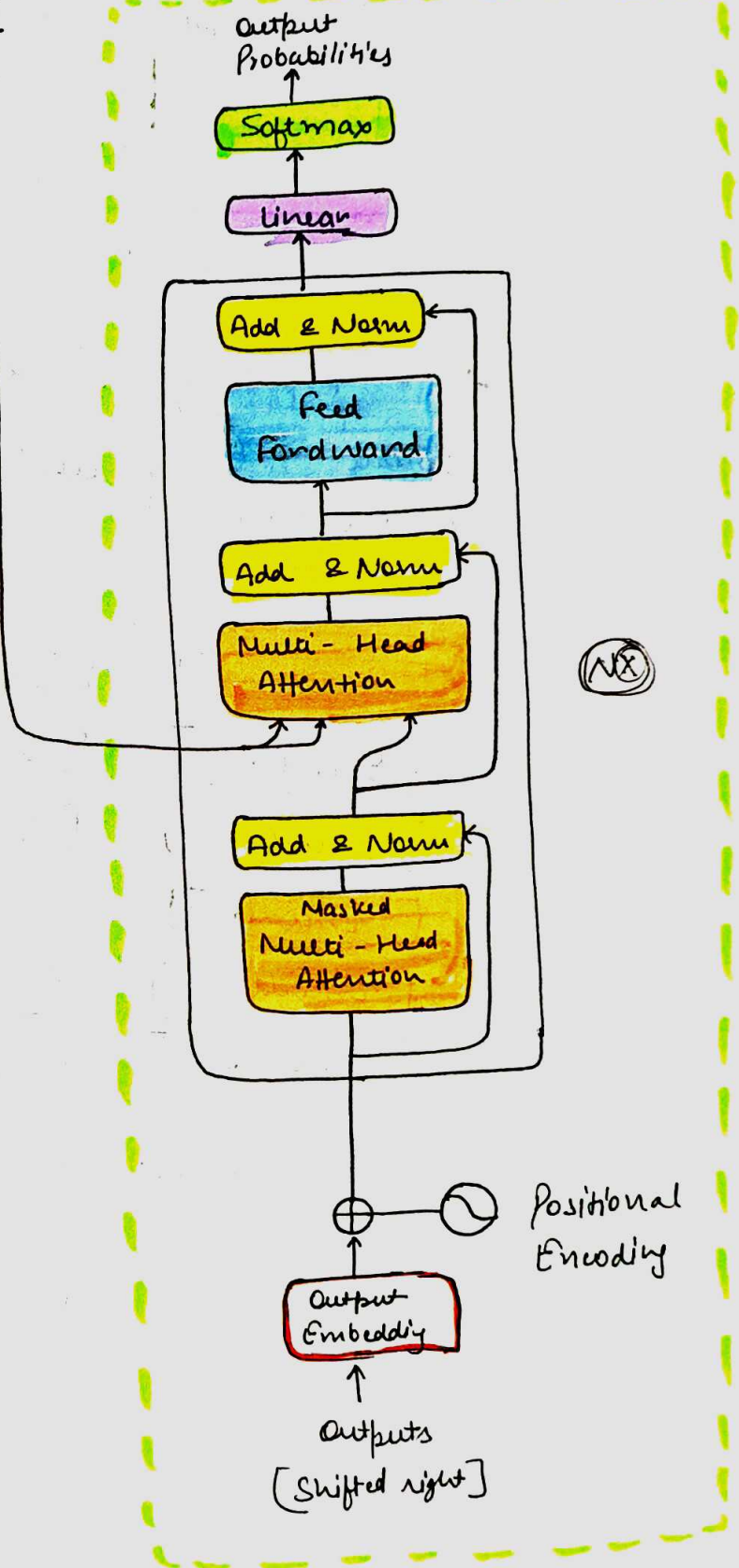
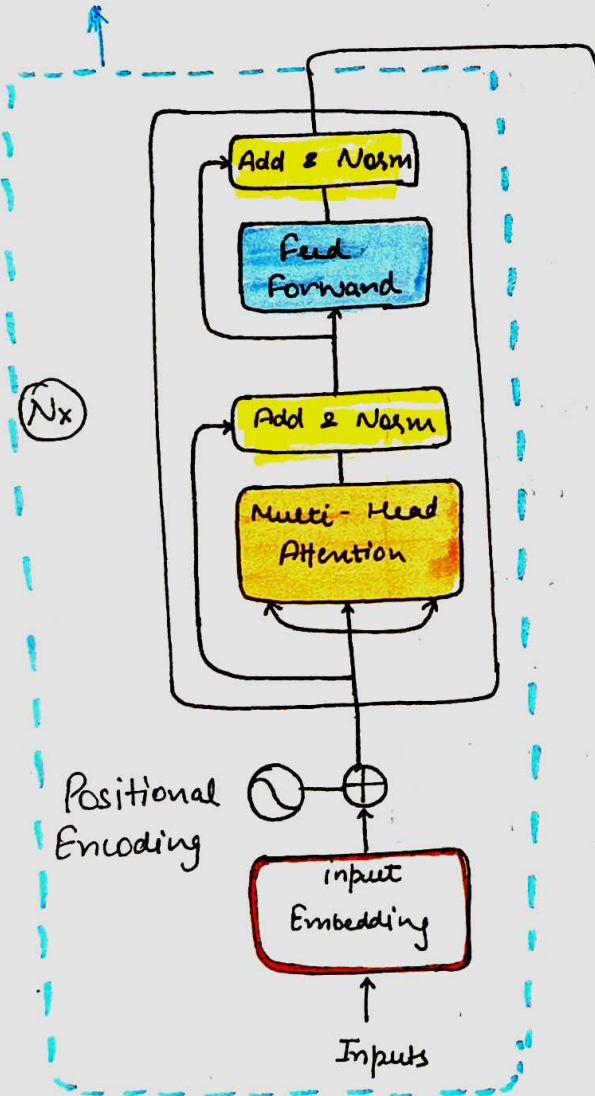


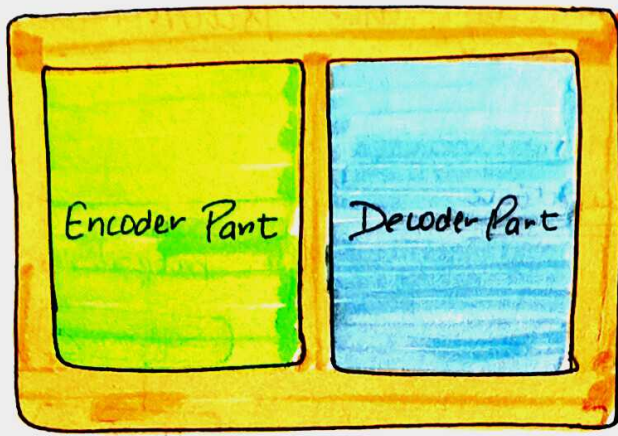
# Transformer Architecture

## Simplified Representation

decoder  
↑

Encoder

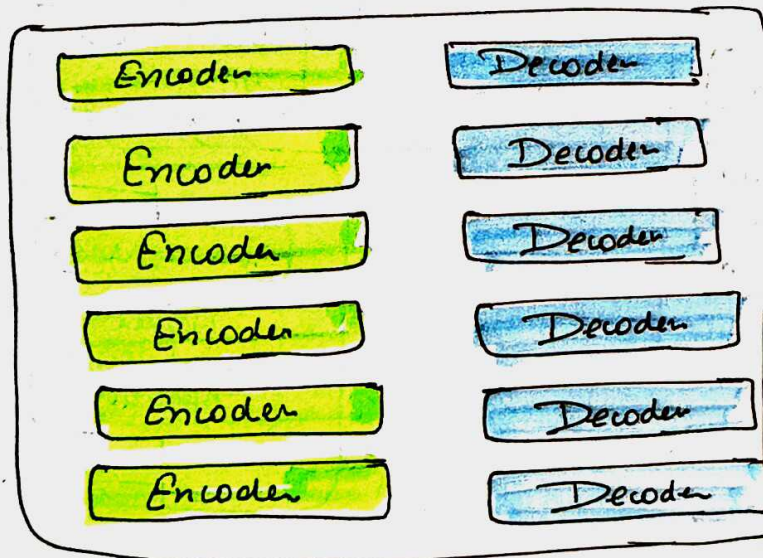




Transformer

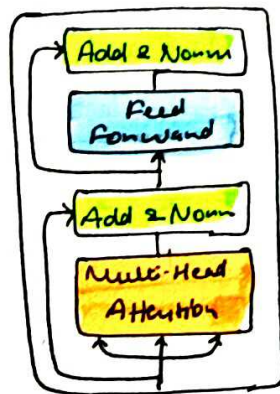
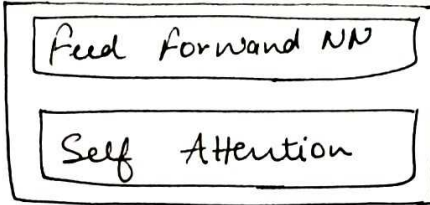
\* In Research paper of Transformation

Encoder  $\rightarrow 6$   
 Decoder  $\rightarrow 6$  } used



Transformer

Encoder



Encoder



Encoder



Encoder



Encoder



Encoder



Encoder



Input

Input



$x_1$



$x_2$



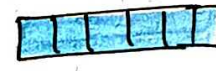
$x_3$



$p_1$



$p_2$



$p_3$

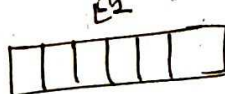


+

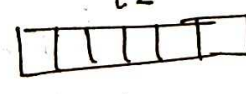
+

+

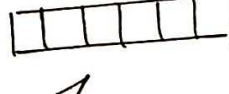
$e_1$



$e_2$



$e_3$



Embedding (512 dim)



How

are

You

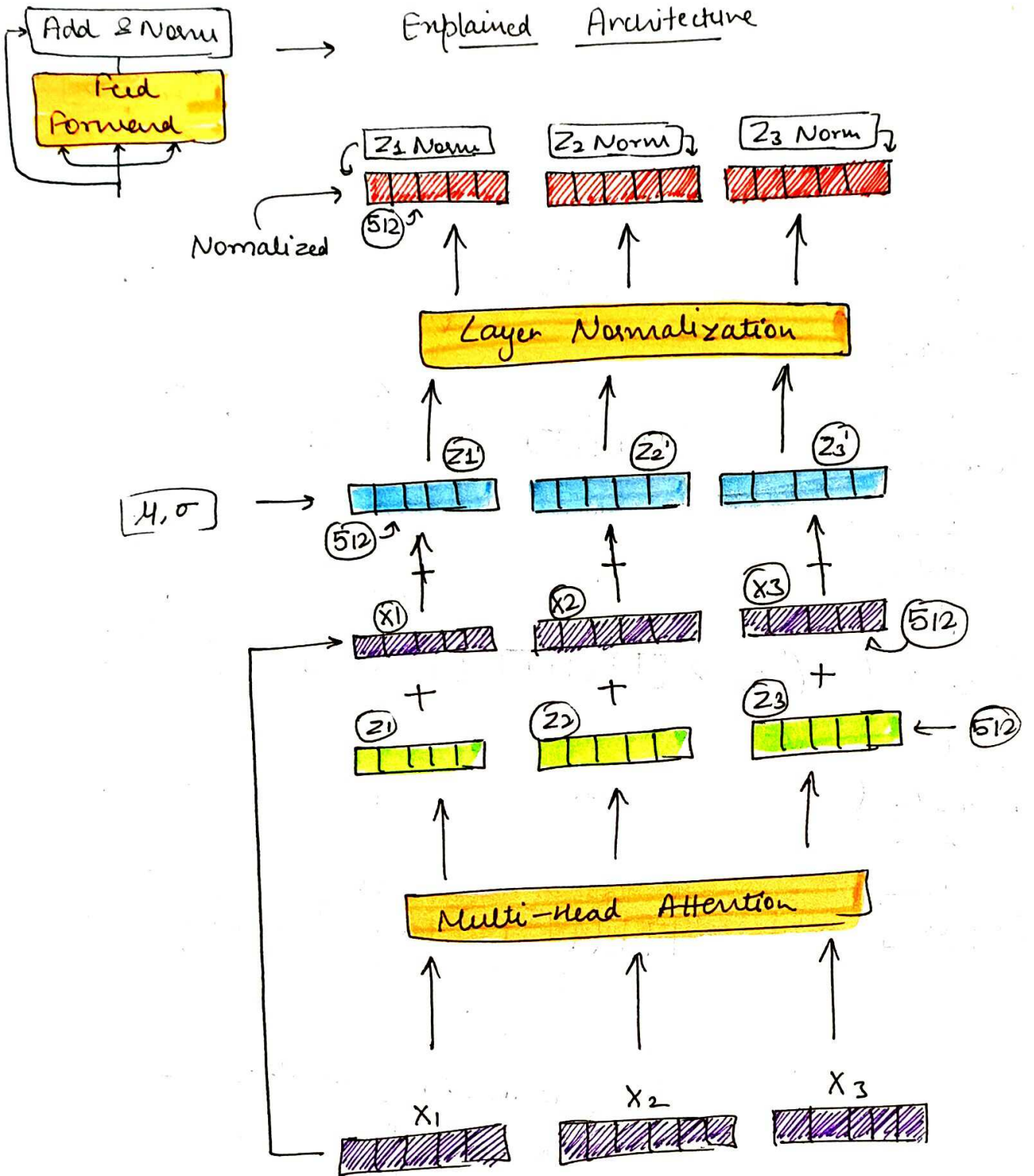
Tokenizer



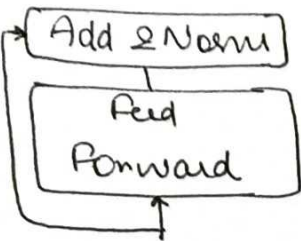
How are You



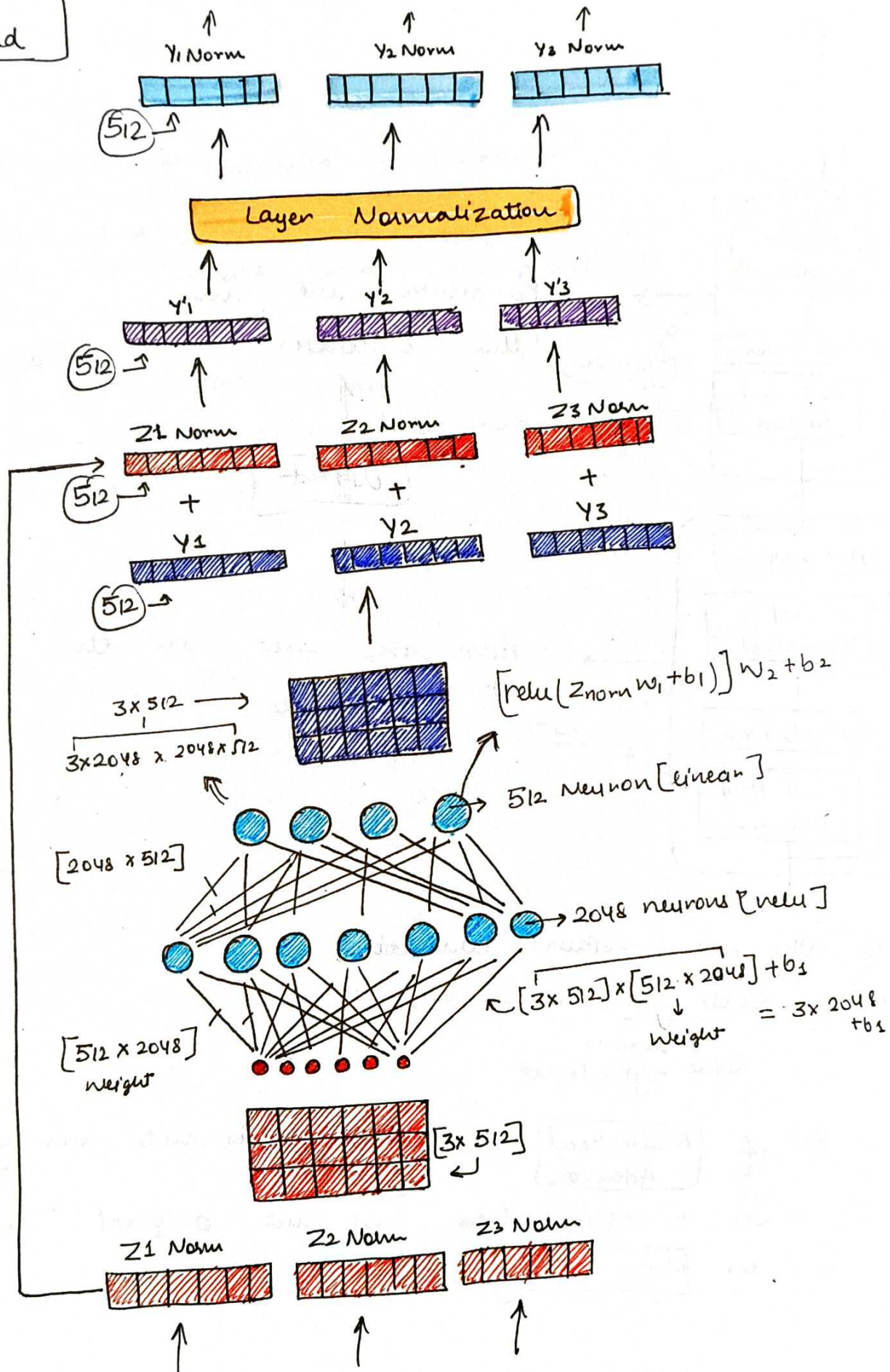
# Explained Architecture



How are you

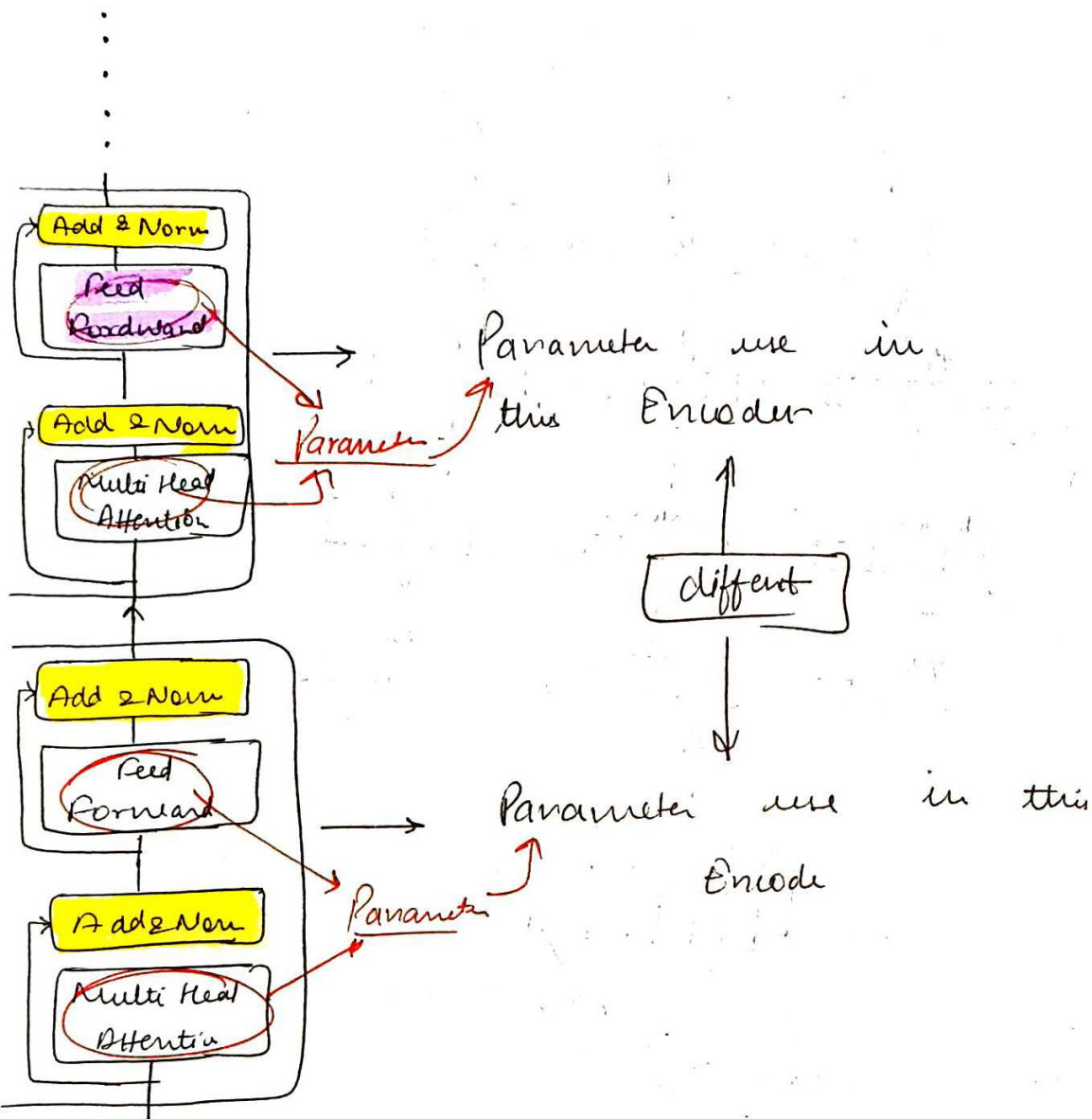


# Explained Architecture



Why increase size from 512 to 2048?

→ Because add non-linearity



Q. Why use residual connection?

Ans 1. Stable Maintaining — deep NN  
variously  
Solve a gradient

2. If Multi-Head Attention and Feed Forward is not working well on embedding then we use original embedding in Norm and Add block



2. Why use a Feed Forward Neural Network?

Sol Capture non-linearity in the process

3. Why use 6 encoder blocks?

Sol language  $\rightarrow$  complex to understand

Researchers use different no. of encoder blocks and 6 encoder give good accuracy.

Multi Encoder  $\rightarrow$  Makes easy to understand language.