\* **Attention**

Lets say $Q$ is a query vector

$K_i$ is a key vector

$V_i$ is a value vector

Every key is associated with a value vector

**Types**

(1) Additive Attention

$$S_i = w_3 \tanh \left( w_2^T Q + w_1^T K_i \right)$$

This equation can be related to the attention seen in RNNs

$$t_{ATT} = v^T \tanh \left( U S_{t-1} + W h_j + b \right)$$

$\downarrow$ decoder state      $\downarrow$ encoder state

(2) Dot - Product Attention

$$S_i = Q^T K_i$$

Here in the above types $S_i$ is the similarity between $Q$ & $K_i$

Now,

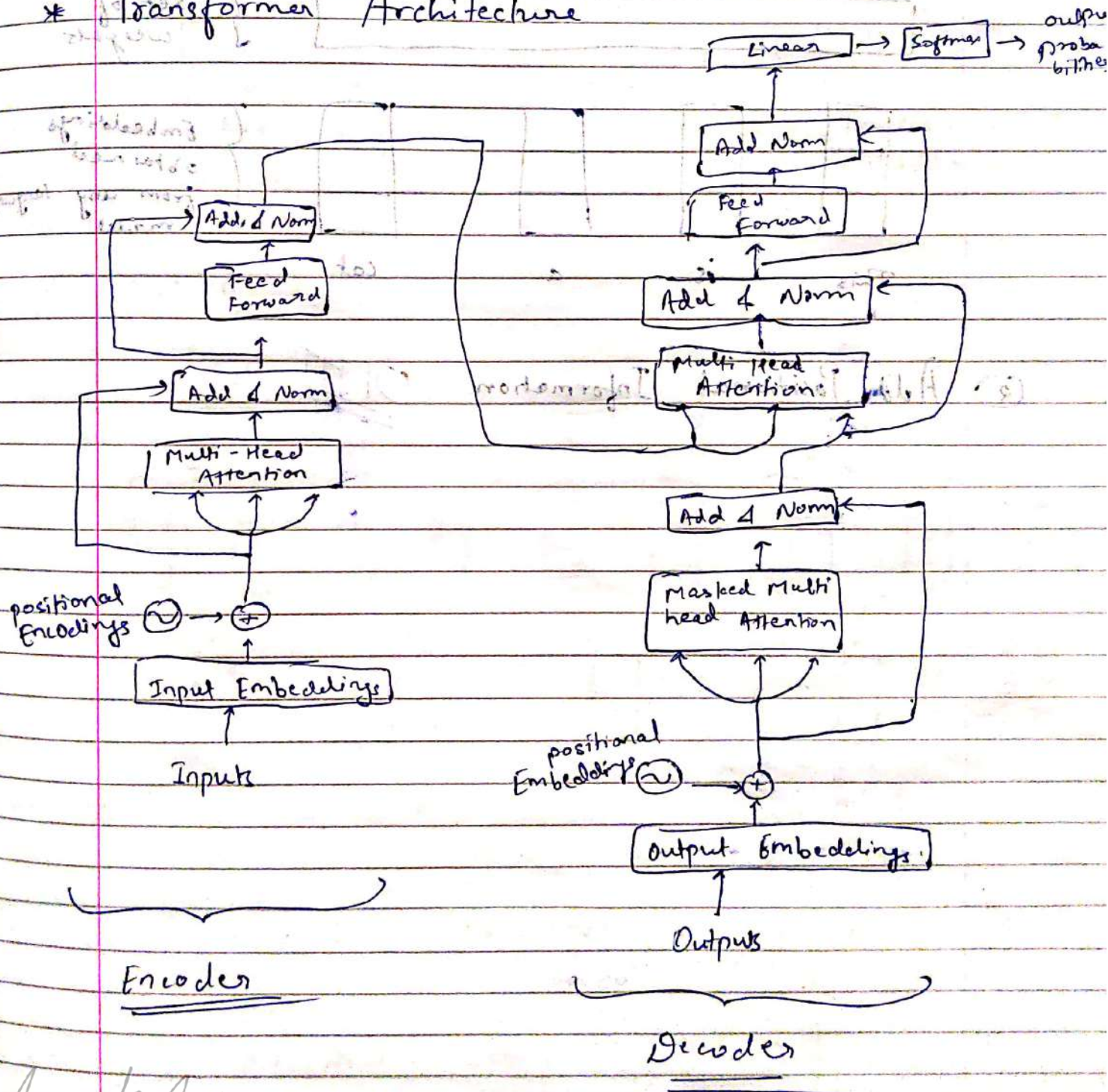$$A(Q, K, V) = \sum_i \frac{e^{S_i}}{\sum_j e^{S_j}} V_i$$

Ishan Modi

$$= \text{softmax} \, (S_i) \cdot V_i$$

softmax of $S_i$ $(Q \, \& \, k_i)$

over $( \, S_1 \, (Q \, \& \, k_1), \, S_2 \, (Q \, \& \, k_2) \, \ldots \ldots )$

**\* Transformer Architecture**

```
                                        Linear → Softmax → outpu
                                                            proba
                                                            bilite
                                        Add Norm
                                        Feed
                                        Forward
      Add & Norm                        Add & Norm
      Feed                              Multi Head
      Forward                           Attention
      Add & Norm
      Multi-Head                        Add & Norm
      Attention
                                        Masked Multi
                                        head Attention
positional                              
Encodings  → +                          
                        positional
  Input Embeddings       Embeddings ~ → +

      Inputs              Output Embeddings

                              Outputs
      Encoder
                              Decoder
```

Ishan Modi

→ Encoder

① Construct input Embedding

This        is        a        dog



} Word Embeddings obtained from any language model

} set of weights

} Transformed Embeddings

② Add Positional Information



} Transformed Embeddings

Adding positional vectors

} Transformed Embeddings with positional information

| Time stamp | 1 | 2 | 3 | 4 |

$$PE\left(pos, 2i+1\right) = cos\left(\frac{pos}{10000^{2i/dmodel}}\right)$$ Odd timestamp

$$PE\left(pos, 2i\right) = sin\left(\frac{pos}{10000^{2i/dmodel}}\right)$$ Even timestamp

→ Another interpretation of positional vectors

sinusoidal



high value

low value

high value

low value

high value

low value

↓      ↓

$\{0,0,0\}$     $[0,0,1]$

If the word is at start of sentence it gets the above positional embedding
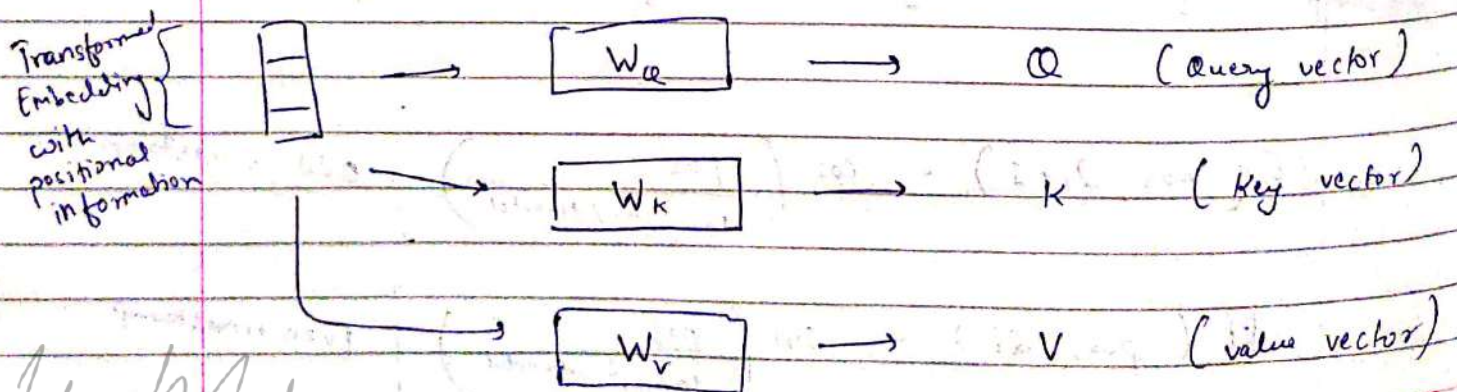
The word is somewhere in middle it gets the above embedding

**Note**

If there are more words more number of sinusoidal waves are used to construct the embeddings
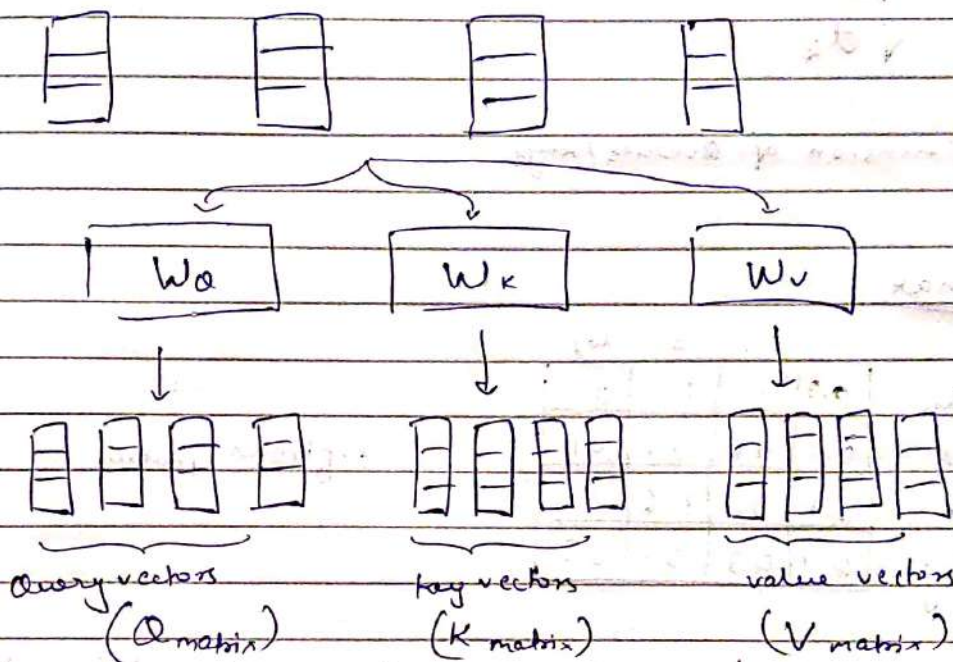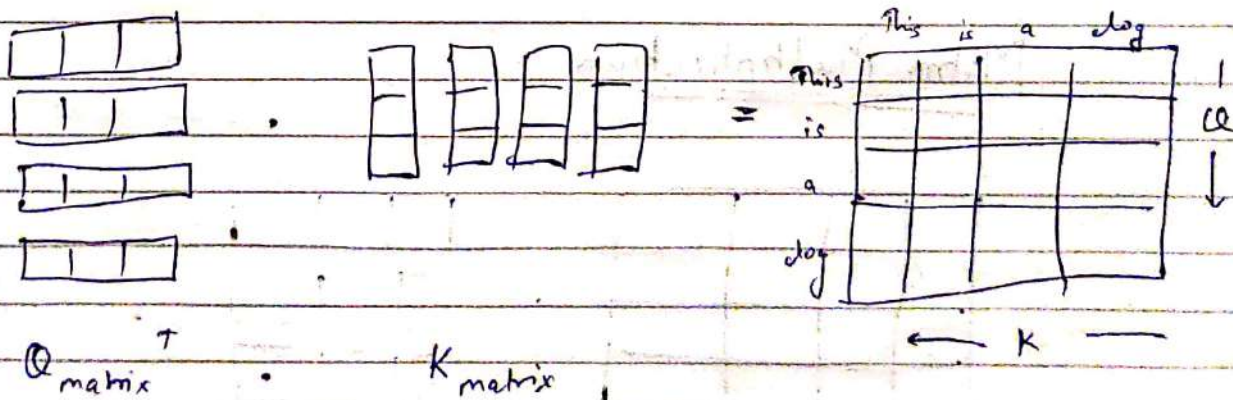
③ Multi-Headed Attention

→ Self Attention

Transformed Embedding with positional information



$W_Q$ → Q (Query vector)

$W_K$ → K (Key vector)

$W_V$ → V (value vector)

Ishan Modi

This is called self attention because Q, K & V are generated from input embeddings

For multiple vectors



Query vectors
(Q matrix)

key vectors
(K matrix)

value vectors
(V matrix)

→ <u>Matrix Multiplication</u>

Every ~~key vector~~ Query vector is multiplied by all key vectors



$Q_{matrix}^T$       $K_{matrix}$

This    is    a    dog

The matrix obtained after multiplication tells us what is importance of key (column) in the query (row)

↦ <u>Scale</u>

$$\frac{\text{Obtained matrix} \left(\boxplus\right)}{\sqrt{d_k}} = \text{scaled Matrix}$$

↙ dimension of Queries/keys

<u>Softmax</u>

|      | This | is  | a   | dog  |
|------|------|-----|-----|------|
| This | 0.7  | 0.1 | 0.1 | 0.1  |
| is   | 0.1  | 0.6 | 0.2 | 0.1  |
| a    | 0.1  | 0.3 | 0.6 | -0.1 |
| dog  | 0.1  | 0.3 | 0.3 | 0.3  |

softmax matrix

Apply softmax horizontally on Query of scaled matrix to get above matrix

$$\text{softmax}(x)_i = \frac{\exp(x_i)}{\sum_j \exp(x_j)}$$

<u>Matrix Multiplication</u>



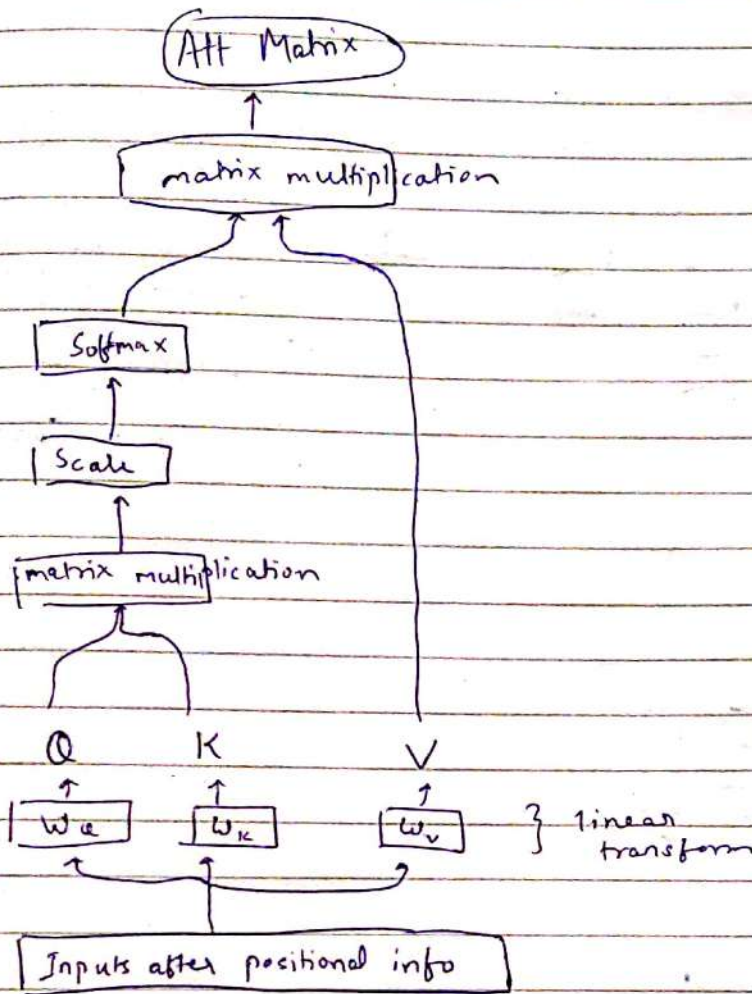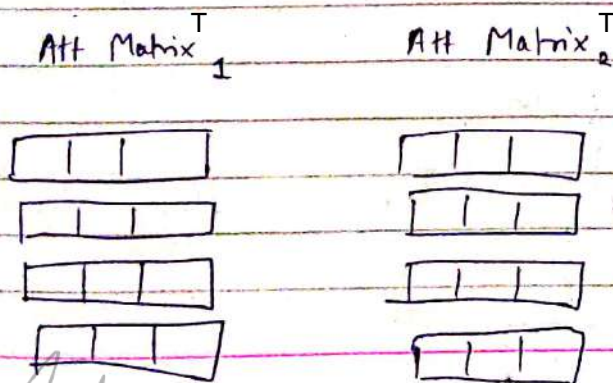softmax matrix     $V_{matrix}^{T}$     $Att_{matrix}^{T}$
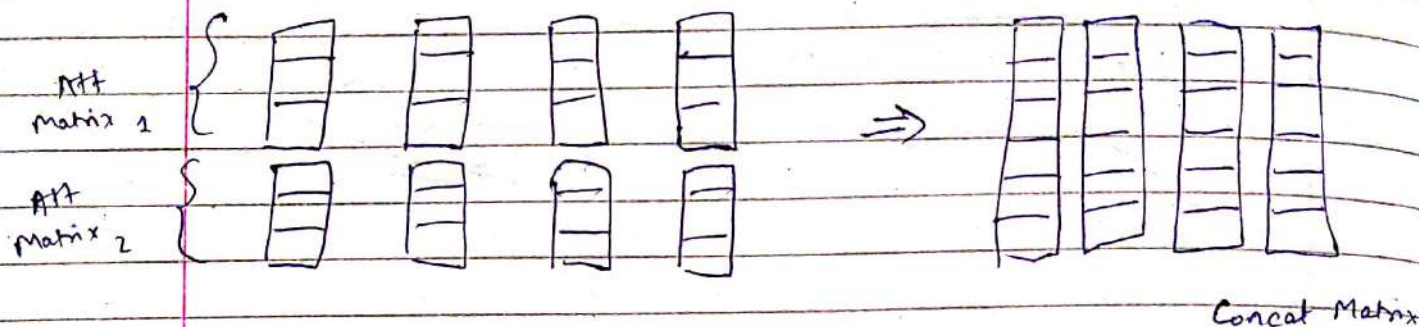
Ishan Modi

Till now,



Figure ①

→ Above is one attention module. Lets say we have multiple $W_Q$, $W_K$, $W_V$ and thus multiple attention modules

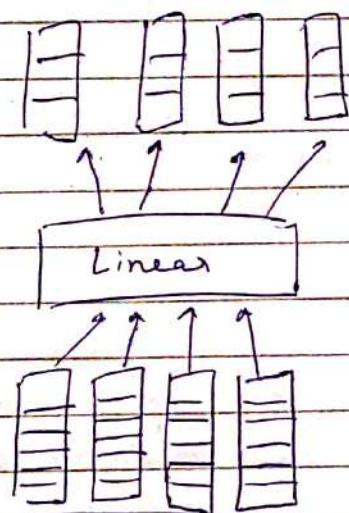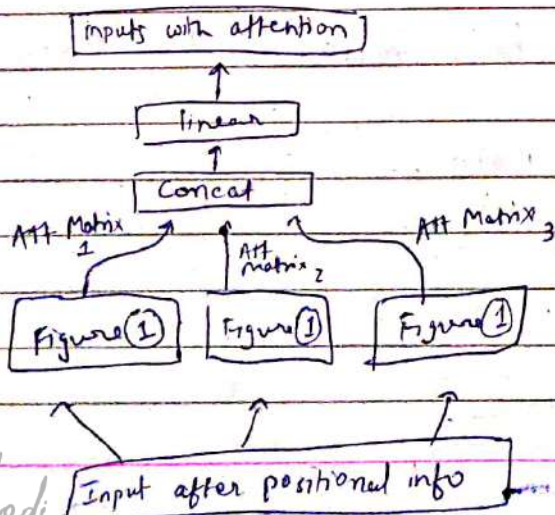In this case we would have multiple <u>Att Matrix</u>.

Att Matrix$^T_1$          Att Matrix$^T_2$



Ishan Modi

# Concatenation



Att matrix 1

Att matrix 2

Concat Matrix

# Linear



Linear

**Thus**

inputs with attention

linear

Concat

Att Matrix 1

Att matrix 2

Att Matrix 3

Figure ①    Figure ①    Figure ①

} multiple attention heads

**Final figure for multiheaded Attention**

Input after positional info

④ Add & Normalize

Layer Norm ( [⊞⊞] + [⊞⊞] ) ⇒ [⊞⊞]

inputs after
positional info.
( Residual
connection)

inputs with
attention

normalized inputs

⑤ Feed Forward

| normalized inputs |

↓

| Linear |
↓
| ReLu |
↓
| Linear |

⎱ Feed Forward Component

↓

⚫

Here We have ~~discussed~~ discussed all the
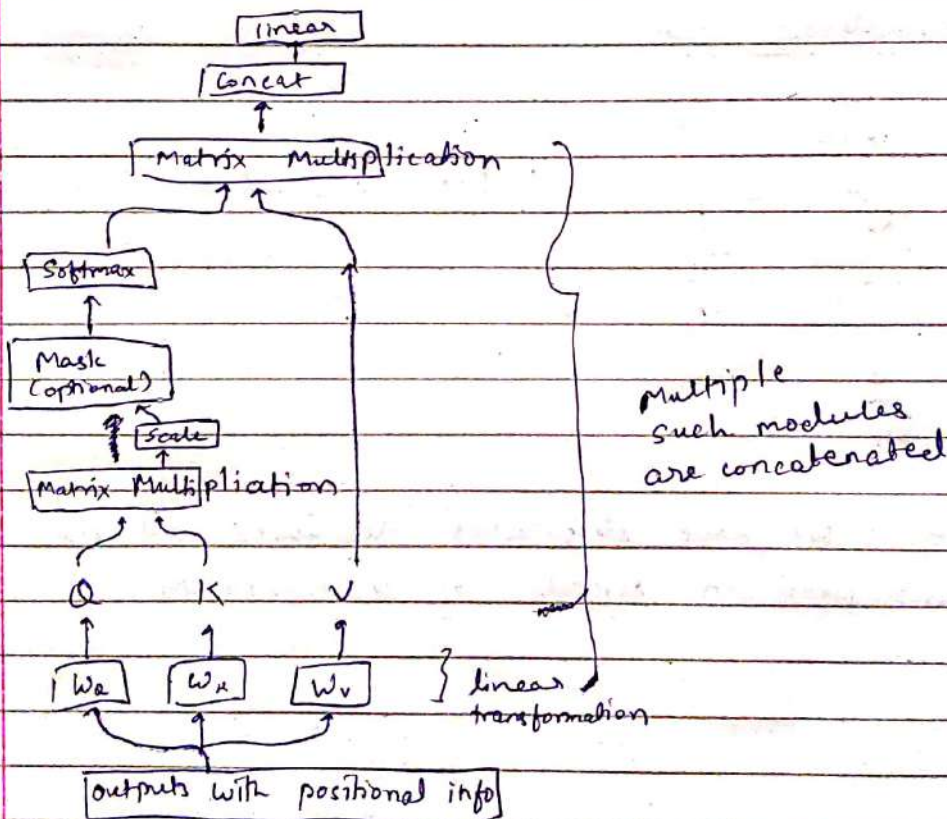components used in encoder of transformers.

→ Decoder

① Construct output Embeddings

similar to step ① of encoder

② Add Positional Information

similar to step ② of encoder

③ Masked Multi Headed Attention



multiple
such modules
are concatenated

linear
transformation

outputs with positional info

The only difference here is the Mask from the Attention structure in encoder

## Mask

Let's say we have the following score matrix after scoring

| 0.7 | 0.1 | 0.1 | 0.1 |
|-----|-----|-----|-----|
| 0.1 | 0.6 | 0.2 | 0.1 |
| 0.1 | 0.3 | 0.6 | 0.1 |
| 0.1 | 0.3 | 0.3 | 0.3 |

$+$

| $0$ | $-\infty$ | $-\infty$ | $-\infty$ |
|-----|-----------|-----------|-----------|
| $0$ | $0$ | $-\infty$ | $-\infty$ |
| $0$ | $0$ | $0$ | $-\infty$ |
| $0$ | $0$ | $0$ | $0$ |

$\underbrace{\qquad\qquad}_{\text{mask values of future}}$

$=$

|     | This | is | a | dog |
|-----|------|-----|-----|-----|
| This | 0.7 | $-\infty$ | $-\infty$ | $-\infty$ |
| is | 0.1 | 0.6 | $-\infty$ | $-\infty$ |
| a | 0.1 | 0.3 | 0.6 | $-\infty$ |
| dog | 0.1 | 0.3 | 0.3 | 0.3 |

$\downarrow Q \downarrow$

$\longrightarrow K \longrightarrow$

Since transformer is auto regressive (it predicts future based on past) we need to hide the future information so that it can learn & predict correctly.

Here in this case,

$softmax \left( \begin{array}{c} \text{Masked} \\ \text{Matrix} \end{array} \right) =$

| 1 | 0 | 0 | 0 |
|------|------|------|------|
| 0.37 | 0.62 | 0 | 0 |
| 0.26 | 0.31 | 0.43 | 0 |
| 0.21 | 0.26 | 0.26 | 0.26 |

Thus attention is not given to future because we want to predict the future.

(4)    Multi Headed Attention

→   Cross Attention
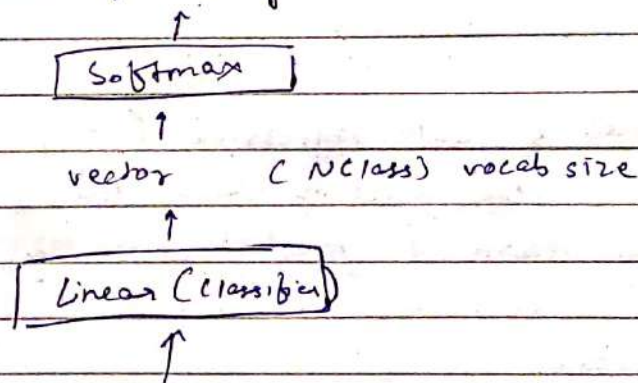
Works similar to self attention.
It is called cross attention because

v. imp   [ • key and values come from encoder
        • query come from masked multi head attention of
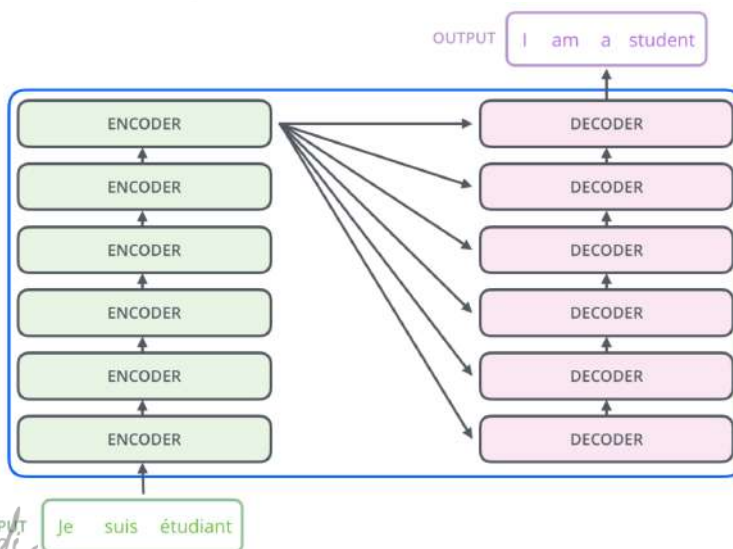          decoder

(5)   Linear Classifier

Decoder is followed by linear classifier & softmax
to give a probability distribution

probability distribution  (NClass) vocab size
            ↑
    ┌─────────────┐
    │  Softmax    │
    └─────────────┘
            ↑
    vector    ( NClass) vocab size
            ↑
    ┌─────────────────┐
    │ Linear (Classifier) │
    └─────────────────┘
            ↑

Note -

OUTPUT  I  am  a  student

ENCODER → DECODER
ENCODER → DECODER
ENCODER → DECODER
ENCODER → DECODER
ENCODER → DECODER
ENCODER → DECODER

INPUT  Je  suis  étudiant