

# Transformer Block

## Encoder Block

- Positional Embeddings
- Input embeddings
- M·H·A
- Projection after M·H·A
- Add & Norm (Residual + layerNorm) After Atten
- feed forward (two linear layers)  $N \times$
- Add & Norm (Residual + layerNorm after feed forward)

→ Repeat N times  
Since Each encoder is  $N$  times

Input & positional embeddings.

$$\left\{ \begin{array}{l} \text{Sequence length} = 1024 = d_L \\ \text{Embed dim} = 512 = d_{\text{model}} \end{array} \right\} \text{Inputs} = \mathbb{R}^{d_L \times d_{\text{model}}}$$

$$\text{attention heads} = 8 = \# \text{heads.}$$

$$\text{head dim} = 512/8 = (64) = \frac{d_{\text{model}}}{\text{heads}}$$

$$d_K = \frac{d_{\text{model}}}{\text{heads}} = \frac{512}{8}$$

$X = \mathbb{R}^{d_L \times d_{\text{model}}}$   $X$  = Input to multi head attention

Now taking the input  $X$  projecting them multi head attention

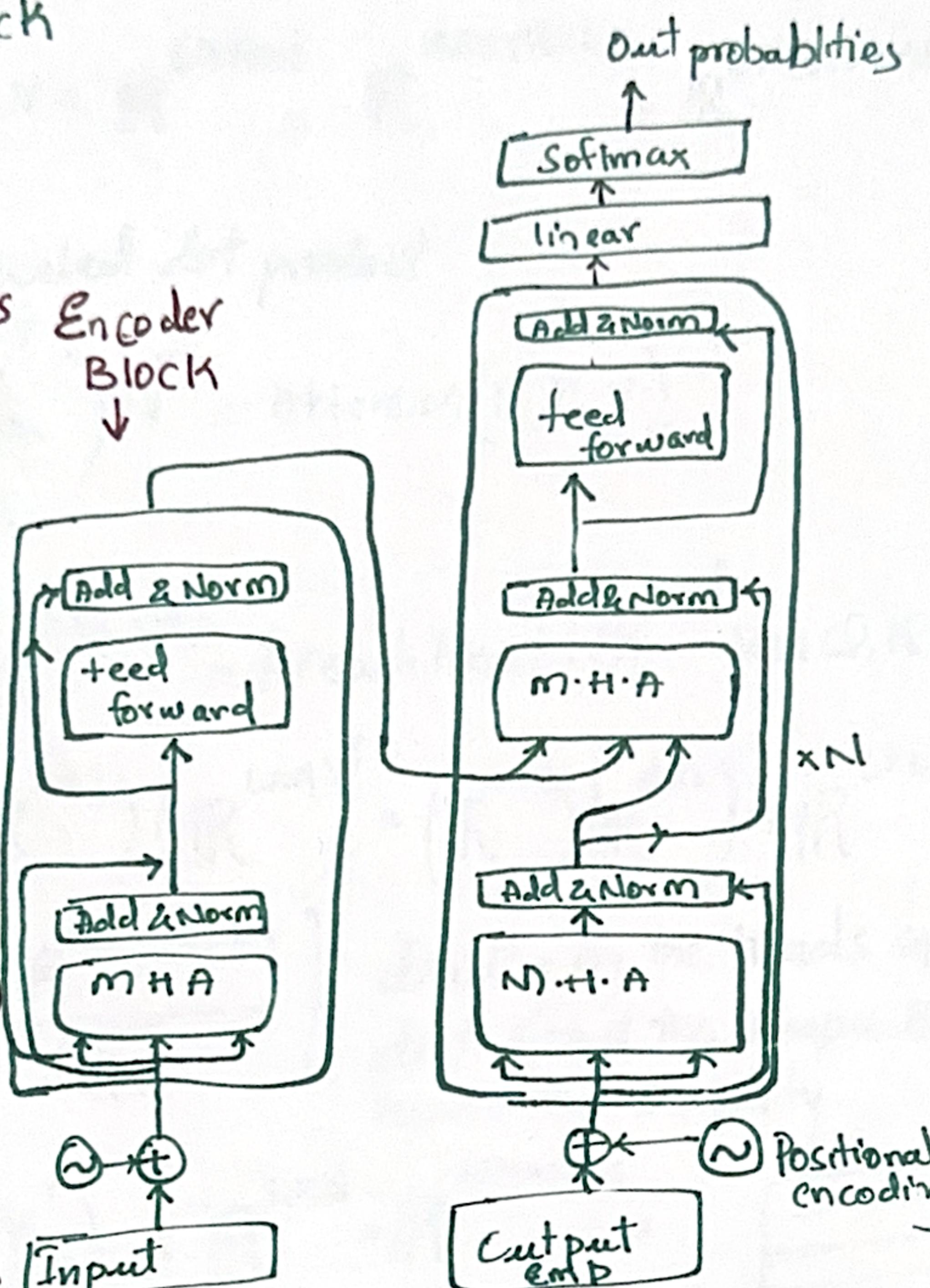
$$Q = X W_Q \quad [W_Q \quad W_K \quad W_V] = \text{learned weights.} \quad \mathbb{R}^{d_{\text{model}} \times d_K} : \text{Vanilla attention}$$

$$K = X W_K \quad [W_K] = \mathbb{R}^{d_{\text{model}} \times d_K} \quad \text{multihead} = d_K = \frac{d_{\text{model}}}{\text{no of head}} \quad \mathbb{R}^{d_{\text{model}} \times d_K} : K, Q, V$$

$$V = X W_V \quad Q, V, K = \mathbb{R}^{1024 \times 64} \quad \mathbb{R}^{1024 \times 64} : \mathbb{R}^{1024 \times 64} \quad \text{all heads} \Rightarrow \mathbb{R}^{1024 \times 64 \times 8} \quad \text{Reshape} \quad 8 \times 1024 \times 64$$

$X$  is constant or same for all  $K, Q, V$

$W$  are three different  $\mathbb{R}$  weights



$$\text{Now } Q, K, V = \mathbb{R}^{L \times h \times d} = \mathbb{R}^{1024 \times 8 \times 64} = \mathbb{R}^{8 \times (1024 \times 64)}$$

Now the scaled dot product

$$\text{Softmax} \left( \frac{QK^T}{\sqrt{d_K}} \right) * V = \text{Attention}(Q, K, V)$$

$$QK^T = \mathbb{R}^{h \times (L \times d)} = \text{foreach head} = \mathbb{R}^{L \times d} \text{ then } Q, KV \text{ becomes } \mathbb{R}^{L \times d}$$

$$QK^T = (\mathbb{R}^{L \times d}) (\mathbb{R}^{L \times d})^T = (\mathbb{R}^{L \times d}) (\mathbb{R}^{d \times L}) = \mathbb{R}^{L \times L} = \mathbb{R}^{1024 \times 1024}$$

$$\frac{QK^T}{\sqrt{d_K}} = \left[ \frac{\mathbb{R}^{1024 \times 1024}}{\sqrt{d_K}} \right]$$

dividing by the heads sqrtroot  
don't affect the shapes But added for numerical stability.

$$\text{Softmax} \left( \frac{QK^T}{\sqrt{d_K}} \right) = \mathbb{R}^{L \times d} = \mathbb{R}^{1024 \times 1024} =$$

Apply softmax  
Row By Row

Convert to probability  
Row wise

$$\left[ \begin{array}{c} \text{Row-0} \left[ \begin{array}{ccccccc} - & - & - & - & - & - & - \\ 0 & & & & & & 1023 \end{array} \right] \\ \text{Row-1} \left[ \begin{array}{ccccccc} 0 & - & - & - & - & - & - \\ 0 & 1 & & & & & 1023 \end{array} \right] \\ \vdots \\ \text{Row-1023} \left[ \begin{array}{ccccccc} - & - & - & - & - & - & - \\ 0 & & & & & & 1023 \end{array} \right] \end{array} \right]$$

1024

$$\text{Attention Score} = \mathbb{R}^{1024 \times 1024} = \text{Softmax} \left[ \frac{QK^T}{\sqrt{d_K}} \right]$$

Now we have attention scores  $\text{Softmax}\left(\frac{QK^T}{\sqrt{d_K}}\right) = \mathbb{R}^{L \times L} = \mathbb{R}^{1024 \times 1024}$

Now we do  $\text{Softmax}\left(\frac{QK^T}{\sqrt{d_K}}\right)^T = \mathbb{R}^{L \times L} * \mathbb{R}^{L \times d_K} = \mathbb{R}^{L \times d_K} = \mathbb{R}^{1024 \times 64}$   $d_K = \frac{d_{\text{model}}}{\text{heads}}$

We can do the above computations in parallel since we have divided our  $d_{\text{model}}$  by heads and then we concatenate these heads.

Concatenate heads  $(\mathbb{R}^{1024 \times 64}, \mathbb{R}^{1024 \times 64})$

$\text{Attention}(Q, K, V)$   $\Rightarrow \mathbb{R}^{1024 \times (n \cdot d_K)} = \mathbb{R}^{1024 \times 8 \times 64} = \underline{\mathbb{R}^{1024 \times 512}}$

$\text{Softmax}\left(\frac{QK^T}{\sqrt{d_K}}\right)^T \Rightarrow$

..... After Concatenating we pass this through a Projection layer. Projection after MHA

$w_p \in \mathbb{R}^{d_{\text{model}} \times d_{\text{model}}} = \mathbb{R}^{512 \times 512}$

Once we have the multi head attention output  $= \mathbb{R}^{L \times d_{\text{model}}} = \mathbb{R}^{1024 \times 512}$   
we pass this concatenated head through a linear projection  $w_p \in \mathbb{R}^{d_{\text{model}} \times d_{\text{model}}}$

$(m \cdot \text{MHA}) * w_p = \mathbb{R}^{L \times d_{\text{model}}} \times \mathbb{R}^{d_{\text{model}} \times d_{\text{model}}} = \mathbb{R}^{L \times d_{\text{model}}} = \underline{\mathbb{R}^{1024 \times 512}}$

Add & Norm (Residual Connection)

Input to MHA  $\uparrow$  Input to MHA + Output of  
Before projecting to QKV Projection After  
MHA

layer normalization don't affect the  
shapes of the output instead it just  
normalizes the values.

$(X + \text{Projection After MHA})$

$(\mathbb{R}^{1024 \times 512} + \mathbb{R}^{1024 \times 512})$  After adding we do layer normalization.  
Shape of output =  $\mathbb{R}^{1024 \times 512} = \mathbb{R}^{L \times d_{\text{model}}}$

Feedforward  
(Two linear layers)

In this layer we consider the output we got after the layernorm

$\mathbb{R}^{1024 \times 512}$

we take this output and pass it through i) linear projection Expand dimensions  
two linear layers ii) ReLU dimension remain same  
(iii) linear projection Project Back to  $d_{model}$

first linear projection:

Output from previous layer \* weight of this linear projection  
 $\mathbb{R}^{L \times d_{model}} \times \mathbb{R}^{d_{model} \times \text{Expand}} = \mathbb{R}^{\text{Expand} \times 2048}$

$\mathbb{R}^{L \times d_{model}} \times \mathbb{R}^{d_{model} \times \text{expand}} = \mathbb{R}^{1024 \times 512} \times \mathbb{R}^{512 \times 2048} = \mathbb{R}^{1024 \times 2048}$

So now we have projected things from  $512 \rightarrow 2048$  dim & then we apply ReLU activation.

This keeps the Dimension Same  $\Rightarrow \mathbb{R}^{L \times \text{expand}} = \mathbb{R}^{1024 \times 2048}$

After the ReLU we apply one more linear transformation and project the ~~weight~~ things Back to  $\mathbb{R}^{L \times d_{model}}$  dimension

$\mathbb{R}^{L \times \text{expand}} \times \mathbb{R}^{\text{expand} \times d_{model}} = \mathbb{R}^{L \times d_{model}}$   
 $\mathbb{R} \times \mathbb{R}^{\frac{||}{\text{expand}}} = \mathbb{R}^{L \times d_{model}}$

Projecting  
 $L \times \text{expand}$  Back to  $\mathbb{R}^{L \times d_{model}}$

The output of the feed forward layers is ~~is~~  $\mathbb{R}^{L \times d_{model}} = \mathbb{R}^{1024 \times 512}$

→ Add and Norm

(Residual Connection + layer normalization)

The output from the feed forward layer (i.e. second projection layer)  
 $\mathbb{R}^{1024 \times 512} = \mathbb{R}^{L \times d_{\text{model}}}$

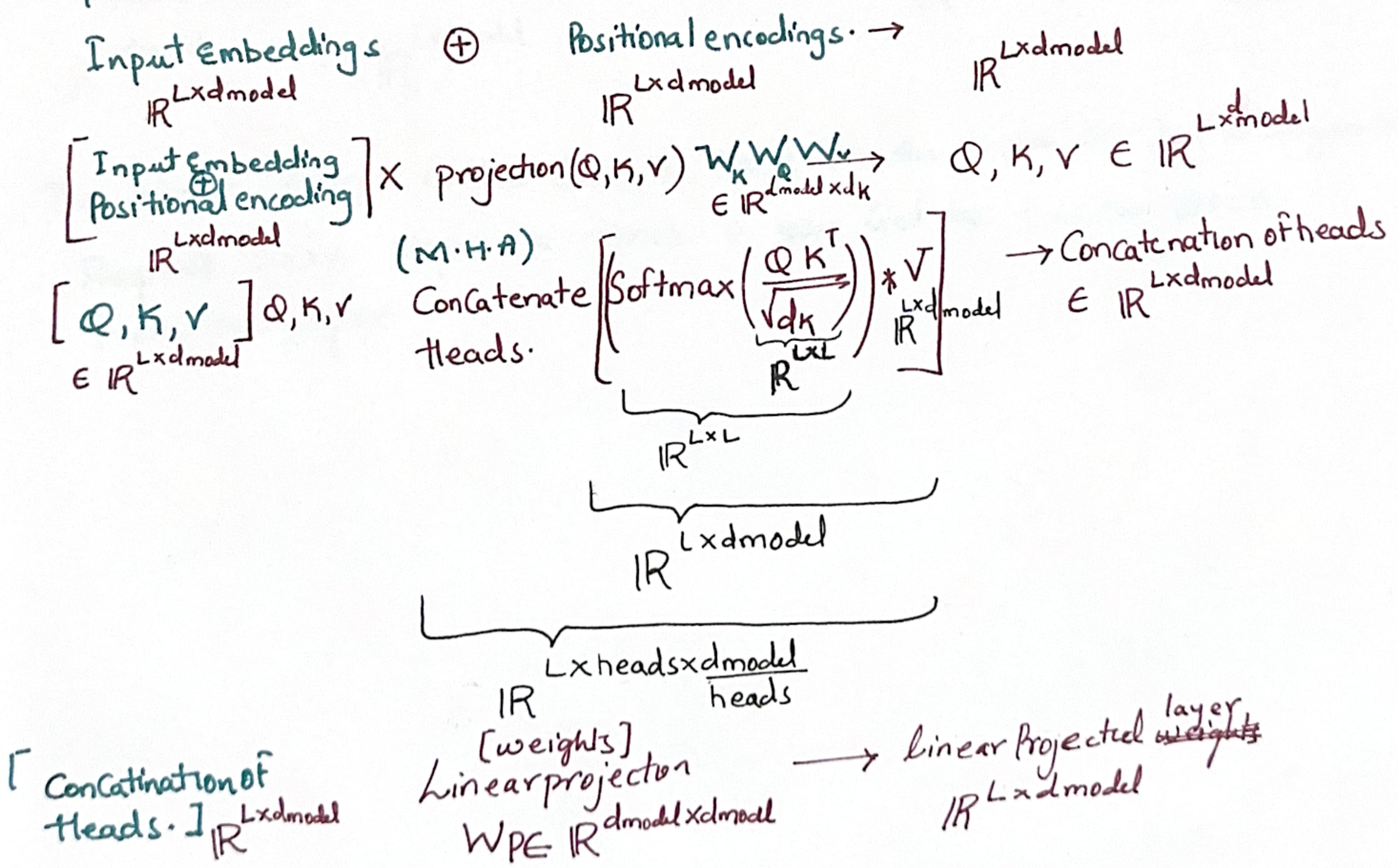
we add the  $X$  the input ~~to the m.t.h.A here~~ to the m.t.h.A here

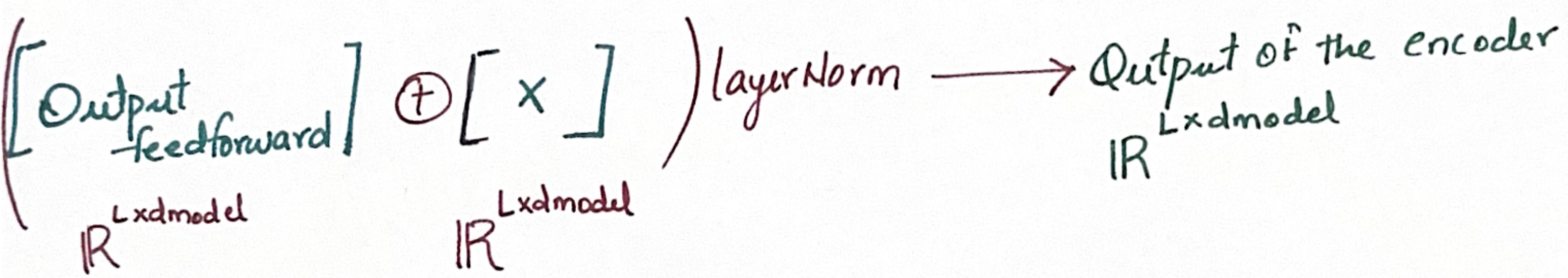
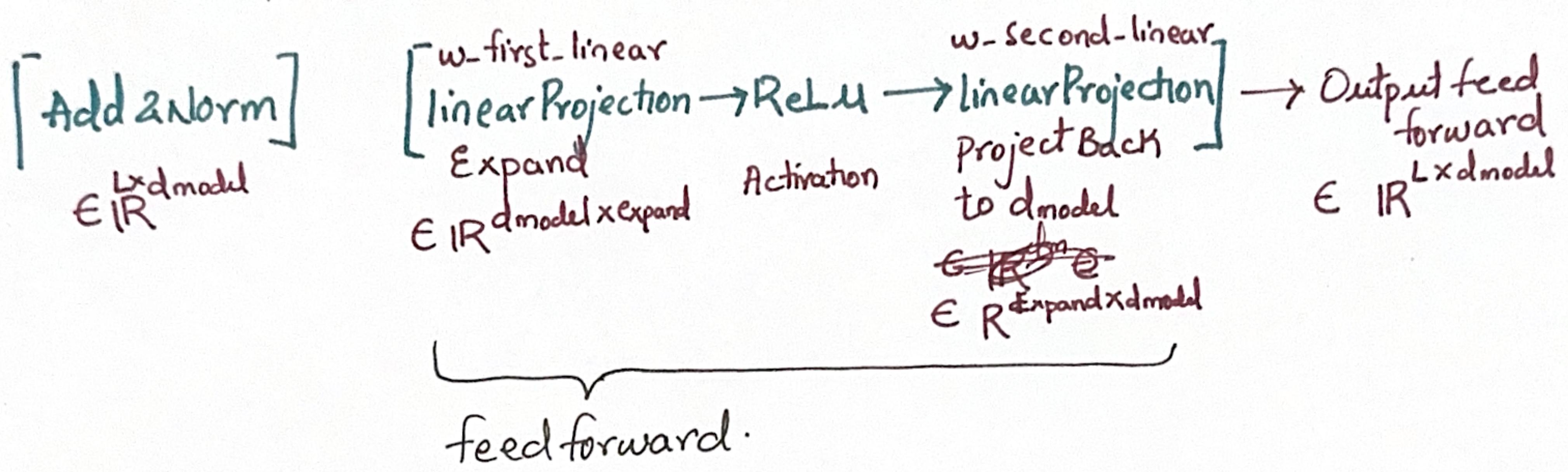
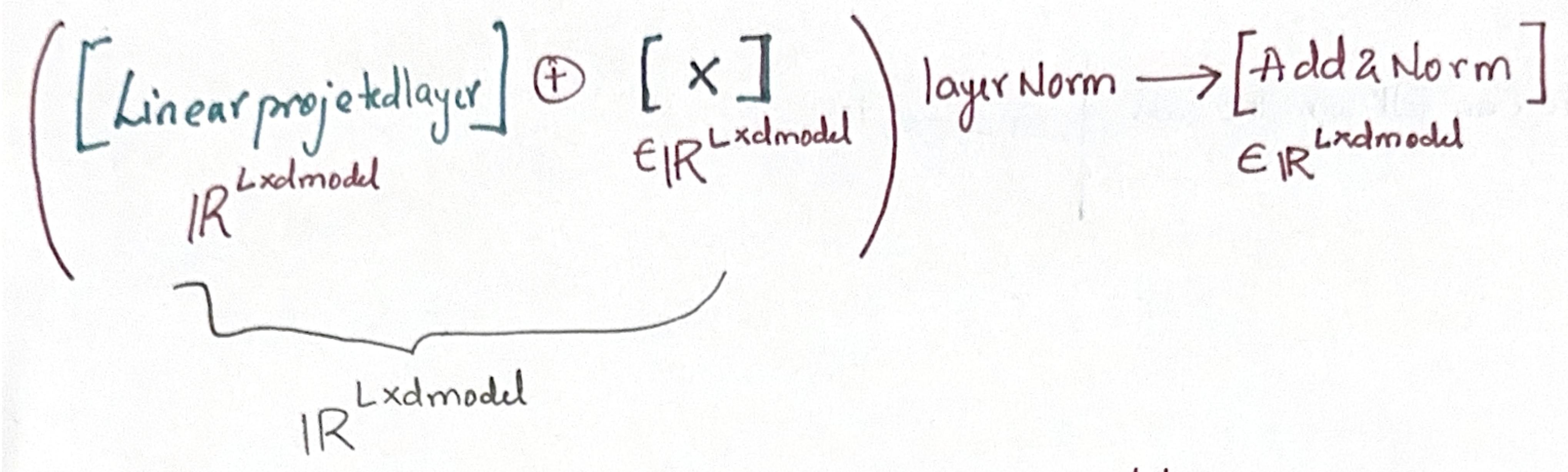
$$X = \mathbb{R}^{L \times d_{\text{model}}} = \mathbb{R}^{1024 \times 512}$$

$$\text{Second linear layer} = \mathbb{R}^{L \times d_{\text{model}}} = \mathbb{R}^{1024 \times 512}$$

$O_p = \mathbb{R}^{L \times d_{\text{model}}}$  then we perform layer norm on the ( $O_p$ ) this don't change the shape of the layer. The shape is  $\mathbb{R}^{L \times d_{\text{model}}} = \mathbb{R}^{1024 \times 512}$

The above layers and steps are the whole Encoder Block that performs operations.





The above 7 steps is the Encoder Block Form the paper Attention is All you need  
 Repeat these steps N times & you Get the Encoder Block.  
 from the paper.