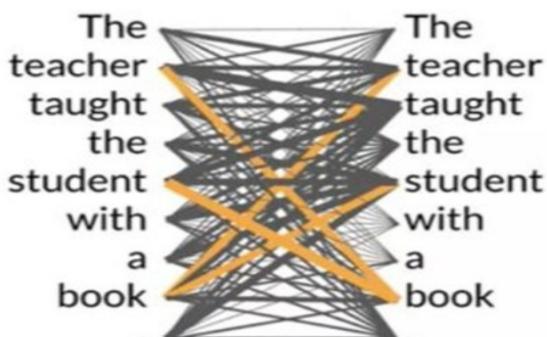
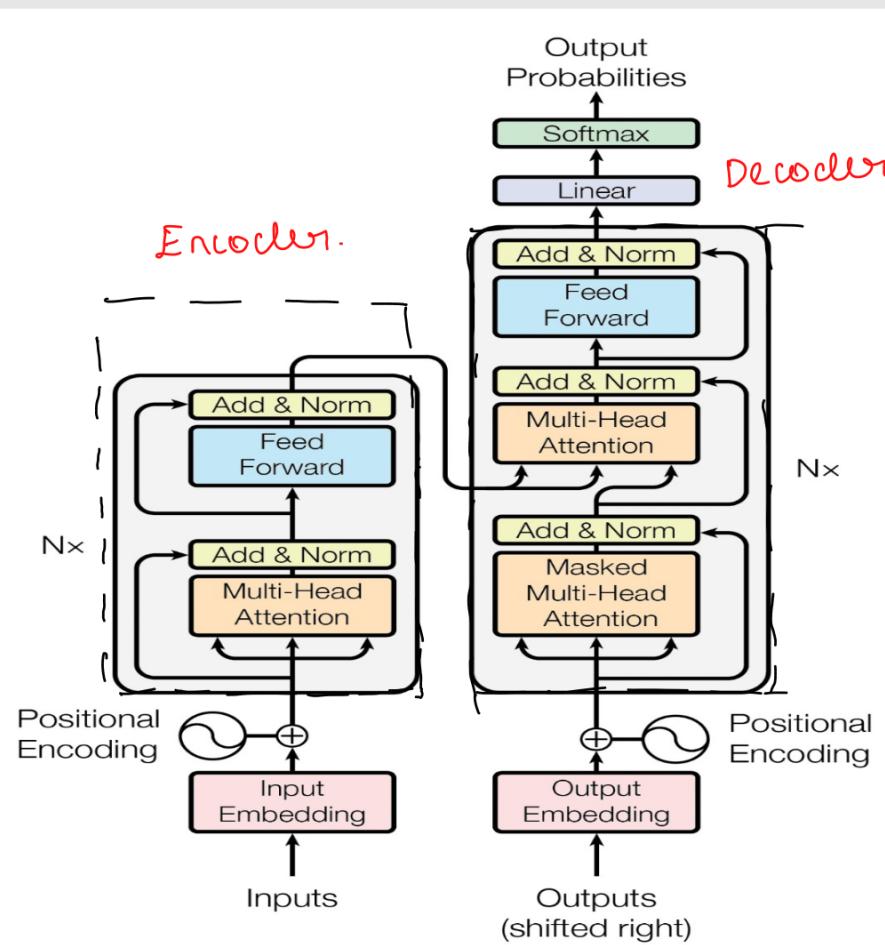


Transformer Model.

- Backbone of Large Language Model.
- The advantage of the transformer model is in its ability to understand the significance and context of every word in sentence.



Working of transformer Model !



A transformer model is a neural network that learns context and thus meaning by tracking relationship in sequential data like the words in the sentences.

1. Tokenization.

Tokenization is the process of breaking down text into smaller units, such as words / phrases for easier processing and analysis.

Raw text: "This is the first step in the NLP pipeline"

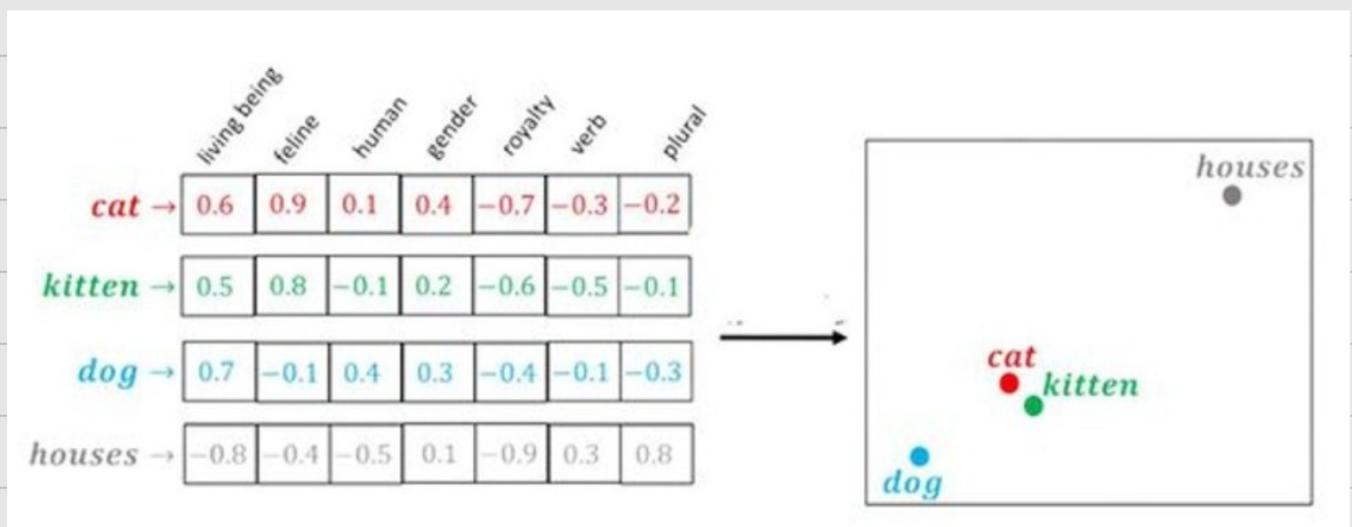


Tokenized text :

'This' 'is' 'the' 'first' 'step' 'in' 'the' 'NLP' 'pipeline'

2. Embeddings

Embedding is the process where each token is then transformed into a vector in a high-dimensional space. This embedding captures the meaning and context of each word.



③ Positional Encoding

Since transformers do not process text sequentially like RNN, they need a way to understand the order of words.

Positional encoding is the process of adding information to a model about the position of elements in a sequences.

Sequence	Index of token	Positional Encoding Matrix			
I	0	P_{00}	P_{01}	...	P_{0d}
am	1	P_{10}	P_{11}	...	P_{1d}
a	2	P_{20}	P_{21}	...	P_{2d}
Robot	3	P_{30}	P_{31}	...	P_{3d}

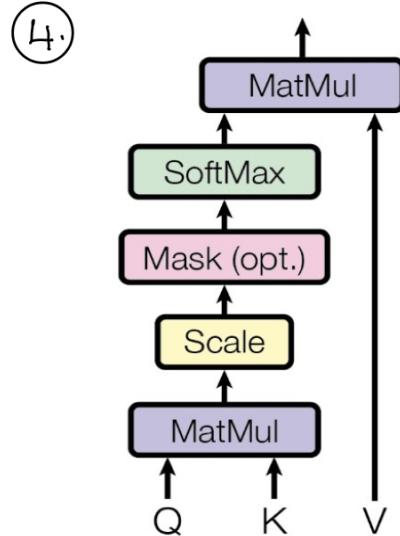
Positional Encoding Matrix for the sequence 'I am a robot'

④ Self Attention

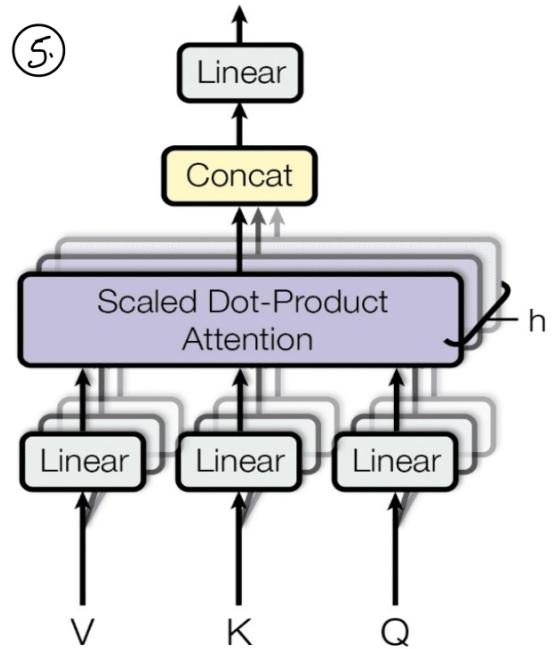
The model calculates attention scores for each word, determining how much focus it should put on other words in the sentences when trying to understand a particular word.

This helps the model capture relationship and context within the text.

Scaled Dot-Product Attention



Multi-Head Attention



⑤ Multi-Headed attention.

Multi-Headed attention is a mechanism in transformer that runs several self-attention processes in parallel, allowing the model to focus on different parts of the input sequences from different perspective at the same time.

⑥ Output

The final layers of the transformer convert the processed data into an output format suitable for the task at hand, such as classifying the text or generating new text.