

⇒ Transformers:

What and why transformers?

Transformers in 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.

↓
sequence-to-sequence task

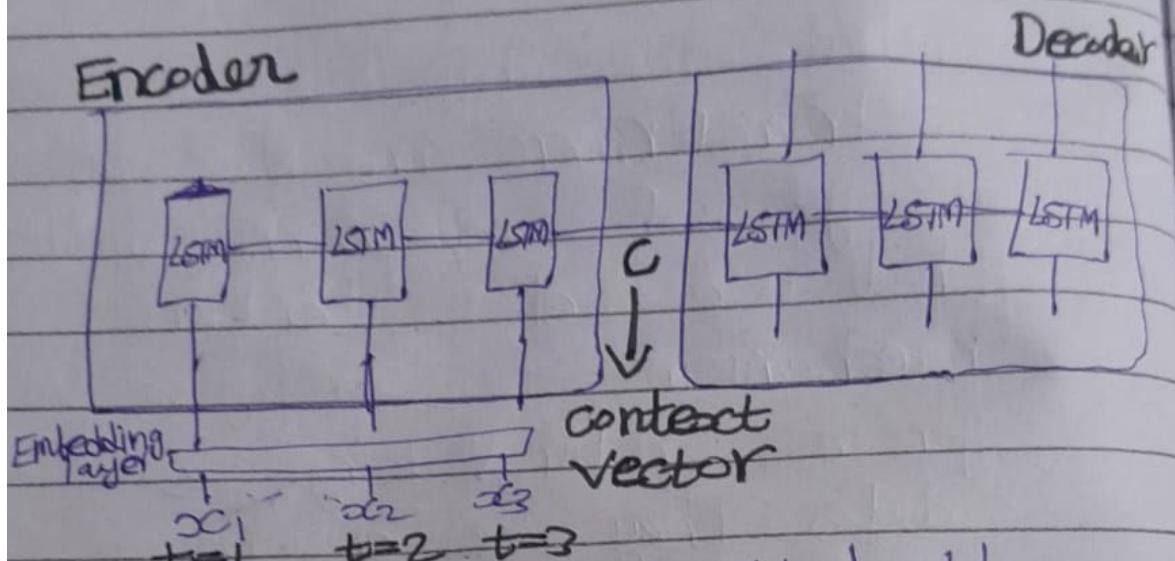
e.g.

Language Translation → Google
many-to-many = Translator

inputs → many outputs → many

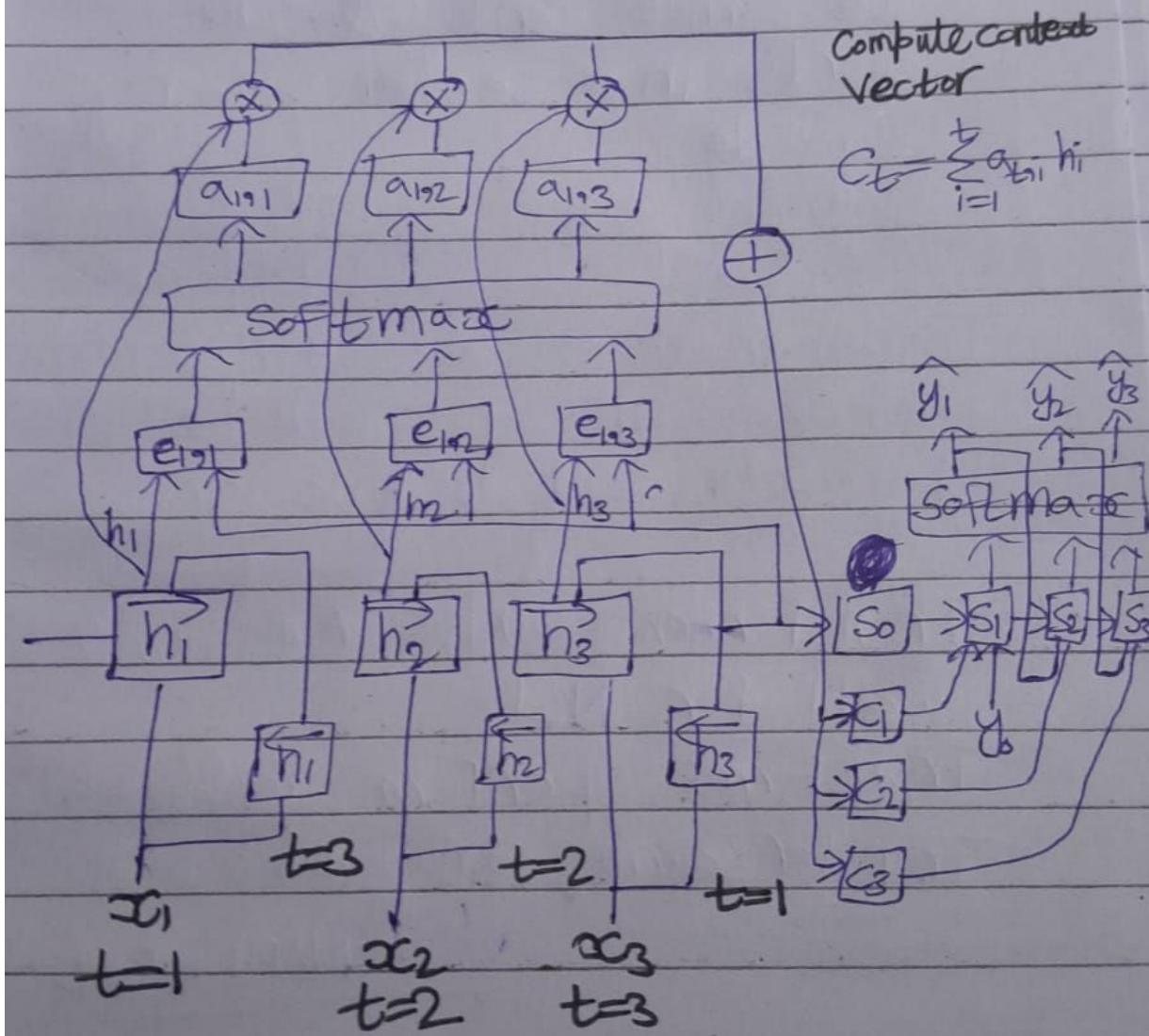
- Sequence of both inputs and outputs is important.
(seq2seq task)

We can also use Encoder-Decoder for this task



⇒ The issue was that the context vector passing has more influence of latest timestamp i.e $t=3$ and as sentence get longer the influence of starting time stamps i.e $t=1$ becomes negligible so context vector is not sufficient and as length of sentence increases the blue score (performance) decreases.

To resolve this issue we used
Attention Mechanism (passing
additional context with context vector)



Problem with Attention Mechanism

① We are sending single word to a single timestamp So, parallelly we cannot send all the words in a sentence So not scalable.

Dataset \rightarrow huge \rightarrow Not scalable w.r.t training

Transformers \rightarrow use Self Attention Module \rightarrow All words will be parallelly sent to Encoder

- Transformers works amazingly well for NLP.
- Transformer performs really well in Transfer Learning \rightarrow Multimodal task
 \Rightarrow NLP + Image

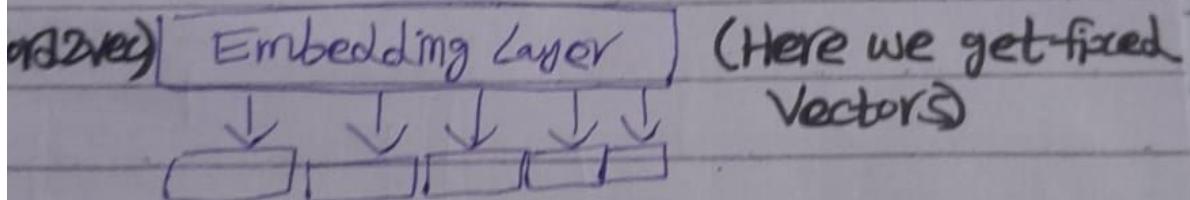
State of Art
↑

Transformers → AI Space → SOTA Models

BERT GPT

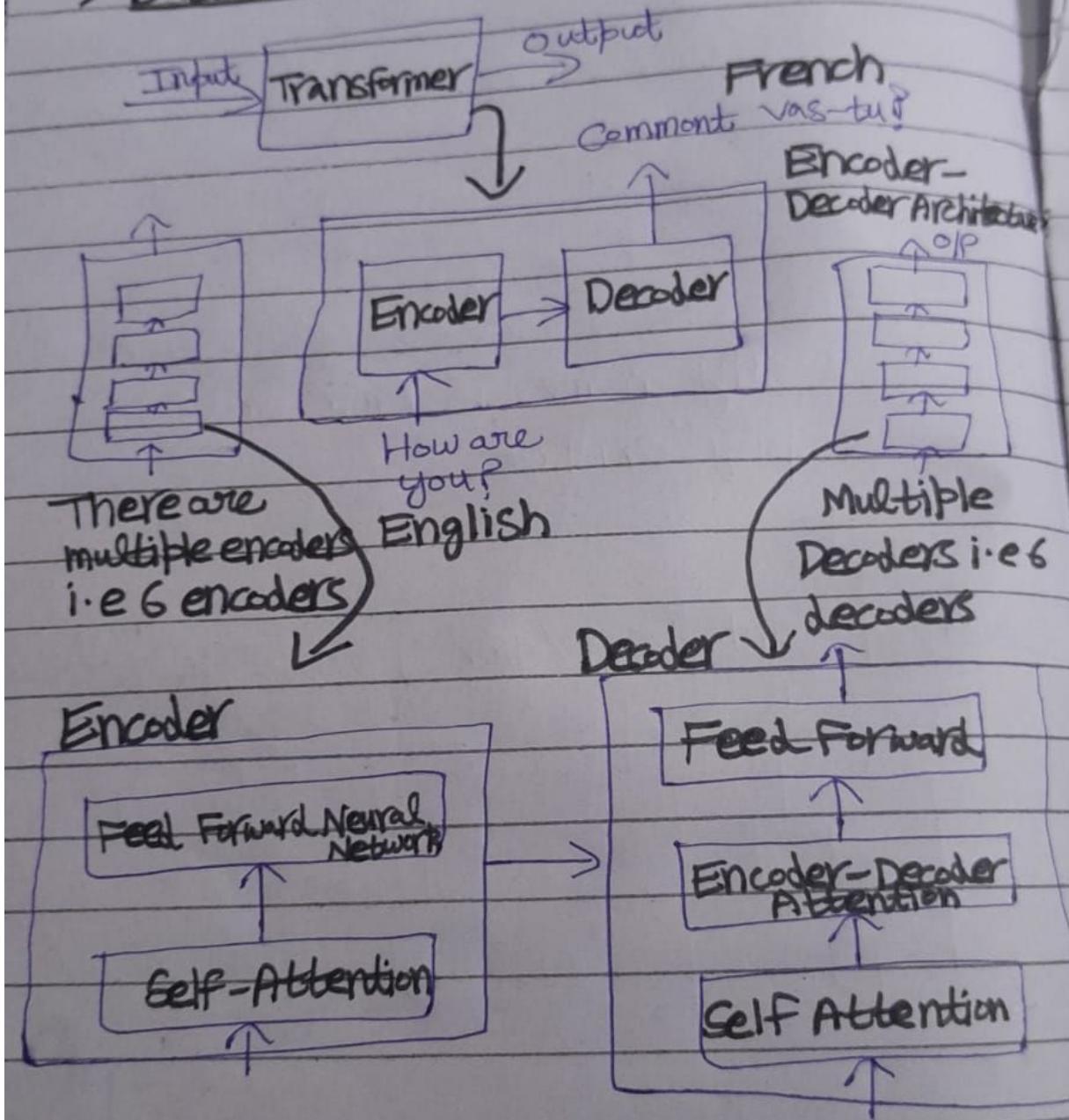
Trained
on huge data

⇒ Contextual Embedding → Self Attention
eg My name is Amir and I
play cricket

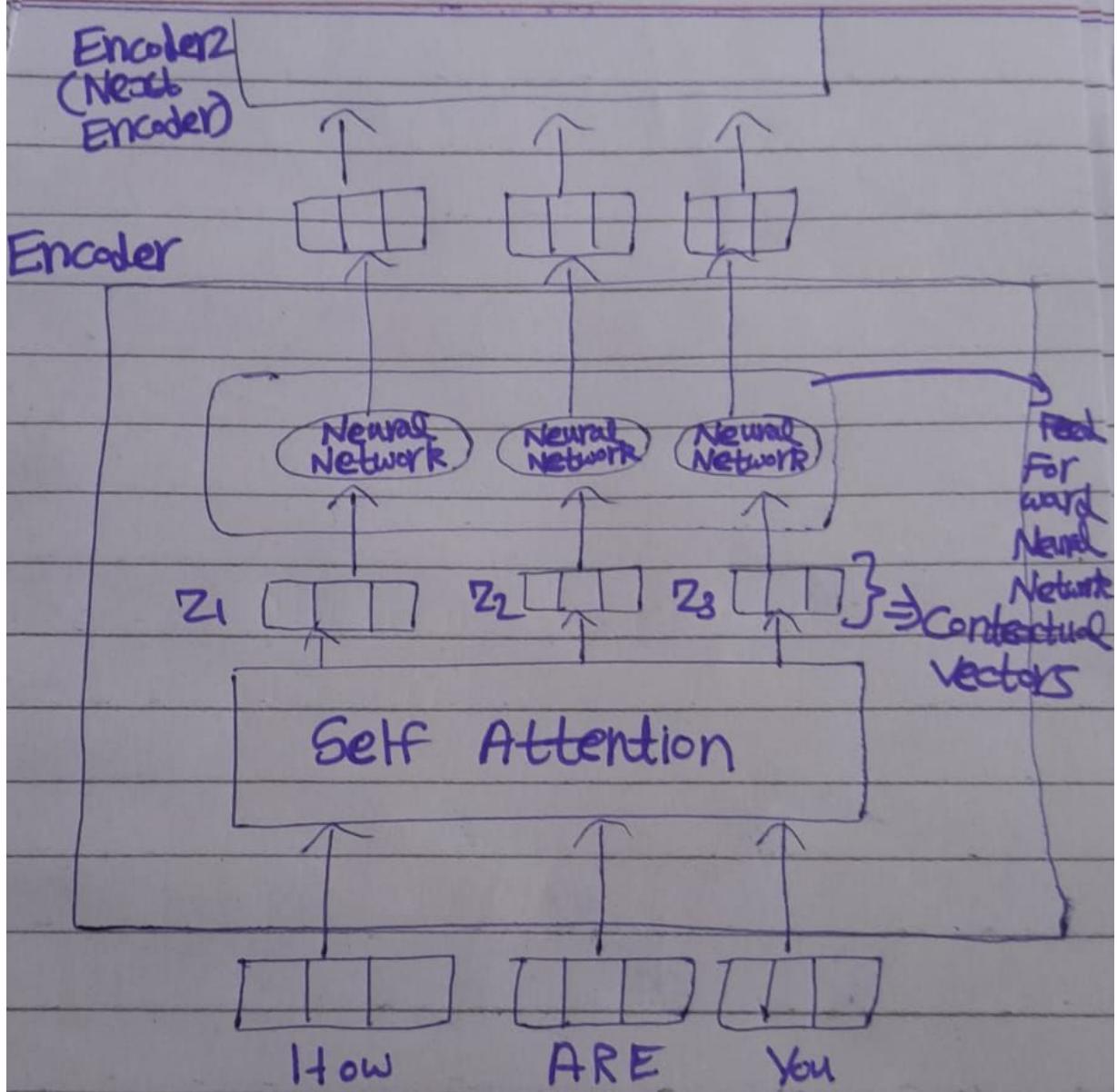


Contextual Vector Embedding (Here we get vectors that show relationship or context b/w words in a sentence using self attention)

⇒ Basic Transformer Architecture:



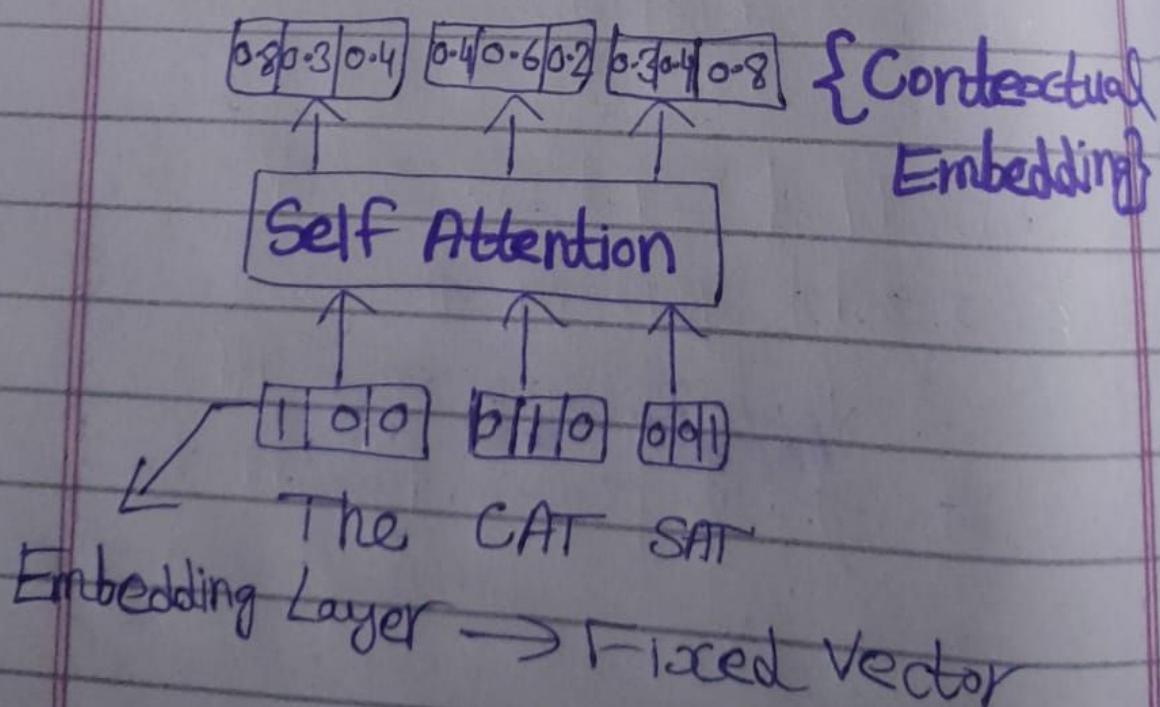
(Remaining on Next Register)



⇒ Self Attention at Higher
and detailed level:

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

Idea



1) Inputs: Queries, Keys, Values

Model \rightarrow Queries, Keys, Values

(i) Query Vectors (q):

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

Importance:

Focus Determination:

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

Contextual Understanding:

Queries contribute to understanding the relationship between current token and rest of sequences which is essential for capturing dependencies and context

(ii) Key Vectors (K):

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

(iii) Value Vectors (V):

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

Example:

Input sequence = ["The", "CAT", "sat"]
embedding size $e = 4$

① Token Embedding (using Embedding Layer)

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

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

$$E_{\text{sat}} = [1 \ 1 \ 1 \ 1]$$

② Linear Transformation:

Let create Q, K, V by
multiplying ^(dot product) the embeddings
by learned weight matrices
 W_Q, W_K and W_V .

Let's consider

$$W_Q = W_K = W_V = I = \begin{bmatrix} 1 & 0 & 0 & 0 \\ 0 & 1 & 0 & 0 \\ 0 & 0 & 1 & 0 \\ 0 & 0 & 0 & 1 \end{bmatrix}$$

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

$$Q_{\text{The}} = \begin{bmatrix} 1 \times 1 + 0 \times 0 + 1 \times 0 + 0 \times 0, 1 \times 0 + 0 \times 1 + 1 \times 1 + 0 \times 0 + 1 \times 0 + 0 \times 1 + 0 \times 1 \\ 0 \times 0 + 1 \times 0 + 0 \times 0 + 1 \times 1 + 0 \times 0 + 1 \times 0 + 0 \times 1 \\ 0 \times 1 \end{bmatrix}$$

As 'I' is identity matrix x offset remain same

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

So, similarly

$$(i) Q_{\text{The}} = K_{\text{The}} = V_{\text{The}} = [1 \ 0 \ 1 0]$$

$$(ii) Q_{\text{CAT}} = K_{\text{CAT}} = V_{\text{CAT}} = [0 \ 1 \ 0 \ 1]$$

$$(iii) Q_{\text{SAT}} = K_{\text{SAT}} = V_{\text{SAT}} = [1 \ 1 \ 1 \ 1]$$

③ Computing Attention scores:

(i) For the token "The"

$$\begin{aligned} \text{Score}(Q_{\text{The}}, K_{\text{The}}) &= [1 \ 0 \ 1 0] \cdot [1 \ 0 \ 1 0] \\ &= [1 \ 0 \ 1 0] \cdot [1 \ 0 \ 1 0] \end{aligned}$$

$$\text{Score}(Q_{\text{The}}, K_{\text{The}}) = [1 \times 1 + 0 \times 0 + 1 \times 1 + 0 \times 0]$$

$$\boxed{\text{Score}(Q_{\text{The}}, K_{\text{The}}) = 2}$$

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

$$\boxed{\text{Score}(Q_{\text{The}}, K_{\text{CAT}}) = 0}$$

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

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

(ii) For the token "CAT"

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

$$\boxed{\text{Score}(Q_{\text{CAT}}, K_{\text{The}}) = 0}$$

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

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

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

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

(iii) For the token "SAT"

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

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

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

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

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

$$\boxed{\text{Score}(Q_{\text{SAT}}, K_{\text{SAT}}) = 9}$$

Scaling:

We take up the scores and scale down by dividing the scores by the $\sqrt{d_K}$

Now

$d_K \Rightarrow$ dimension of key vector = 4

$$\boxed{\sqrt{d_K} = 2}$$

⇒ Scaling in the attention mechanism is crucial to prevent the dot product from growing too large to ensure stable gradient during training.

Problems here are:

① Gradient Exploding

② Softmax saturation (Vanishing

Now Let,

Gradient Problem)

$$Q = [2 \ 3 \ 4]$$

$$K_1 = [1 \ 0 \ 1 \ 0]$$

$$K_2 = [0 \ 1 \ 0 \ 1]$$

Without scaling:

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

$$[Q \cdot K_1^T = 6]$$

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

$$[Q \cdot K_2^T = 4]$$

Attention scores = [6, 4]

Let we don't applied scaling

$$\text{Softmax}([6, 4]) = \left[\frac{e^6}{e^6 + e^4}, \frac{e^4}{e^6 + e^4} \right]$$

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

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

$$\text{Softmax}([6, 4]) \approx [0.88, 0.12]$$



Most of the attention weights
is assigned to 1st key vector
and very little to the 2nd vector

- There is hardly difference of 2 still
so much difference in softmax so
if difference become larger then problem

With scaling:

$[6, 4] \Rightarrow$ Attention scores
To scale divide by $\text{ESR} = 2$
 $[6/2, 4/2]$

$[3, 2] \Rightarrow$ Attention scores
after scaling

Now apply Softmax

$$\begin{aligned}\text{Softmax}([3, 2]) &= \left[\frac{e^3}{e^3 + e^2}, \frac{e^2}{e^3 + e^2} \right] \\ &= \left[\frac{e^3}{e^2(1+e)} , \frac{e^2}{e^2(e+1)} \right] \\ &= \left[\frac{1}{(1+e)} , \frac{1}{(e+1)} \right]\end{aligned}$$

$$\boxed{\text{Softmax}([3, 2]) = [0.73, 0.27]}$$

↓
Attention
Weights

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

④ Scaling:

We know

$$d_K = 4$$

$$\sqrt{d_K} = 2$$

For token "The":

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

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

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

Similar scaling will be done for all other tokens.

Apply Softmax:

For "The"

$$\text{Softmax}([1, 0, 1]) = \left[\frac{e^1}{e^1 + e^0 + e^1}, \frac{e^0}{e^1 + e^0 + e^1}, \frac{e^1}{e^1 + e^0 + e^1} \right]$$

$$\text{Softmax}([1, 0, 1]) = [0.4223, 0.1554, 0.4223]$$



Attention Weights
for "The"

507

⑥ Weighted Sum of Values:

We multiply the attention weight by corresponding value vectors i.e actual vectors
 For token "The":

$$\text{Output}_{\text{The}} = 0.4223 \times V_{\text{The}} + 0.1554 \times V_{\text{CAT}} + 0.4223 \times V_{\text{SATT}}$$

$$\text{Output}_{\text{The}} = 0.4223 \times [1 \ 0 \ 1 \ 0] + \\ 0.1554 \times [0 \ 1 \ 0 \ 1] + \\ 0.4223 \times [1 \ 1 \ 1 \ 1]$$

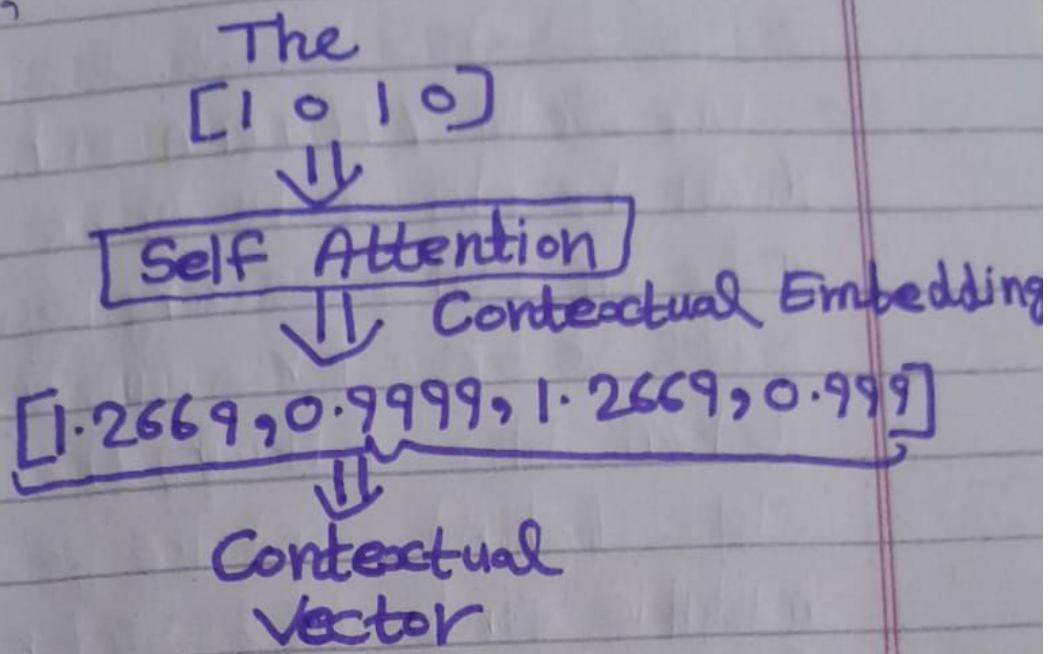
$$\text{Output}_{\text{The}} = [0.4223, 0, 0.4223, 0] + \\ [0, 0.1554, 0, 0.1554] + \\ [0.4223, 0.4223, 0.4223, 0.4223]$$

$$\boxed{\text{Output}_{\text{The}} = [1.2669, 0.9999, 1.2669, 0.9999]}$$



Contextual
Vector

So,



• Summarized Steps:

- ① Calculate Q, K, V
- ② Compute Attention Scores
- ③ Scale Attention Scores
- ④ Perform Softmax
- ⑤ Weighted sum of Values

⇒ Multi-head Attention: ⇒

- In Self Attention with Multi-heads actually we make multiple query, key and value vectors for a single token like following the steps behind and we get ~~multiple~~ ^{attention} contextual vectors (heads) for single token.

- It actually expands model's ability to focus on different positions of tokens

How we create multiple Q, K and V:

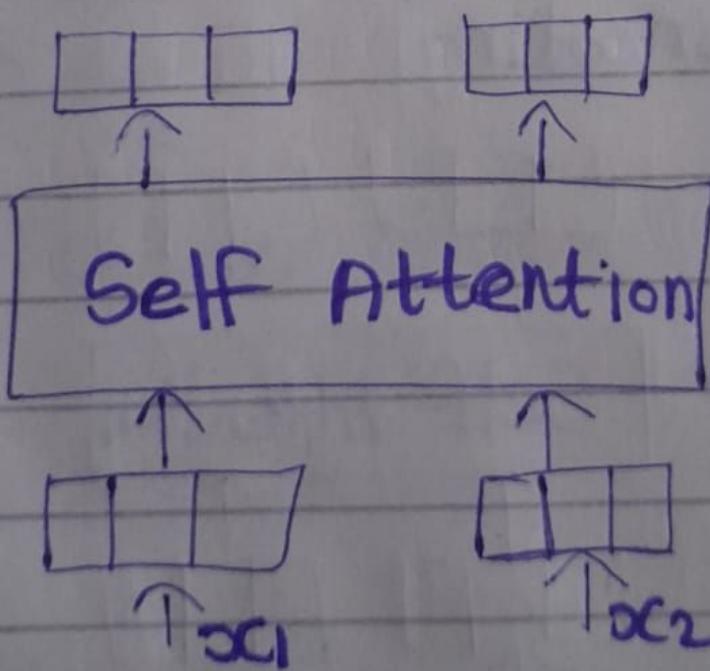
By initializing multiple different weights W_Q , W_K and W_V .

(For remaining information see the repo)

Feed Forward Neural Network (In Repo)

⇒ Positional Encoding
(Representing order of sequence)

- Major Advantage of transformers is that it can process multiple words tokens parallelly.
- Major Drawback of transformers is that it lack the sequential structure (order) of words.



e.g

① Lion kills tiger

② Tiger kills lion

Both sentences are opposite.

but as transformer we

pass all words parallelly

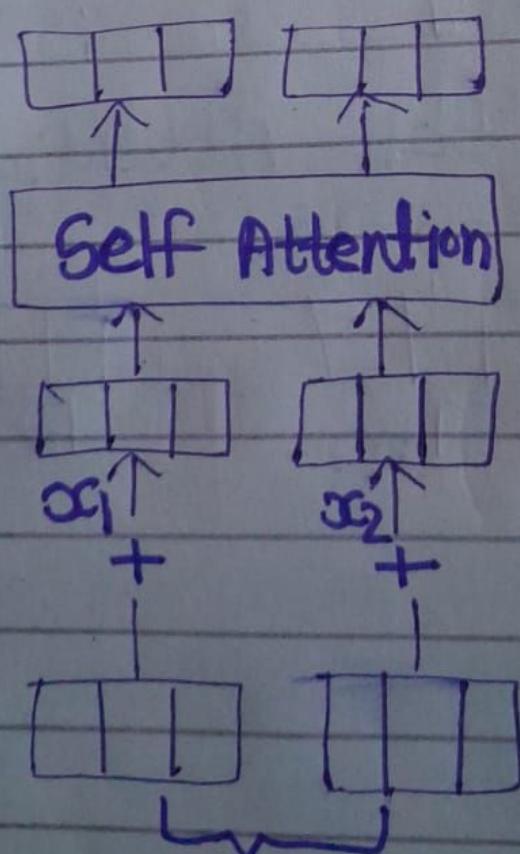
and lack sequence ordering

- so it will take both as same.

⇒ To resolve this issue

we use Positional

Encoding



Positional Encoded
vectors (tell about ordered)

Types of Positional Encoding:

- 1) Sinusoidal Positional Encoding
- 2) Learned Positional Encoding

Positional Encodings are learned during training

1) Sinusoidal Positional Encoding:

It uses sine and cosine functions of different frequencies to create positional embeddings.

Formula:

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

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

where,

$\text{pos} \Rightarrow$ position

$i \Rightarrow$ is the dimension

$\text{dimodel} \Rightarrow$ is the dimensionality of embeddings

Example:
"The cat sat"

The \rightarrow $[0.1 \ 0.2 \ 0.3 \ 0.4]$
cat \rightarrow $[0.5 \ 0.6 \ 0.7 \ 0.8]$
sat \rightarrow $[0.9 \ 1.0 \ 1.1 \ 1.2]$

Now these are embedded
direct vectors before passing
them to self attention we
have to compute positional
encoded vectors using
sinusoidal positional encoding.

We know,

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

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

For our example

dmodel (dimensionality) = 4

For position pos=0 i.e "The"

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

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

$$P.E(0,2) = \sin\left(\frac{0}{10000^{2\pi}}\right) = 0$$

$$P.E(0,3) = \cos\left(\frac{0}{10000^{2\pi}}\right) = 1$$

P.E = [0, 1, 0, 1] \Rightarrow "The"

For position pos=1 i.e "Cat"

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

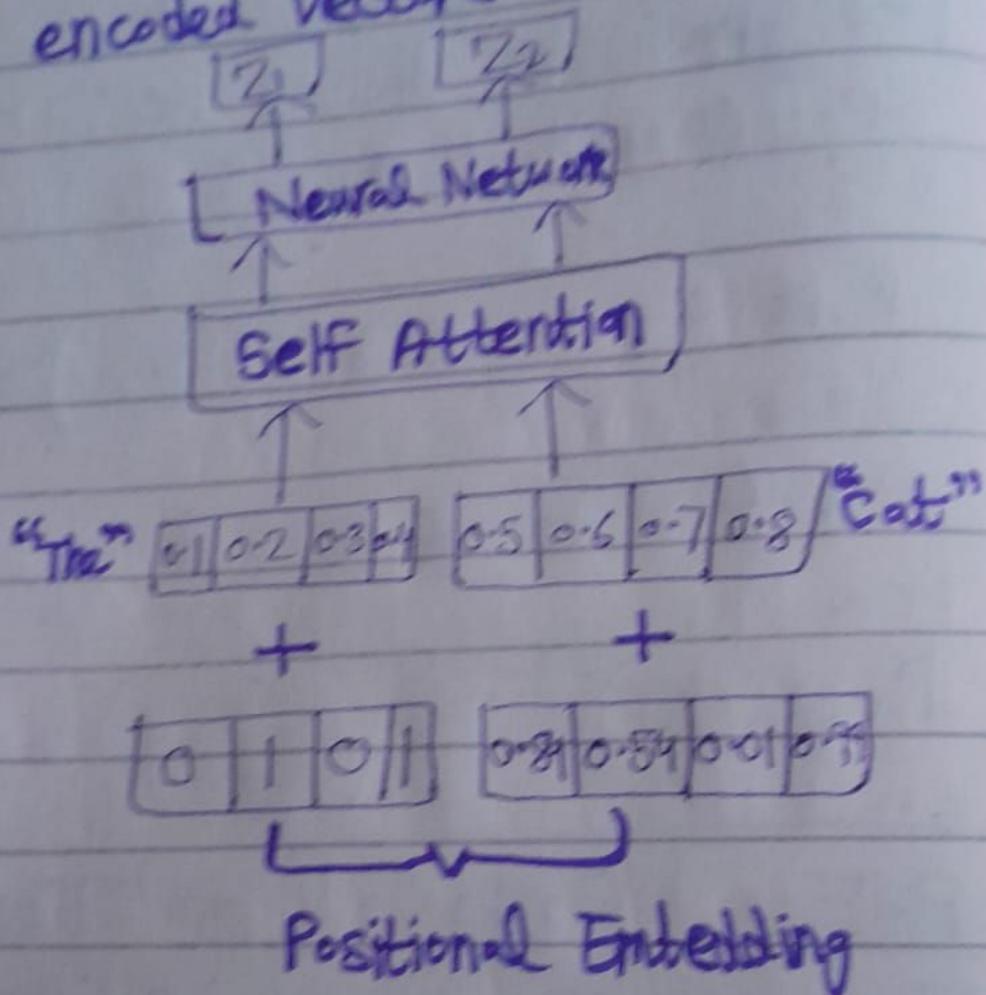
$$P.E(1,1) = \cos\left(\frac{1}{10000\pi}\right) = 0.5403$$

$$P.E(1,2) = \sin\left(\frac{1}{10000^{2\pi}}\right) = 0.01$$

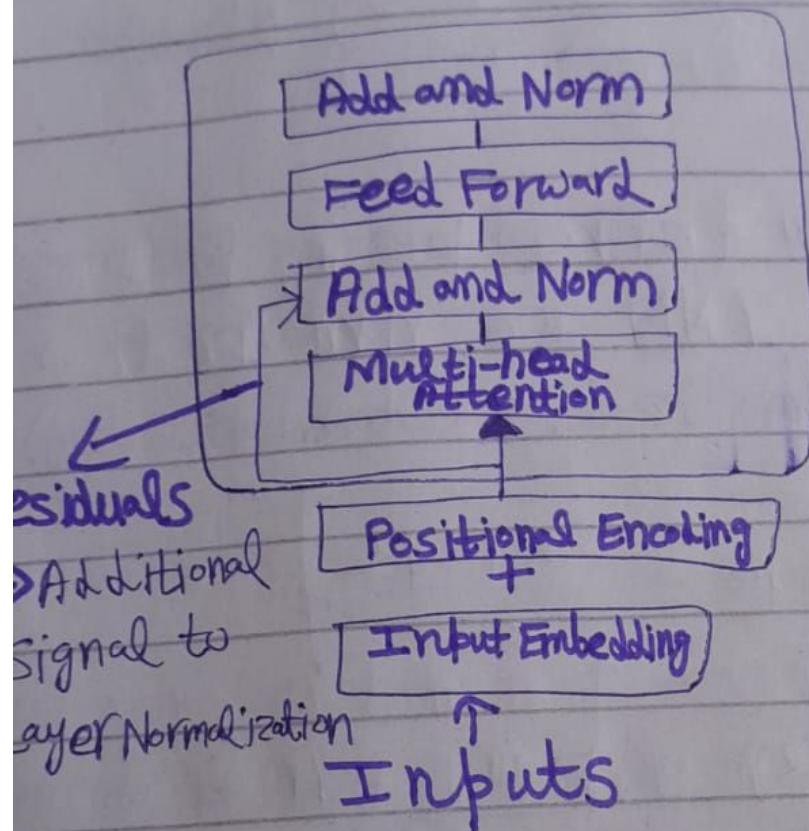
$$P.E(1,3) = \cos\left(\frac{1}{10000^{2\pi}}\right) = 0.9995$$

E = [0.8415, 0.5403, 0.01, 0.9995] \Rightarrow "Cat"

Now we have find positional
encoded vectors



⇒ Layer Normalization in
transformers:

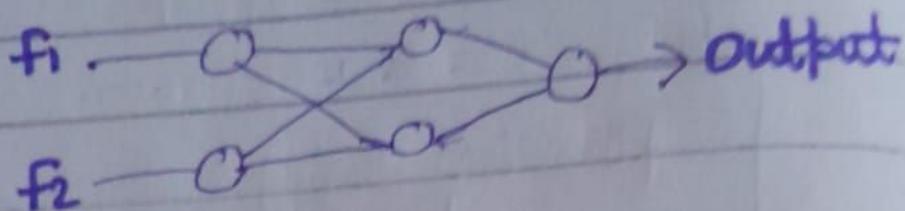


Normalization \Rightarrow Batch Normalization
Layer Normalization

Let we have dataset

f_1 (House Size)	f_2 (No. of Rooms)	Price
1200	2	45
1500	3	70
2000	3.5	80

Let ANN is :



- Now before passing features to ANN we normalize them
- Normalization (Standard Scaling):

$$Z\text{-score} = \frac{x_i - \mu}{\sigma}$$

After applying normalization,

$$\mu=0, \sigma=1$$

⇒ Advantages of Normalization:

- 1: Improved training stability
(as we will not face vanishing or exploding gradient problem)
- 2: Faster convergence (as data is zero-centered) (stable update during back propagation)

\Rightarrow Batch Normalization:

In batch Normalization we take whole column like f_i (house size) compute μ (mean) and σ (standard deviation) and apply formula

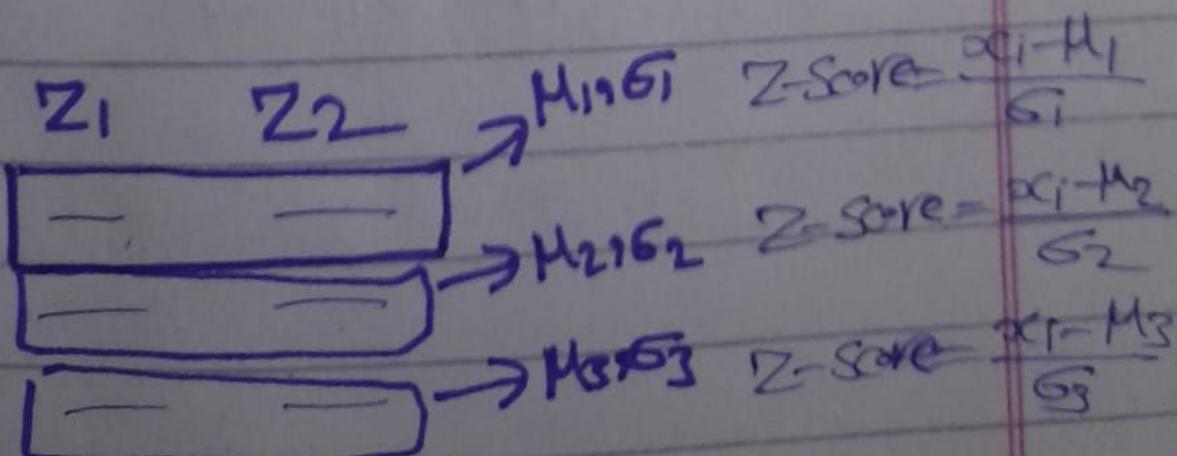
$$z\text{-score} = \frac{x_i - \mu}{\sigma}$$

to every value in column.

This is usually used in ANN.

- But if the column contains multiple zeros then batch normalization will also affect Non-zero values.
- That's why, in Transformers we use Layer Normalization.

\Rightarrow Layer Normalization



- In Layer Normalization, we compute μ and σ for every layer (row) and apply

$$Z\text{-Score} = \frac{x_i - \mu}{\sigma}$$

to the entries of that row.

$\gamma, \beta \rightarrow$ Learnable parameters
scale and shift
parameters

- These parameters control layer normalization that how much normalization we want to do or whether do or not.

$$y = \gamma \left[\frac{z_i - \mu_i}{\sigma_i} \right] + \beta$$

• Example:

(1) "CAT" = [2.0, 4.0, 6.0, 8.0]

(2) Parameters:

$$\gamma = [1.0, 1.0, 1.0, 1.0] \rightarrow \text{Learned Scale}$$

$$\beta = [0.0, 0.0, 0.0, 0.0] \rightarrow \text{Shift}$$

↓
Scale and
Shift parameters

Step 1: Compute the mean

$$\mu = \frac{\sum x_i}{n} = \frac{2.0 + 4.0 + 6.0 + 8.0}{4}$$

$$\boxed{\mu = 5.0}$$

Step 2: Compute the variance

$$\sigma^2 = \frac{\sum (x_i - \bar{x})^2}{n}$$

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

$$\boxed{\sigma^2 = 5.0}$$

Step.03 Normalize the input

Step.4) S

$$\hat{x}_1 = \frac{x_1 - \mu}{\sqrt{\sigma^2 + \epsilon}}$$

Let

$\epsilon = 1e^{-5} \Rightarrow$ To avoid division by zero

$$\sqrt{\sigma^2 + \epsilon} = \sqrt{5.0 + 1e^{-5}} = \sqrt{5.0001}$$

$$\boxed{\sqrt{\sigma^2 + \epsilon} = 2.236}$$

$$\hat{x}_1 = \frac{2.0 - 5.0}{2.236} \approx -1.34$$

$$\hat{x}_2 = \frac{4.0 - 5.0}{2.236} \approx -0.45$$

$$\hat{x}_3 = \frac{6.0 - 5.0}{2.236} \approx 0.45$$

$$\hat{x}_4 = \frac{8.0 - 5.0}{2.236} \approx 1.34$$

$$\boxed{\hat{x} = [-1.34, -0.45, 0.45, 1.34]}$$

Normalized Vector

Step:4) Scale and Shift:

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

$$\text{We know } \gamma = [1.0, 1.0, 1.0, 1.0]$$

$$\beta = [0.0, 0.0, 0.0, 0.0]$$

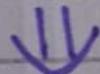
$$\hat{x} = [-1.34, -0.45, 0.45, 1.34]$$

$$y_1 = \gamma_1 \hat{x}_1 + \beta_1 = 1.0(-1.34) + 0.0 = -1.34$$

Similarly calculate for y_2, y_3, y_4

The final Normalized Vector:

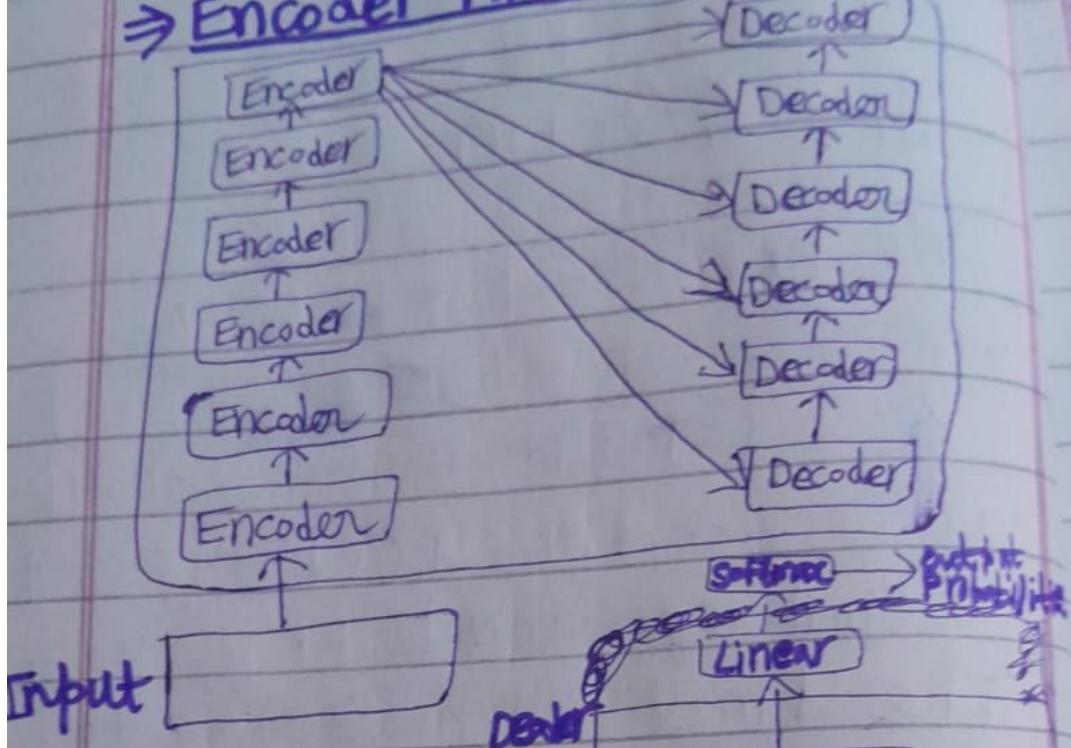
$$y = [-1.34, -0.45, 0.45, 1.34]$$



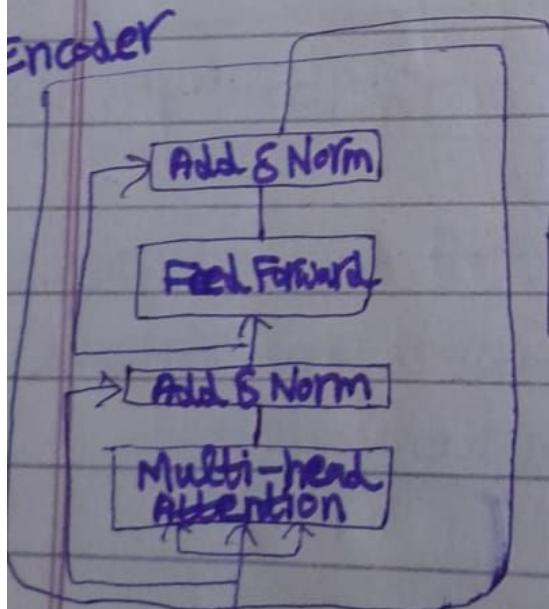
Output of Layer Normalization

(will pass to Feed Forward
Neural Network)

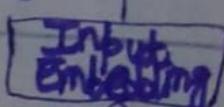
Encoder Architecture



Encoder



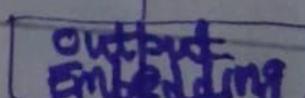
sitional
encoding



Inputs



Positional
Encoding



Outputs

According to research paper,
Embedding vector dimensions = 512
No. of encoders & Decoders = 6
 $Q, K, V = 64$
Head Attention = 8

⇒ Residual Connection
(skip connection): (Why to use)

1) Addressing the vanishing
gradient problem:

Residual connection creates
a short paths for gradients
to flow directly through
the network so gradient
remains sufficiently large.

2) Improve gradient flow

Convergence will be
faster

3) Enables training of
deeper networks

⇒ Feed Forward Neural Network (Why to use).

⇒ D

res

se

er

g

- ① Adding Non-linearity
- ② processing each position independently

We know

Self Attention captures relationships

Feed Forward Network

↓
Each token representation
Independently
↓

Transforming these representations further and allows the model to learn richer representations.

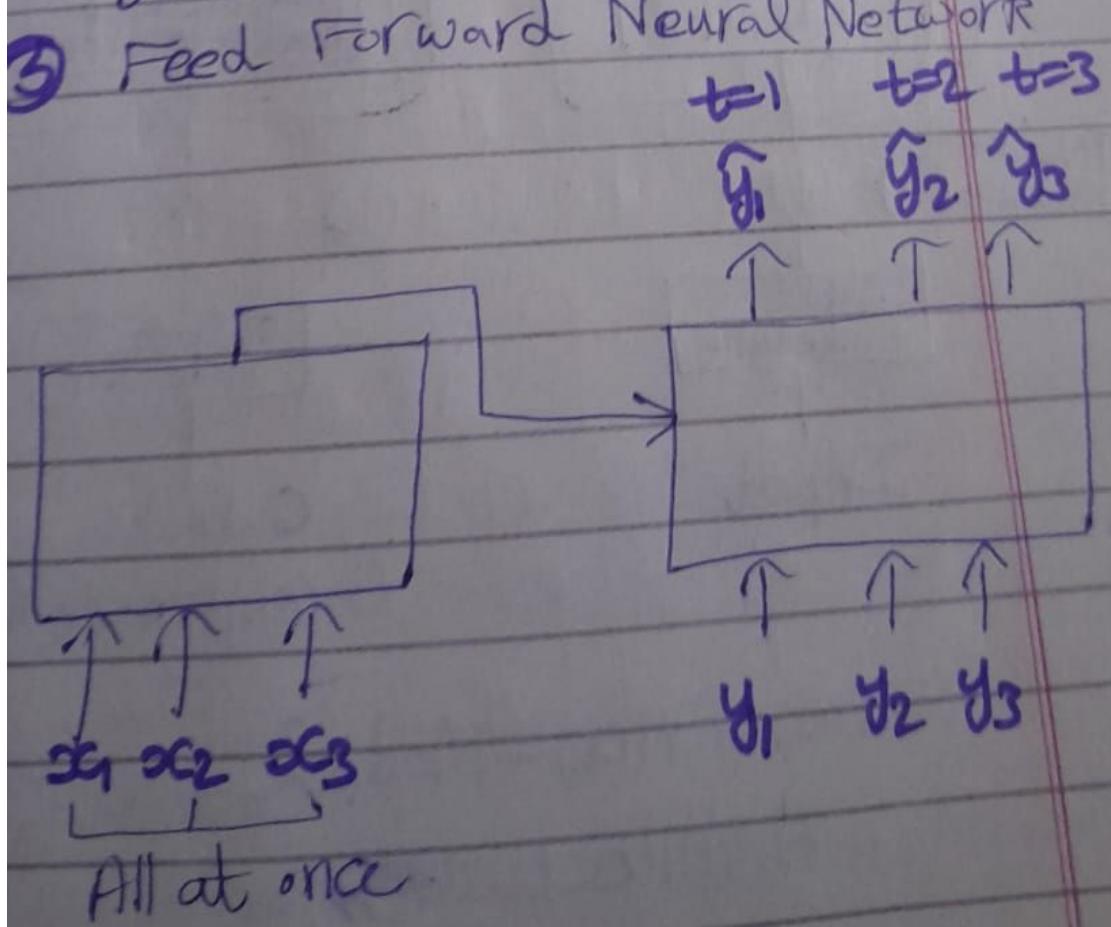
- ③ Adds Depth to the model as More depth more learnings and performance

→ Decoders in Transformers:

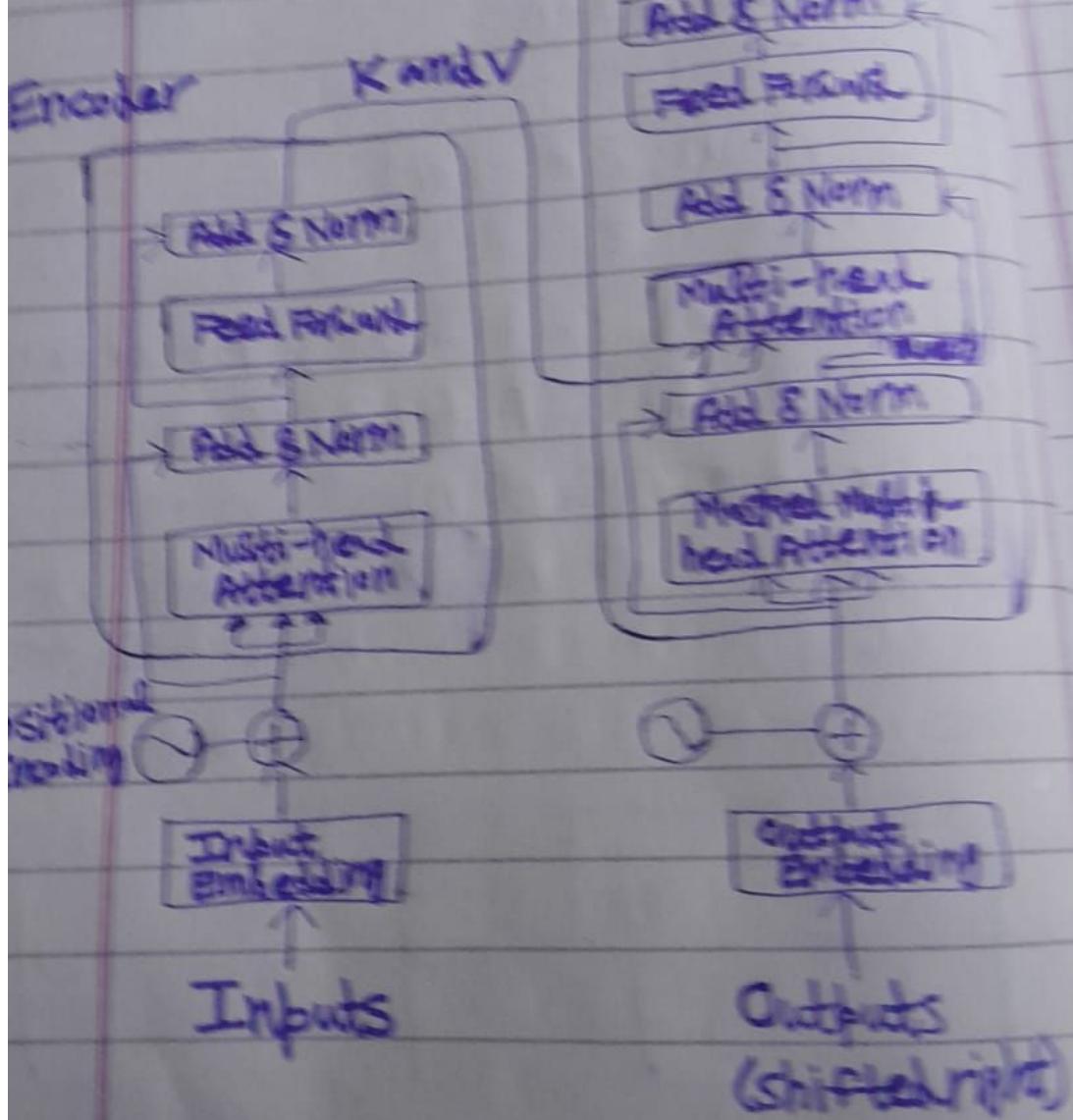
The transformer decoder is responsible for generating the output sequence one at a time using encoder's output and previously generated tokens

3 main components of decoder are:

- ① Masked Multi-head Attention
- ② Multi-head Attention (Encoder-Decoder Attention)
- ③ Feed Forward Neural Network



→ Masked Multi-head
Self Attention



Transformer-Model Architecture

Let we have a dataset:

English

$\langle x_1, x_2, x_3 \rangle$

Urdu

$\langle y_1, y_2 \rangle$



We have to
pass this to
decoder



We will do

zero padding for
making sequence length
equal i.e

$\langle y_1, y_2, 0 \rangle$ (Shifted
Right)

Now this output will be zero

padded.

i) Input Embedding and
Positional Encoding:

Let

Input = $[4 \ 5 \ 6 \ 7]$

Output = $[1 \ 2 \ 3] \Rightarrow$ We will pad it

Output = $[1 \ 2 \ 3 \ 0]$

Let every single element in output is represented by 4 dimensional vector.

Output embedding:

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

- Now we will add **Positional encoding** to it. Let positional encoded vectors are all 0 so output embedding will remain the same.

② Linear Projection for Q, K and V

Create Query, Key and Value vectors

Let $W_Q = W_K, W_V = I$ (Identity Matrix)

Sor
Q =
K =
V =
AS V
Sor
Q =

③

So,

$$Q = \text{output Embedding} \bullet W_a$$

$$K = \text{output Embedding} \bullet W_k$$

$$V = \text{output Embedding} \bullet W_v$$

$$\text{As } W_a = W_k = W_v = I \text{ (Identity)}$$

So,

$$Q = K = V = \begin{bmatrix} [0.1, 0.2, 0.3, 0.4], \\ [0.5, 0.6, 0.7, 0.8], \\ [0.9, 0.1, 0.11, 0.12], \\ [0.0, 0.0, 0.0, 0.0] \end{bmatrix}$$

⑤ Scaled Dot product Attention Calculation:

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

∴ Dividing by $\sqrt{d_K}$ for scaling

$$\therefore \sqrt{d_K} = 2$$

$$\text{Scores} = Q^* K^T / 2$$

We know how to do this from encoder

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

⇒ Masked Application:

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

Reasons:

① Handelling variable length sequences with padding mask.

Purpose:

- (i) To handle sequences of different length in batch
- (ii) To ensure that padding sequences/tokens which are added to make sequences of uniform length don't affect model prediction.

e.g

Input \rightarrow Sequence 1 $\rightarrow [1, 2, 3]$

Output \rightarrow Sequence 2 $\rightarrow [4, 5]$

\downarrow zero padding

[4, 5, 0]

\downarrow Influence the
attention mechanism

\downarrow Lead to incorrect or
biased predictions

Padding Mask ensures that
padding tokens are ignored

Masking \Rightarrow Padding Mask
 \Rightarrow Look Ahead Mask

Padding Mask for sequence 1: [1 | 0]
(As no zero
padding)

Padding Mask for sequence 2: [1 | 0 0]
(As last one
is zero padded)

→ Look ahead mask

- * Maintain auto-regressive property
- * To ensure that each position in the decoder output sequence can only attend to the previous position but no future position.

e.g

↑ AD Padding mask

$$[4, 5, 0] \rightarrow [1, 1, 0]$$

↓ convert 1D
to 2D MASK

Token1 attends

1, 2

Token2 attends

1, 2

Token3

attends
anything

1	1	0
1	1	0
0	0	0

↓
Padding
MASK

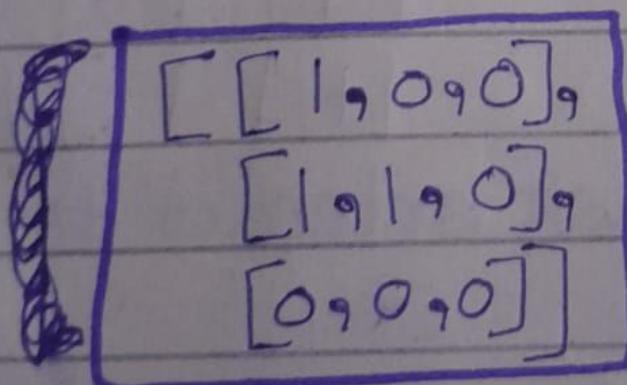
* Look Ahead Mask:

$$\begin{bmatrix} [1, 0, 0] \\ [1, 1, 0] \\ [1, 1, 1] \end{bmatrix} \quad (\text{As we don't need future context})$$

⇒ Combine padding and look ahead mask:

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

(Do element wise multiplication)


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

⇒ Combine Mask

• Mask:

Now let's compute mask
for the example we are
doing

Scores: $\begin{bmatrix} [0.3, 0.7, 1.9, 0.0] \\ [0.7, 1.9, 3.1, 0.0] \\ [1.1, 3.1, 5.1, 0.0] \\ [0.0, 0.0, 0.0, 0.0] \end{bmatrix}$
 $\begin{bmatrix} [1 \ 0 \ 0 \ 0] \end{bmatrix}$

Look ahead Mask = $\begin{bmatrix} 1 & 1 & 0 & 0 \\ 1 & 1 & 1 & 0 \\ 1 & 1 & 1 & 1 \end{bmatrix}$

Padding Mask 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
We will do element wise multiplication

$$[[1, 0, 0, 0],$$

$$[1, 1, 0, 0],$$

$$[1, 1, 1, 0],$$

$$[0, 0, 0, 0]]$$

↓ (Wherever there is 0
mark -∞) (Add -∞ to

0) (As it zero
out the influence
when softmax
is applied)

$$[[1, -\infty, -\infty, -\infty],$$

$$[1, 1, -\infty, -\infty],$$

$$[1, 1, 1, -\infty],$$

$$[-\infty, -\infty, -\infty, -\infty]]$$

Masked Score:

Now we will do element wise
multiplication of scores and mask

$$[[0.3, -\infty, -\infty, -\infty],$$

$$[0.7, 1.9, -\infty, -\infty],$$

$$[1.1, 3.1, 5.1, -\infty],$$

$$[-\infty, -\infty, -\infty, -\infty]]$$

⇒ Apply softmax

Softmax Score =
Softmax (masked scores)

$$\text{Softmax score} = \begin{bmatrix} [1.0, 0.0, 0.0, 0.0] \\ [0.3, 0.7, 0.0, 0.0] \\ [0.1, 0.3, 0.6, 0.0] \\ [0.0, 0.0, 0.9, 0.0] \end{bmatrix}$$

⇒ Weighted sum of values:

Attention output = Softmax
scores * \downarrow

\downarrow
Value
Vectors

\Rightarrow Encoder Decoder Multi-head
Attention:

- ① Encoder output \rightarrow Sending Attention vectors K and V.
- ② Masked Multihead \rightarrow Sending attention vector Q

So, encoder-decoder multi-head attention will get K and V (key and value vectors) from encoder and Q (query vector) from masked multi-head attention, and then it will perform all the things similarly as performed by encoder multi-head attention and output passed to Add & Norm layer.

At timestamp = 1 some token (start token) will get passed to masked multi-head attention and output will be generated, for t=2 previous output will get passed to decoder.

Linear

⇒ The final and Softmax Layer:

- We get vectors as output from decoder. In this layer we will see how we will convert output vectors to words.

Linear Layer:

The Linear Layer is a simple fully connected neural network that projects the vector produced by decoders.

- It generates a large vector called **logits vector**
- Logits vector size = model vocabulary size

Softmax Layer:

Softmax Layer turns those scores into probabilities.

- The cell with highest probability

is chosen and the word associated with it is produced as output.

- Then we will calculate loss function and back propagate, to reduce the loss.