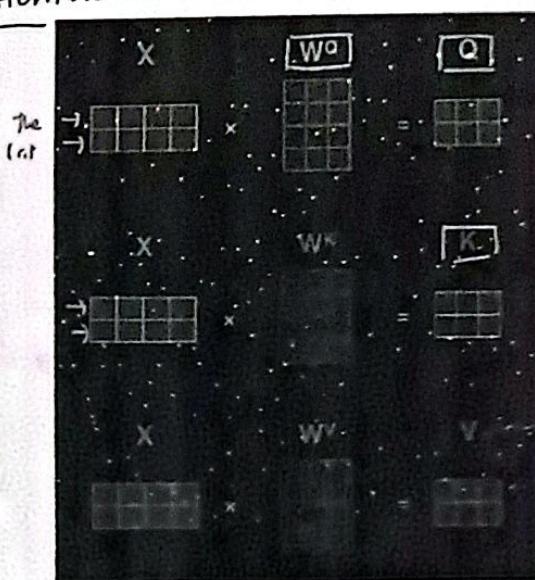
$$= [1.2669, 0.9999, 1.2669, 0.9999]$$

The $\boxed{1 \ 1 \ 0 \ 1 \ 1 \ 0}$ \Rightarrow Self Attention $\Rightarrow [1.2669, 0.9999, 1.2669, 0.9999]$

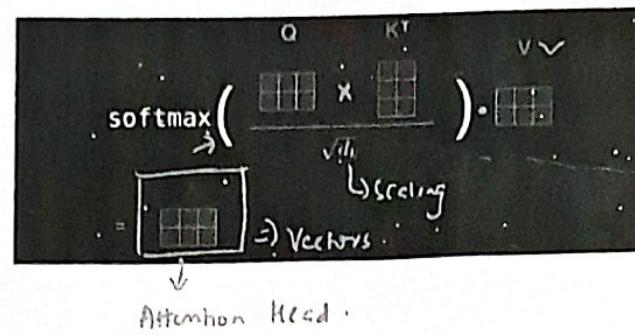
↓
Context vector

① $\hookrightarrow Q, K, V [w^Q, w^K, w^V]$
② \hookrightarrow Attention Score \hookrightarrow there are trained weight
③ \hookrightarrow Scaled
④ \hookrightarrow Softmax
⑤ \hookrightarrow Weighted sum of value (Softmax \times V)

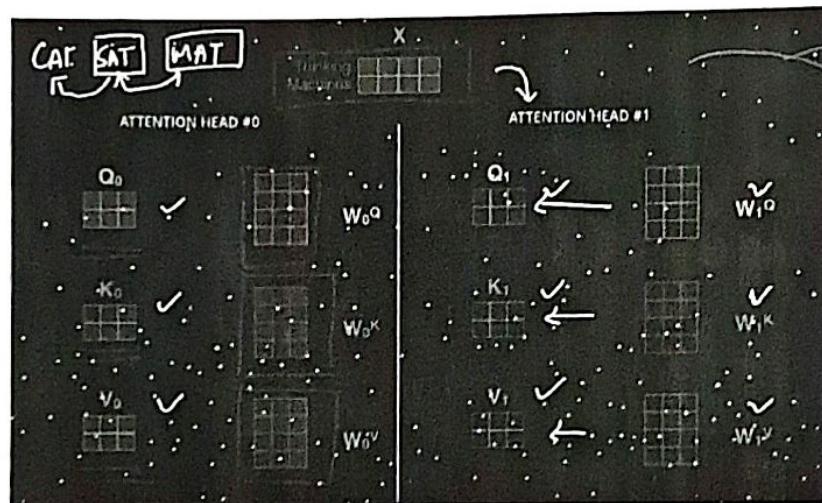
④ Multi Head Attention



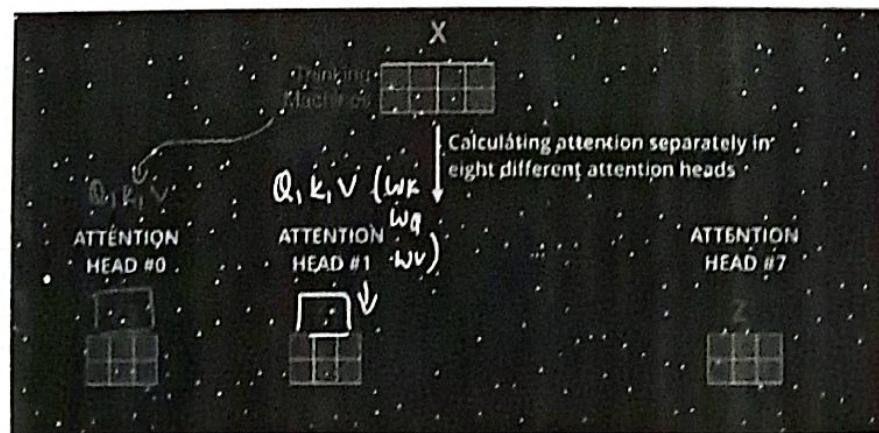
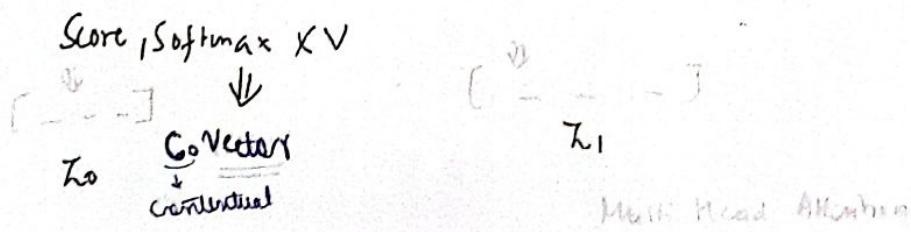
for every word we
get separate Q, K, V
matrix



→ Self Attention with Multi Heads



It expands Model's ability to focus on different portion of token

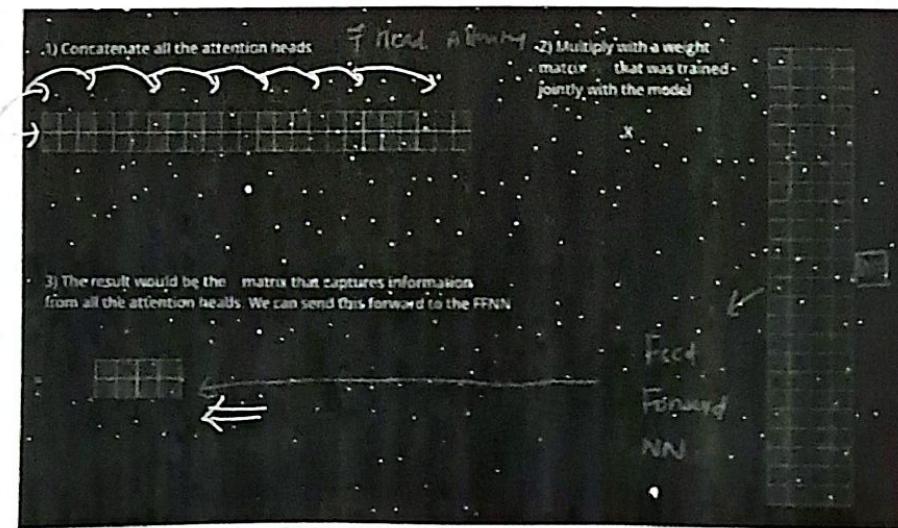
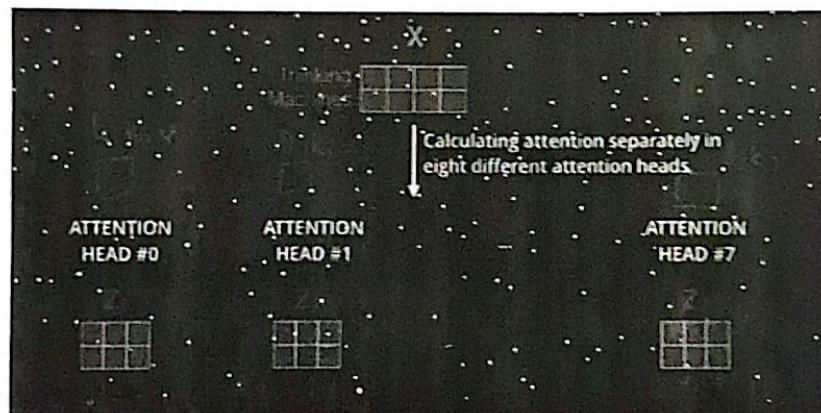
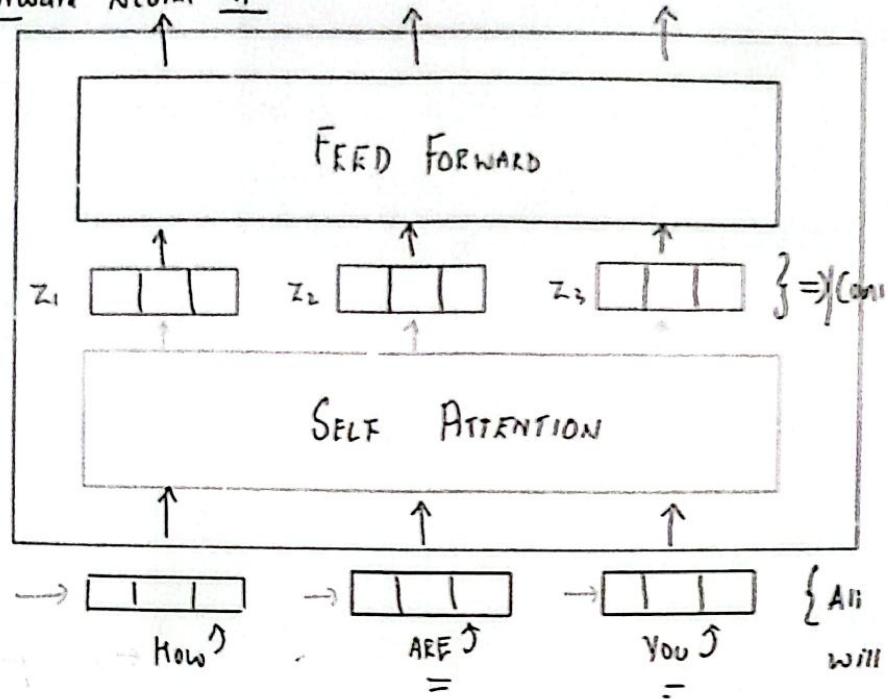


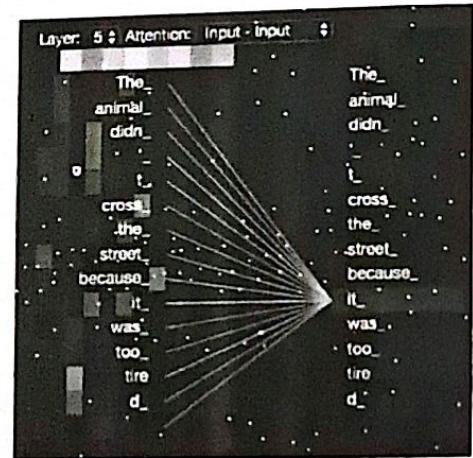
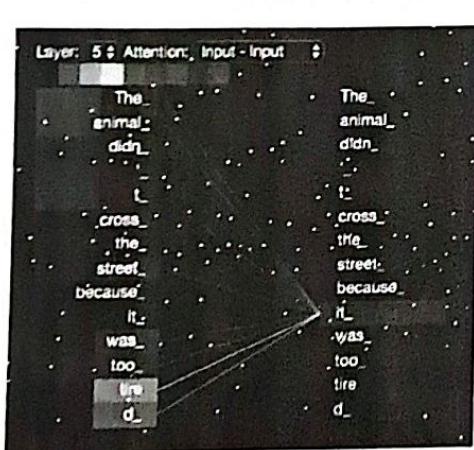
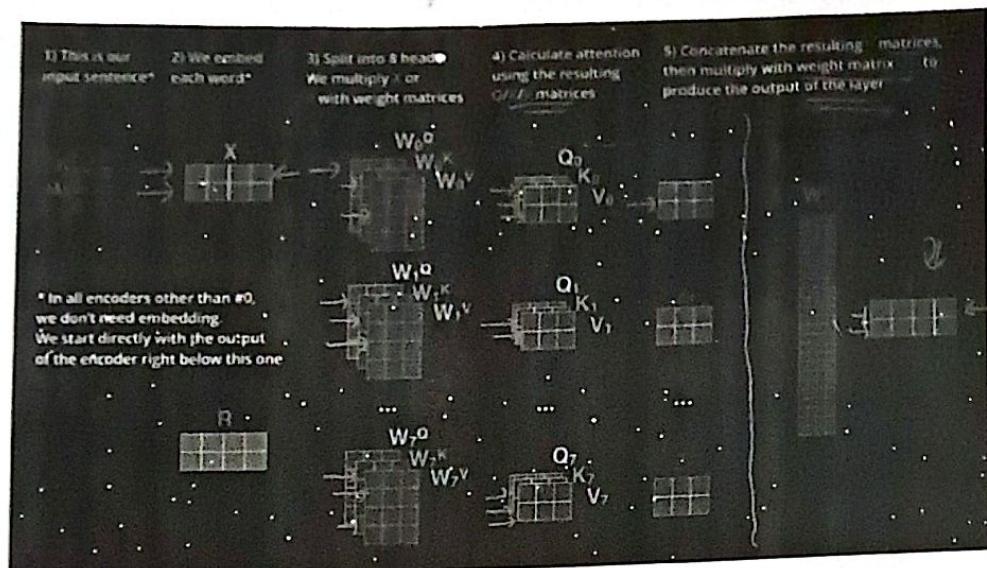
z_0

z_1

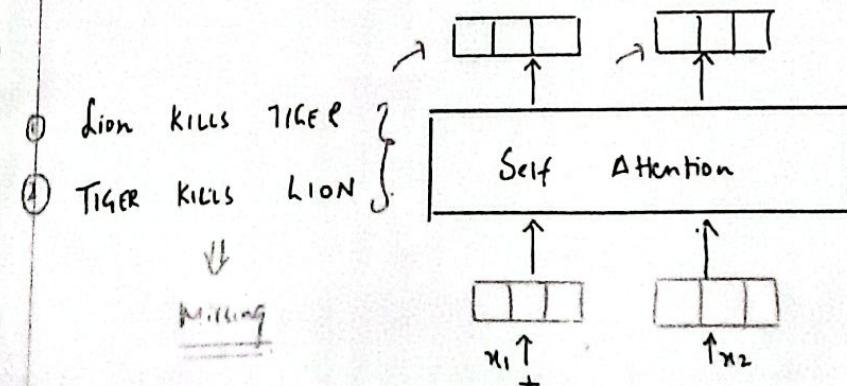
z_2

④ Feed Forward Neural NW





③ Positional Encoding - Representing Order of Sequence



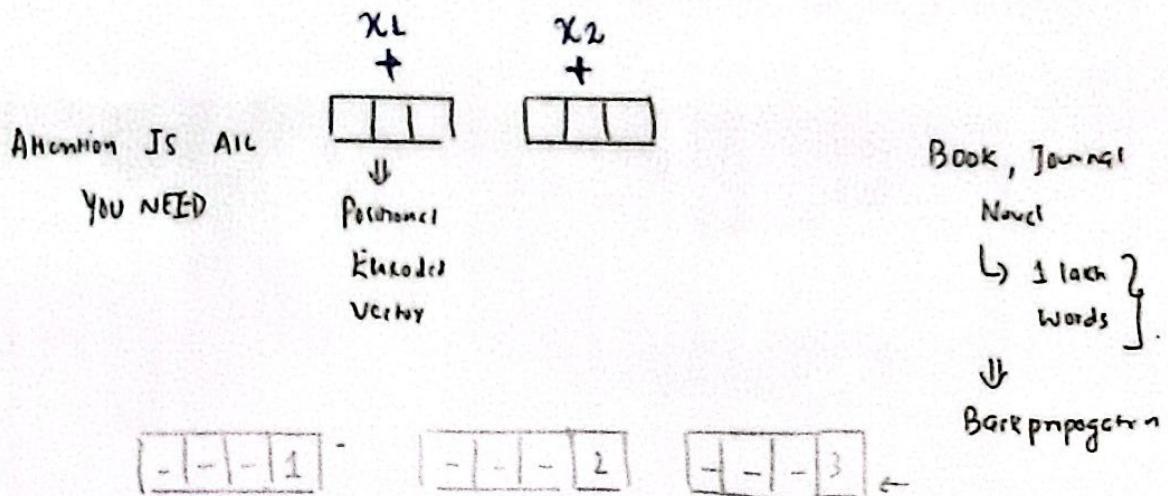
Advantages

- ① Word Tokens it can process parallelly



DRAWBACK

Lack the sequential structure of the words {order}

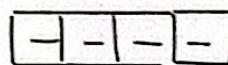


Types of Position Encoding

- 1) Sinusoidal Position Encoding → In Research paper
- 2) Learned Positional Encoding ⇒ Positional Encoding Are learned during Training

① Sinusoidal Positional Encoding = It uses sine and cosine functions of different frequencies to create positional encodings

Formula



$$P.E(pos, 2i) = \sin\left(\frac{pos}{10000^{\frac{2i}{d_{model}}}}\right)$$

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

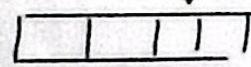
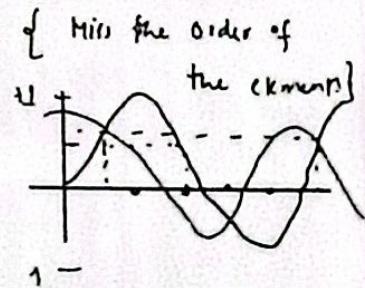
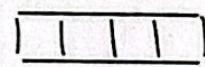
Where $pos \rightarrow$ is the position
 $i \rightarrow$ is the dimension
 $d_{model} \rightarrow$ is the dimensionality
 of the embeddings.

Eg: The Cat Sat

The → $[0.1 \ 0.2 \ 0.3 \ 0.4]$

CAT → $[0.5 \ 0.6 \ 0.7 \ 0.8]$

SAT → $[0.9 \ 1.0 \ 1.1 \ 1.2]$



sin

Calculate
in group
of two

cos

→ one sin & one cos

$$PE_{(pos, 2i)} = \sin\left(\frac{pos}{10000^{2i}/\text{dmodel}}\right) \quad PE_{(pos, 2i+1)} = \cos\left(\frac{pos}{10000^{2i+1}/\text{dmodel}}\right)$$

For our example $d_{\text{model}} = 4$

For position $pos = 0$ → The Ke key

$$P.E(0) = \sin\left(\frac{0}{100000/4}\right) = \sin(0) = 0$$

$$P(E(0,1)) = \cos\left(\frac{\pi}{2}\right) = \cos(0) = 1$$

$$P.E(0,2) = \sin \left(\frac{0}{1.0000} \right) = \sin(0) = 0$$

$$P.E(0,3) = \frac{1}{3} (0) = 1$$

$$P \cdot F = [0, 1, 0, 1]$$

$$P.E_{(pos, 2i)} = \sin \left(\frac{pos}{(1+20)^{2i/dimsize}} \right)$$

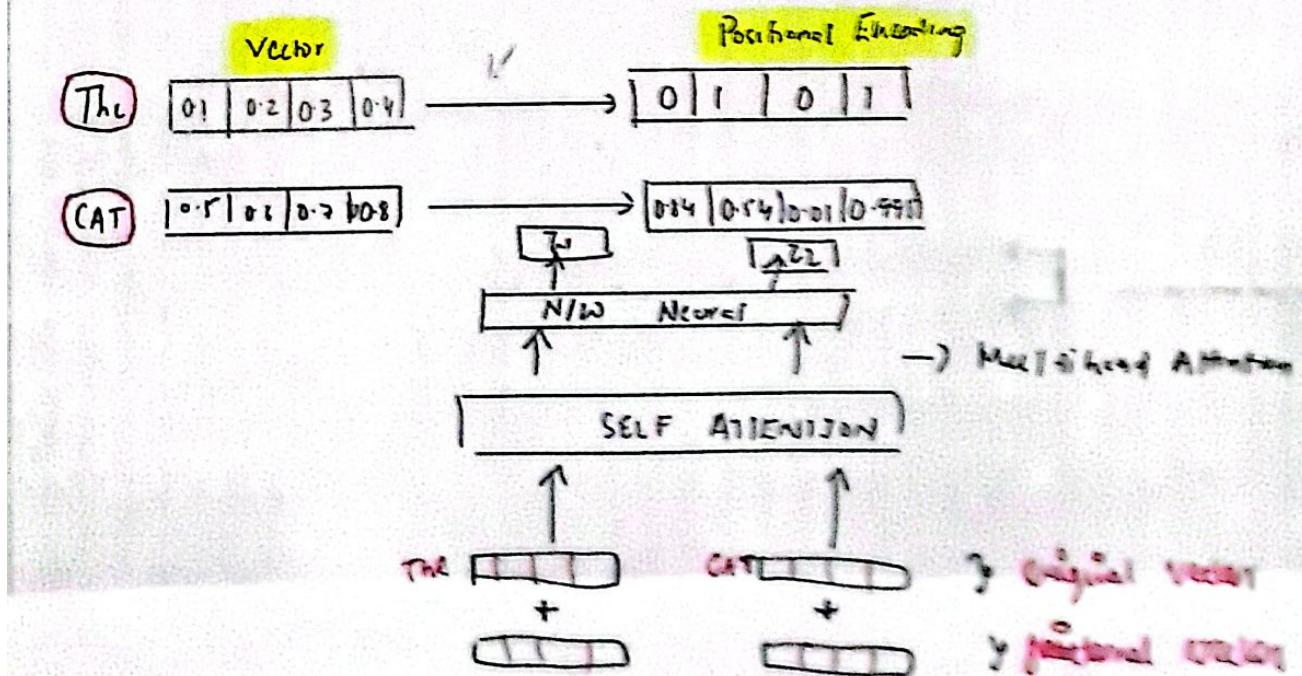
For pos = 1 ~~at~~ ~~ke~~ ~~line~~

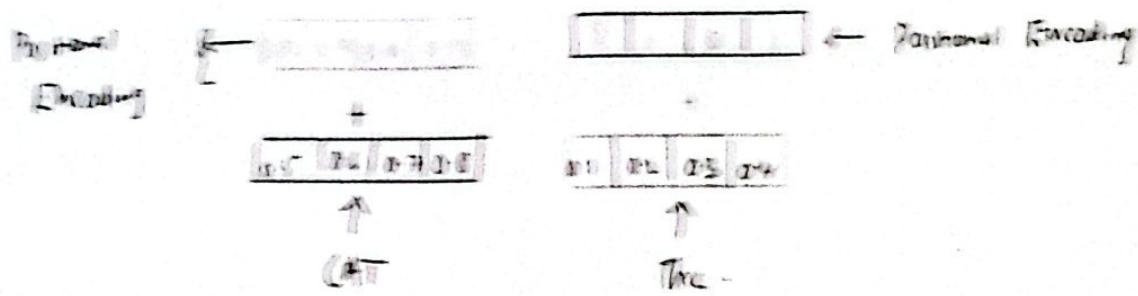
$$P.E(1,0) = \sin\left(\frac{1}{10000}\right) = \sin(1) = 0.8415$$

$$PE(1,1) = \cos\left(\frac{1}{10000^{2/4}}\right) \approx 0.5403$$

$$PE(1,2) = \sin\left(\frac{1}{100000000}\right) \approx 0.01$$

$$PE(1,3) = \cos \approx 0.99995$$

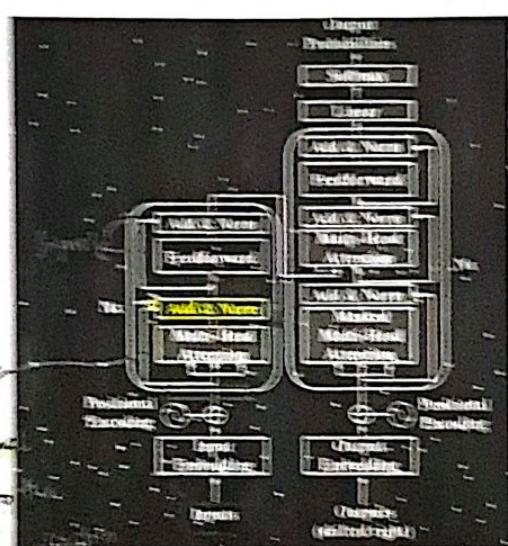




Multi-Head Attention Layer Implementation

Transformations

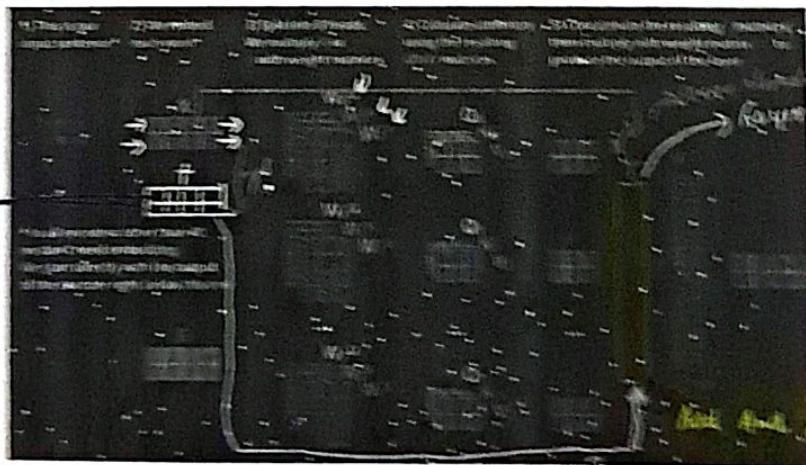
Residuals



- ① Self Attention layer
- ② Multi-head Attention
- ③ Positional Encoding
- ④ Layer Normalization

Add & Norm

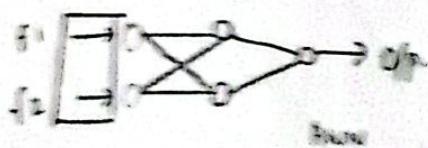
Finalized output



Additional Information

Normalization

Batch Normalization
Layer Normalization



Model Type	Model Name	Perplexity
Transformer	GPT-2	15.3
Transformer	GPT-3	1.3

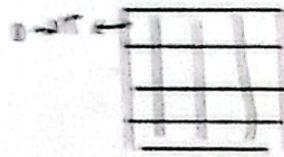
Normalization, Standard Scaling

$$z_{\text{norm}} = \frac{x - \mu}{\sigma}$$

$$\{A \in \mathbb{R}, \Gamma = 1\}$$

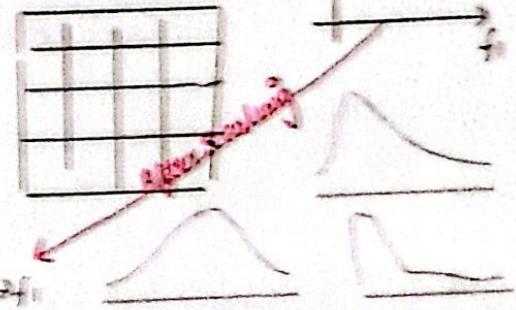
$$f_1 \rightarrow f_1' \leftarrow \mu \text{ or } \sigma$$

One scaling for input \Rightarrow We can learn



\Rightarrow Mean Max (Scale) \Rightarrow

$$\begin{matrix} z_{\text{norm}} \\ \text{Current} \end{matrix} \xrightarrow{\mu \text{ or } \sigma}$$



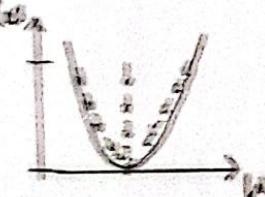
Advantages

① Improved Training Stability

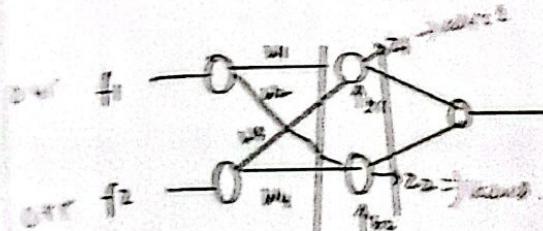


Vanishing and exploding gradient problem

② Faster Convergence



\Rightarrow Back propagation \Rightarrow Stable update



f_1	f_2	Number of layers
Mean	Mean	
0.45	0.45	1
0.60	0.20	
-	-	

$$z_1 = \Gamma \left[(0.45 + b_1 + 0.45 + b_2) + b_3 \right] = \text{Value 1} \quad \text{Due to distribution change}$$

$$z_2 = \Gamma (\text{Value 1}) = \text{Value 2} \quad \frac{\{A \in \mathbb{R}, \Gamma = 1\}}{A, \Gamma} \quad f_1, f_2$$

18,300 after back propagation of calculation of z_1, z_2
are distribution with change
60,40, again apply
Normalisation on it.

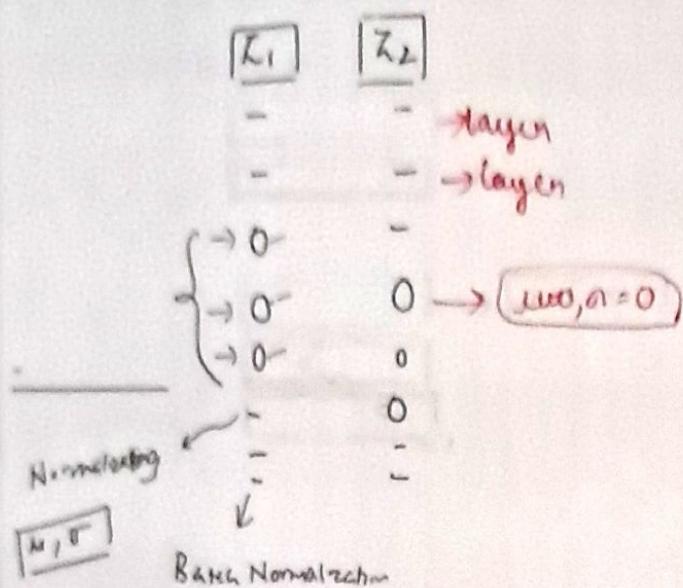
Batch Normalization vs Layer Normalization

f_1	f_2	z_1	z_2
-	-	-	-
-	-	-	-
-	-	-	-

BATCH Normalization

layer
Normalisation

Gramma
Beta
 $\gamma, \beta \rightarrow$ Learnable parameters



Normalization



γ, β

$$z_1 = \sigma [w^T x + b_1]$$

$$y = \gamma \left[\frac{z_1 - \mu_1}{\sigma_1} \right] + \beta$$

Scale And Shift parameters

① Example to understand this Theory normalized.

$$① \text{ "CAT" } = [2.0, 4.0, 6.0, 8.0] \leftarrow \{ \text{Vectors} \}$$

$$② \text{ Parameters } = \gamma = [1.0, 1.0, 1.0, 1.0] \rightarrow \text{Learned scale } \quad \left. \begin{array}{l} \gamma \\ \beta \end{array} \right\} \Rightarrow \text{Scale And} \\ \beta = [0.0, 0.0, 0.0, 0.0] \rightarrow \text{shift} \quad \left. \begin{array}{l} \beta \end{array} \right\} \text{Shift param}$$

① Compute the mean

$$z_{score} = \frac{x_i - \mu}{\sigma}$$

$$\begin{aligned} \mu &= \frac{1}{4} (2.0 + 4.0 + 6.0 + 8.0) \\ &= \frac{20.0}{4} = 5.0 \end{aligned}$$

② Compute the variance (σ^2)

$$\sigma^2 = \frac{1}{4} \left[(2.0 - 5.0)^2 + (4.0 - 5.0)^2 + (6.0 - 5.0)^2 + (8.0 - 5.0)^2 \right] = 5.0$$

③ Normalize the i/p

$$\hat{x}_i = \frac{x_i - \mu}{\sqrt{\sigma^2 + \epsilon}} \quad (\epsilon = 1e^{-5} \Rightarrow \text{Avoid division by 0})$$

$$\sqrt{\sigma^2 + \epsilon} = \sqrt{5.0 + 1e^{-5}} \approx \sqrt{5.00001} = 2.236$$

$$\begin{aligned}
 \hat{x}_1 &= \frac{2.0-5.0}{2.236} \approx -1.34 & \text{Normalized vector} \\
 \hat{x}_2 &= \frac{4.0-5.0}{2.236} \approx -0.45 & \hat{x} = \left[\begin{array}{c} -1.34, -0.45, 0.45, 1.34 \end{array} \right] \\
 \hat{x}_3 &= \frac{6.0-5.0}{2.236} \approx 0.45 \\
 \hat{x}_4 &= \frac{8.0-5.0}{2.236} \approx 1.34
 \end{aligned}$$

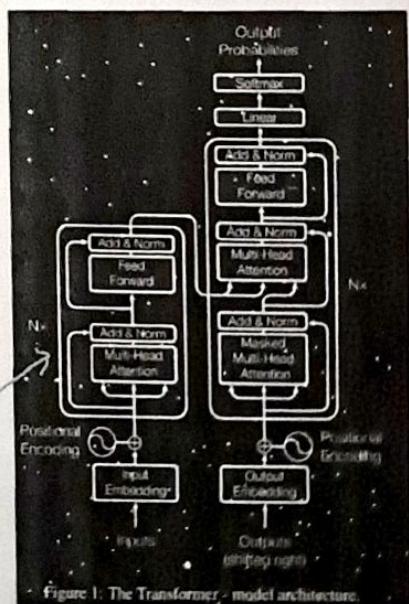
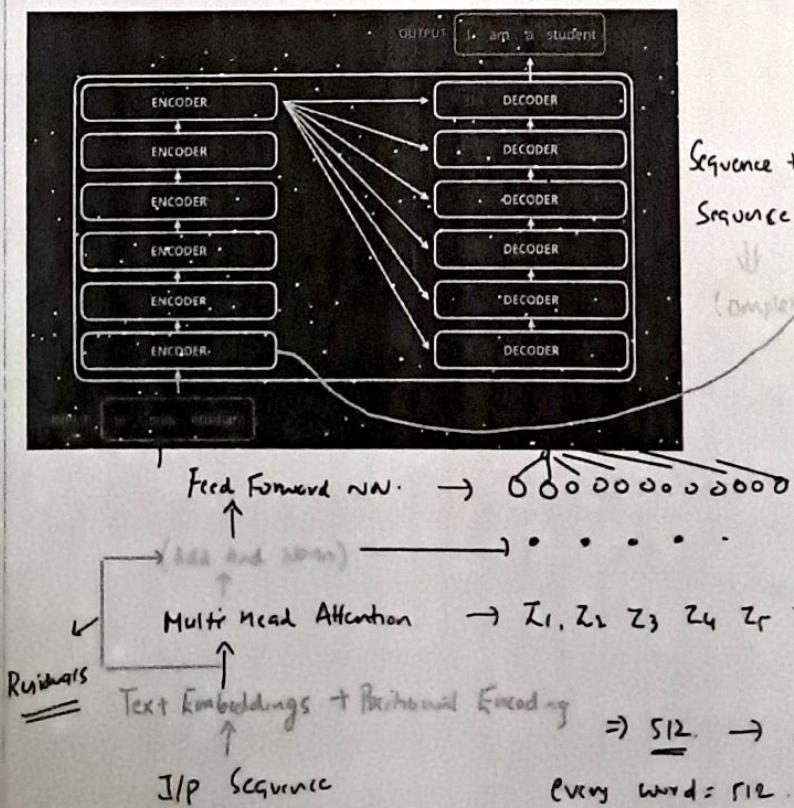
④ Scale And Shift

$$y_i = \gamma_i \hat{x}_i + \beta_i$$

$$\gamma = [1.0, 1.0, 1.0, 1.0] \quad \beta = [0.0, 0.0, 0.0, 0.0]$$

$$y = [-1.34, -0.45, 0.45, 1.34]$$

⑤ Encoder Architecture [Research Paper]



$$d_h = 64$$

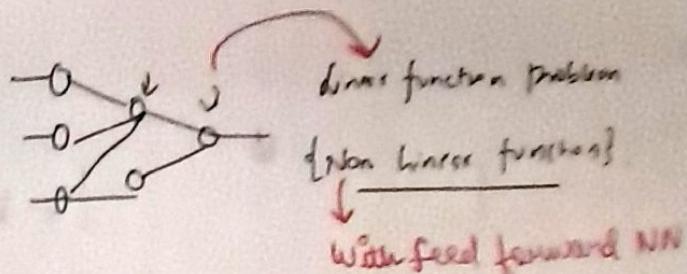
$$K = 64$$

$$V = 64$$

$$\sqrt{64} = 8$$

flow

① Feed Forward NN Use



① Adding Non-linearity

② Processing Each position Independently.

Self Attention → (capture relationships b/w one another)

FFN → Each token representation Independently



Transforming these representation further

and allows the model to learn

Richer Representation

FFN =>

③ FFN → Decoder => Adds Depth to the Model

Depth ↑ => More Learnings → DATA

④ Decoders In Transformers

3 main Components

The transformer decoder is responsible for generating the output sequence one token at a time, using the encoder's output and the previously generated tokens.

- ① Masked Multi Head Self Attention ✓
- ② Multi Head Attention (Encoder Decoder Attention)
- ③ Feed Forward Neural Network.

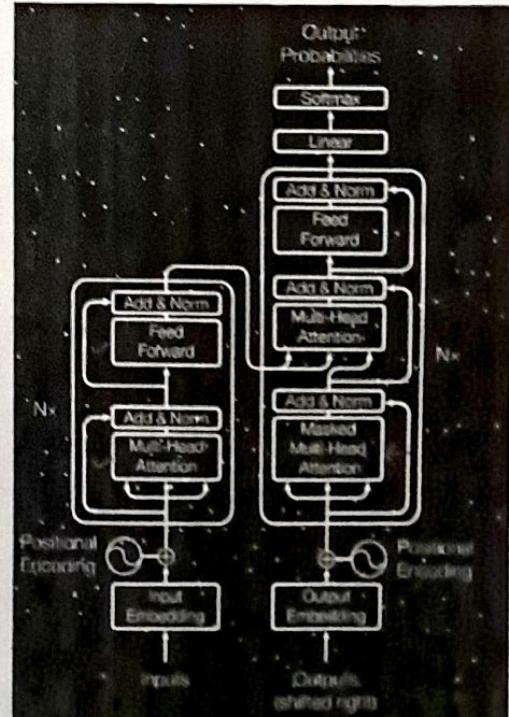
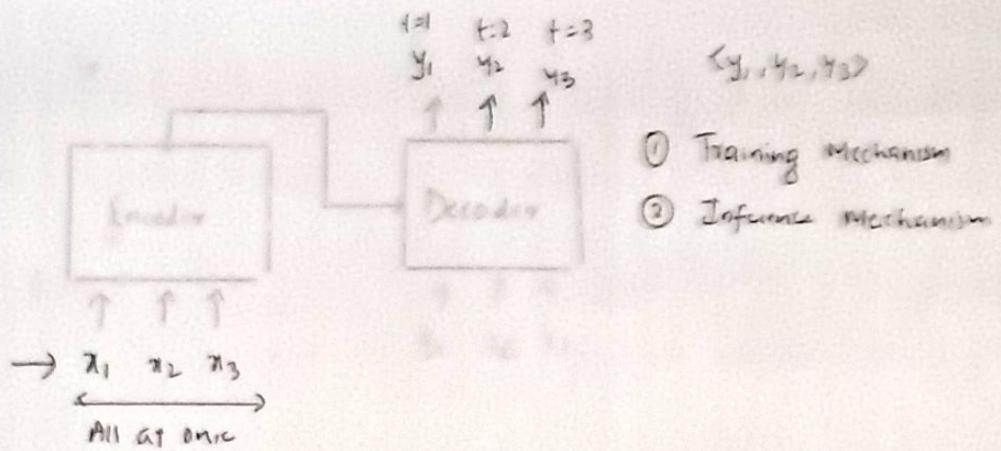
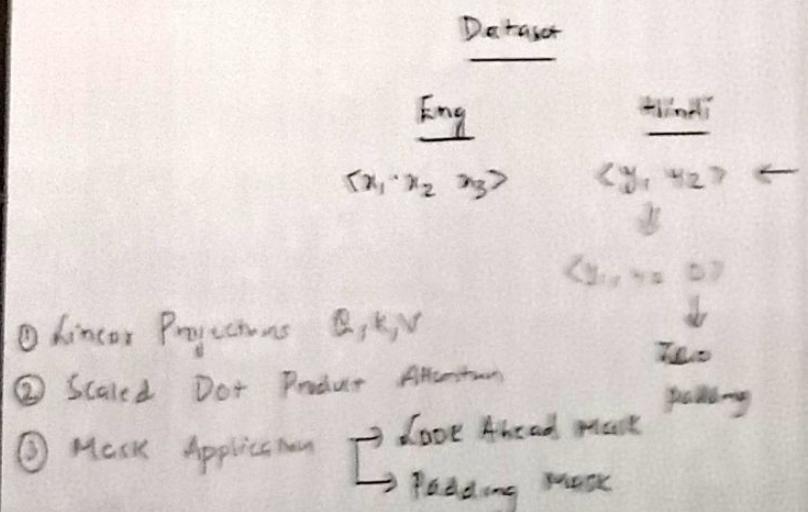
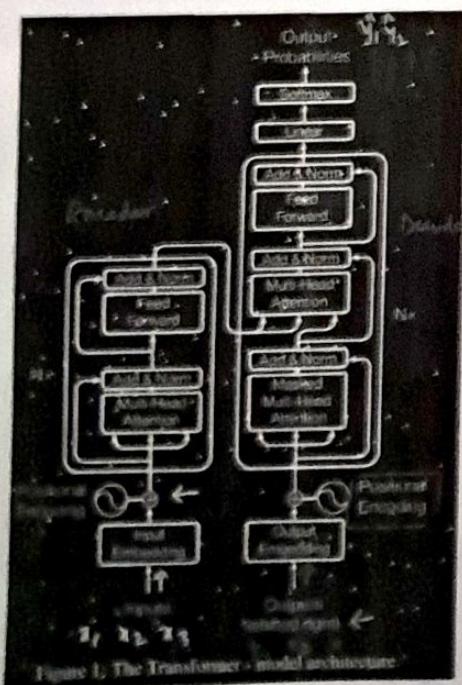


Figure 1: The Transformer - model architecture.



① Masked Multi Head Self Attention

- ① I/P Embedding And Positional Embedding \rightarrow Two padding \rightarrow Sequence Length Equal
- ② Linear Projection for Q, K, V
- ③ Scaled Dot Product Attention
- ④ Mask Application \Rightarrow Try to understand the imp.
- ⑤ Multi Head Attention
- ⑥ Concatenation And Final Linear Projection
- ⑦ Residual Connection And Layer Normalization



Masked

Multi-Head
Attention

$$\begin{bmatrix} 1 & 2 & 3 & 0 \end{bmatrix} \quad \begin{bmatrix} 1 & 2 & 3 & 0 \end{bmatrix}$$

4 dimension vector

3) Input Embedding and Positional Encoding

Output Embedding [STEP 1]

$$\begin{bmatrix} [0.1, 0.2, 0.3, 0.4], \\ [0.5, 0.6, 0.7, 0.8], \\ [0.9, 1.0, 1.1, 1.2], \\ [0.0, 0.0, 0.0, 0.0] \end{bmatrix}$$

Step 2: Linear Projection for Q, K, and V

$$WQ = WK = WV = I$$

Create query (Q), Key(K) and value(V) vectors

$$Q = \text{Output Embedding} + W_Q = \text{Output Embedding}$$

$$K = " + WK = " "$$

$$V = " + WV = " "$$

$$Q = K = V = \begin{bmatrix} [0.1, 0.2, 0.3, 0.4], \\ [0.5, 0.6, 0.7, 0.8], \\ [0.9, 1.0, 1.1, 1.2], \\ [0.0, 0.0, 0.0, 0.0] \end{bmatrix}, \begin{bmatrix} 0.1 \\ 0.2 \\ 0.3 \\ 0.4 \end{bmatrix} K^T \begin{bmatrix} 0.5 \\ 0.6 \\ 0.7 \\ 0.8 \end{bmatrix} \begin{bmatrix} 0.9 \\ 1.0 \\ 1.1 \\ 1.2 \end{bmatrix} \begin{bmatrix} 0 \\ 0 \\ 0 \\ 0 \end{bmatrix}$$

③ Scaled Dot Product Attention Calculation

$$\text{Scores} = Q * K^T / \sqrt{dk}$$

$$= Q * K^T / 2$$

$$\begin{aligned} & 0.1 * 0.1 + 0.2 * 0.2 + 0.3 * 0.3 + 0.4 * 0.4, \\ & 0.1 * 0.5 + 0.2 * 0.6 + 0.3 * 0.7 + 0.4 * 0.8, \\ & 0.1 * 0.9 + 0.2 * 1.0 + 0.3 * 1.1 + 0.4 * 1.2 \\ & 0.1 * 0 + 0.2 * 0 + 0.3 * 0 + 0.4 * 0 \end{aligned}$$

$$\text{Scores} : \left[\begin{bmatrix} 0.3, 0.7, 1.1, 0.0 \end{bmatrix} \quad \left[\begin{bmatrix} 0.7, 1.9, 3.1, 0.0 \end{bmatrix} \quad \left[\begin{bmatrix} 1.1, 3.1, 5.1, 0.0 \end{bmatrix} \quad \left[\begin{bmatrix} 0.0, 0.0, 0.0, 0.0 \end{bmatrix} \right] \right] \right]$$

④ Masked Application

It helps manage the structure of the sequences being processed and ensures the model behaves correctly during training and inference.

Reasons

① Handling Variable length Sequences with Padding MASK

Purpose

- ① To handle sequences of different length in batch
- ② To ensure that padding tokens, which are added to make sequences of uniform length, do not affect the model prediction

Eg: $\text{inp} \leftarrow \text{Sequence 1} \quad [1, 2, 3]$ Op
 $\rightarrow y_1, y_2, 0.0000$

$\text{inp} \leftarrow \text{Sequence 2} \quad [4, 5, \boxed{0}]$ 0 is the padding token

$\uparrow \rightarrow$ Influence the Attention Mechanism

\Downarrow

lead to Incorrect or biased predictions.

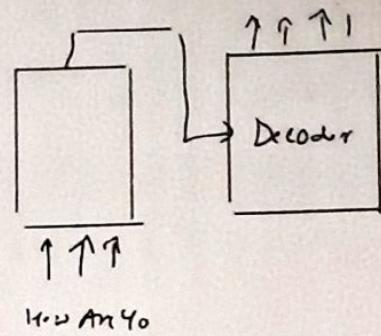
A padding mask \Rightarrow The tokens are ignored.

Masking \rightarrow Padding Mask } Padding Mask $\begin{bmatrix} 1 & 1 & 1 \end{bmatrix}$

Look Ahead Mask . } $\begin{bmatrix} 1 & 1 & 0 \end{bmatrix}$

Mask \rightarrow Maintain Auto Regressive Property

+ each position in the decoder
we can only attend to
position, but no future position



1 1 1
0 0 1

Language Modelling, Translation

$\begin{bmatrix} 1, 0 \end{bmatrix} \rightarrow \begin{bmatrix} 1, 1, 0 \end{bmatrix} \rightarrow$ 1D Mask

Token 1 attends
to token 1, 2
1, 2

$$\begin{bmatrix} 1 & 1 & 0 \\ 1 & 1 & 0 \\ 0 & 0 & 0 \end{bmatrix}$$

Convert 1D to 2D Mask

for each token in the sequence
the mask should indicate
which tokens it can attend
to

Mask \rightarrow Decoder Output

$\begin{bmatrix} 0 \\ 0 \\ 1 \end{bmatrix}$

Padding And Looking Ahead Mask

use multiplication of 2 mask

$\begin{bmatrix} [1, 0, 0] \\ [1, 1, 0] \\ [0, 0, 0] \end{bmatrix}$

④ MASK

$$\underline{\text{Scores}} : \begin{bmatrix} [0.3, 0.7, 1.1, 0.0] \\ [0.7, 1.9, 3.1, 0.0] \\ [1.1, 3.1, 5.1, 0.0] \\ [0.0, 0.0, 0.0, 0.0] \end{bmatrix}$$

Look Ahead Mask

$$\begin{bmatrix} [1, 0, 0, 0] \\ [1, 1, 0, 0] \\ [1, 1, 1, 0] \\ [1, 1, 1, 1] \end{bmatrix}$$

Padding Masking [extended to 2D Format]

$$\begin{bmatrix} [1, 1, 1, 0] \\ [1, 1, 1, 0] \\ [1, 1, 1, 0] \\ [0, 0, 0, 0] \end{bmatrix}$$

Combined Mask = Look Ahead Mask + Padding Mask

$$\begin{bmatrix} [1*1, 0*1, 0*1, 0*0] \\ [1*1, 1*1, 0*1, 0*0] \\ [1*1, 1*1, 1*1, 0*0] \\ [1*1, 1*1, 1*1, 1*0] \end{bmatrix} = \begin{bmatrix} [1, 0, 0, 0] \\ [1, 1, 0, 0] \\ [1, 1, 1, 0] \\ [1, 1, 1, 0] \end{bmatrix} \Rightarrow \begin{bmatrix} [1, -\infty, -\infty, -\infty] \\ [1, 1, -\infty, -\infty] \\ [1, 1, 1, -\infty] \\ [1, 1, 1, -\infty] \end{bmatrix}$$

Masked Score

$$\begin{bmatrix} [0.3, -\infty, -\infty, -\infty] \\ [0.7, 1.9, -\infty, -\infty] \\ [1.1, 3.1, 5.1, -\infty] \\ [0.0, 0.0, 0.0, -\infty] \end{bmatrix}$$

Zero out the influence when the Softmax is applied.

Attention weights.

④ Softmax

$$\begin{aligned}\text{Softmax Score} &= \text{softmax}(\text{masked Scores}) \\ &= \left[\begin{bmatrix} 1.0, 0.0, 0.0, 0.0 \end{bmatrix}, \right. \\ &\quad \left. \begin{bmatrix} 0.3, 0.7, 0.0, 0.0 \end{bmatrix}, \right. \\ &\quad \left. \begin{bmatrix} 0.1, 0.3, 0.6, 0.0 \end{bmatrix}, \right. \\ &\quad \left. \begin{bmatrix} 1.0, 0.0, 0.0, 0.0 \end{bmatrix} \right]\end{aligned}$$

⑤ Weight Sum of Value

$$\text{Attention O/p} = \text{Softmax Scores} \cdot V$$

Masking

Masking in the transformer architecture is essential for several reasons. It helps manage the structure of the sequences being processed and ensures the model behaves correctly during training and inference. Here are the key reasons for using masking:

1. Handling Variable-Length Sequences with Padding Mask

Purpose

To handle sequences of different lengths in a batch.

To ensure that padding tokens, which are added to make sequences of uniform length, do not affect the model's predictions.

2. Maintaining Autoregressive Property with Look-Ahead Mask

Purpose

To ensure that each position in the decoder output sequence can only attend to previous positions and itself, but not future positions.

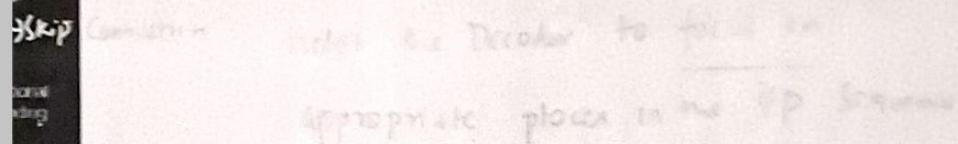
This is crucial for sequence generation tasks like language modeling and translation, where the model should not have access to future tokens when predicting the current token.

Multi-Head Attention

- Encoder O/p \rightarrow Set of Attention vector K & V
 - Model's Matrixt \rightarrow Attention vector Q {Query Vector}

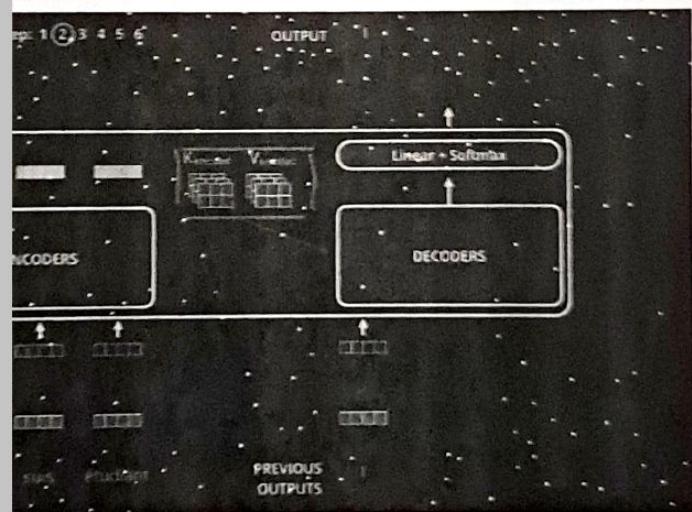
These are to be used by each decoder in its

"Encoder-decoder" Attention layer

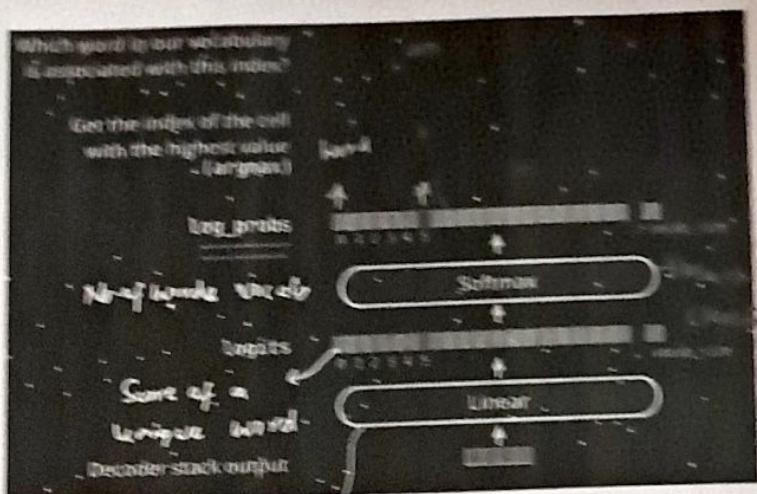
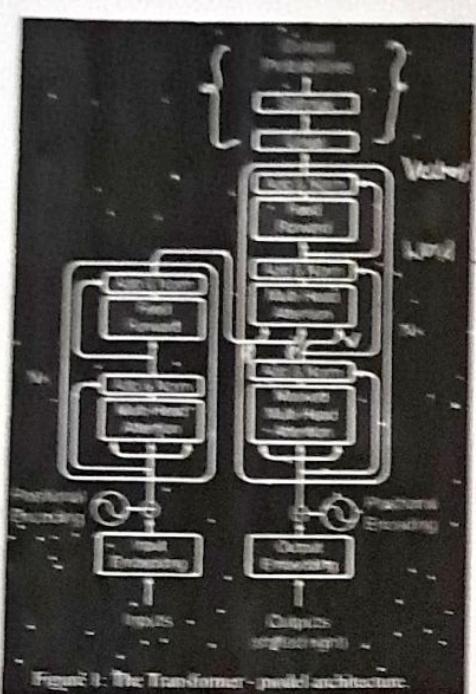


Dakar

$$\begin{array}{cc} i/p & o/p \\ \langle x_1, x_2, x_3 \rangle & \langle y_1, y_2, y_3 \rangle \end{array}$$



④ The Final Linear And Sigmoid Layer ($z \rightarrow \text{softmax}$)



Linear \Rightarrow The linear layer is a simple fully connected
matrix with their weights are also produced
by the stack of decoder \Rightarrow begin with
 \Rightarrow

Model = 10,000 \Rightarrow Vocabulary \Rightarrow logit vector = 10000 (0.15 wide)

③ Softmax layer turns those scores into probabilities (all add up to 1.0)

The cell with the highest probability is chosen, and the word associated with it is produced as the $o(p, i)$ time stamp.

Recap of Training

WORD	1	2	3	4	5
INDEX	0	1	2	3	4

Output Vocabulary					
WORD	a	am	i	thanks	student
INDEX	0	1	2	3	4
One-hot encoding of the word 'am'					
	0.0	1.0	0.0	0.0	0.0

Hence \rightarrow Theorems

Untrained Model Output					
0.2	0.2	0.1	0.2	0.2	0.1

2) $\sin \theta = \frac{4}{5}$.

Back Propagation

Output Vocabulary	am	is	thanks	student	very	
position #1	0.0	0.0	1.0	0.0	0.0	0.0
position #2	0.0	1.0	0.0	0.0	0.0	0.0
position #3	1.0	0.0	0.0	0.0	0.0	0.0
position #4	0.0	0.0	0.0	0.0	1.0	0.0
position #5	0.0	0.0	0.0	0.0	0.0	1.0

5 ← Loss

I am a Student

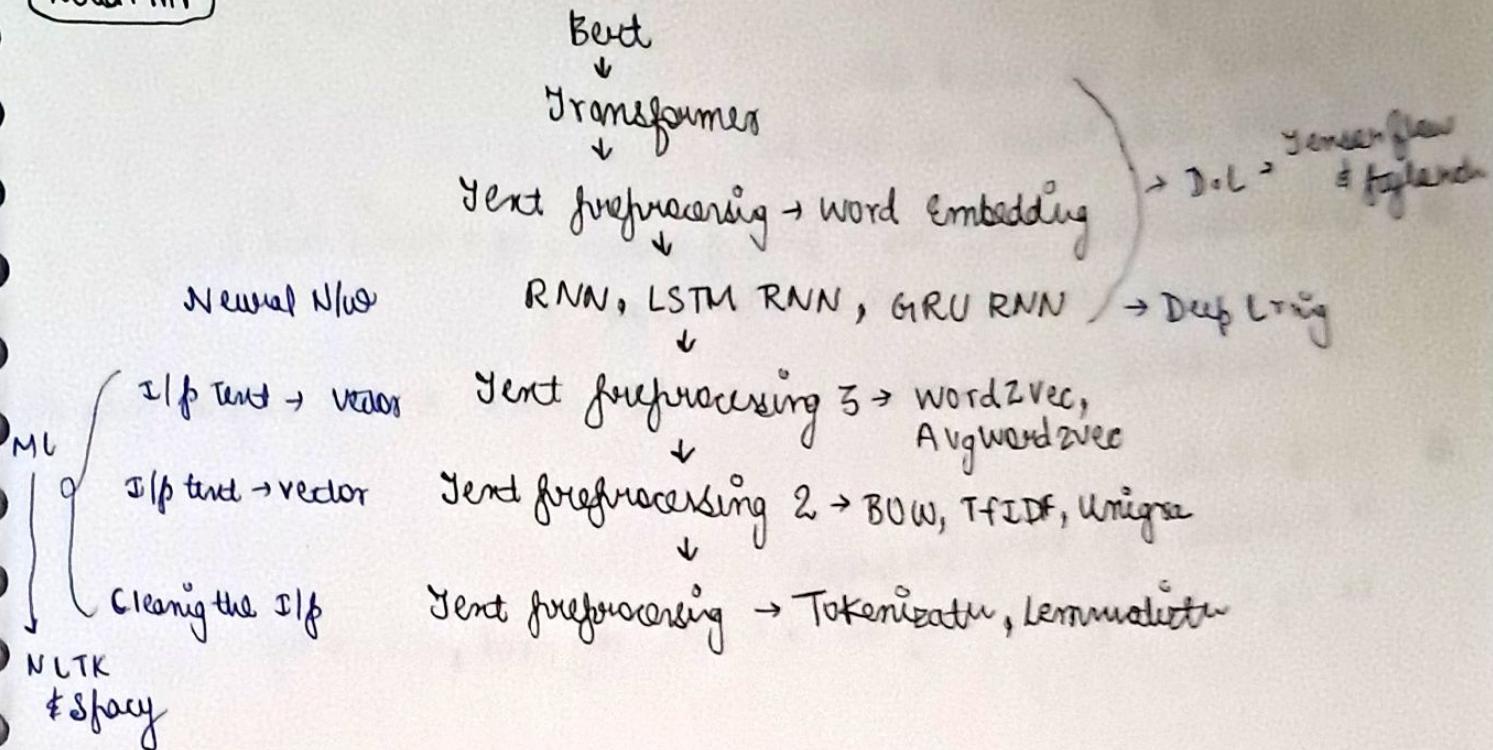
I am a student
↓ ↓ ↓ ↓
ONE THE ONE ONE

Trained Model Outputs						
Output Vocabulary	a	an	the	was	is	are
position #1	0.01	0.02	0.03	0.01	0.03	0.01
position #2	0.01	0.8	0.1	0.05	0.01	0.03
position #3	0.99	0.001	0.001	0.001	0.002	0.001
position #4	0.001	0.002	0.001	0.002	0.04	0.01
position #5	0.01	0.01	0.001	0.001	0.001	0.00

1

① NLP in Machine Learning

Roadmap



② Tokenization in NLP

- ① Corpus → Paragraph
- ② Documents → Sentences
- ③ Vocabulary → Unique word
- ④ words →

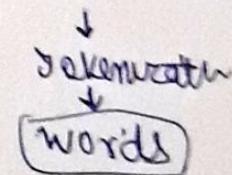
Tokenizer

"My name is Krish and I have a interest in teaching ML, NLP and DL. I am also a YouTuber"

↓

Tokens of sentences

- ① My name is - - - - - DL
- ② I am also a YouTuber



[I like to drink Apple Juice. My friend likes mango juice]

11 words

↓
Tokenization
Token (sentence)

- ① I like to drink Apple Juice
- ② My friend likes Mango juice
↓
if Remove

→ Unique words → 10 word
↓
Vocabulary
↓
45 remain
9 words

- NLP V/S Spacy in Vs Code.
- Stemming
 - ↳ Word to its word stem
 - ↳ root word known as lemma
- Lemmatization (POS Tag \rightarrow part of speech, eg \rightarrow noun ^N, verb ^V)
 - ↳ To solve problem come with Stemming
 - ↳ wordNet Lemmatizer class be check to the replace to the
- Stopwords
 - ↳ Available for many languages
 - ↳ eg \rightarrow In english \rightarrow They are \rightarrow i, me, my, and, are, why