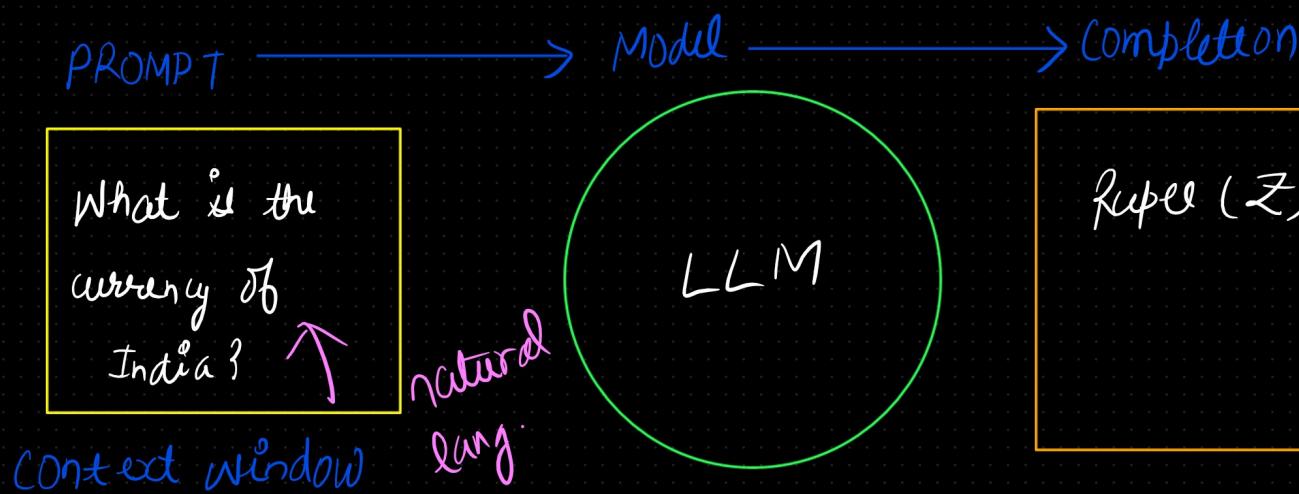


Large Language Model (LLM)

What is LLM?

LLM are machine learning models that have learned from massive datasets of human-generated content, finding statistical patterns to replicate human-like abilities.



This process is called
"Inference"

ChatGPT → Application



LLM Model → GPT3, GPT4, GPT5

Google - Gemma, Gemini

What are the Use Cases of LLMs?

- Chatbots
- Writing → Essay, Email, Report
- Summarization → Summarize long content into a meaningful shorter length.
- Code
- Language Translation
- Information Retrieval
- Augmented LLM

TRANSFORMERS



Model Architecture



2017 → Research paper



"Attention is All You Need"



revolutionized generative AI

How was text generation done before
Transformers?

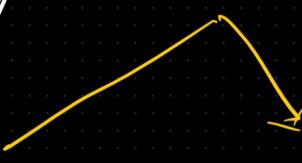
Recurrent Neural Network
(RNN)



predict next words

→ The next word was predicted looking
at the previous few words

Eg → I took my money to the bank.



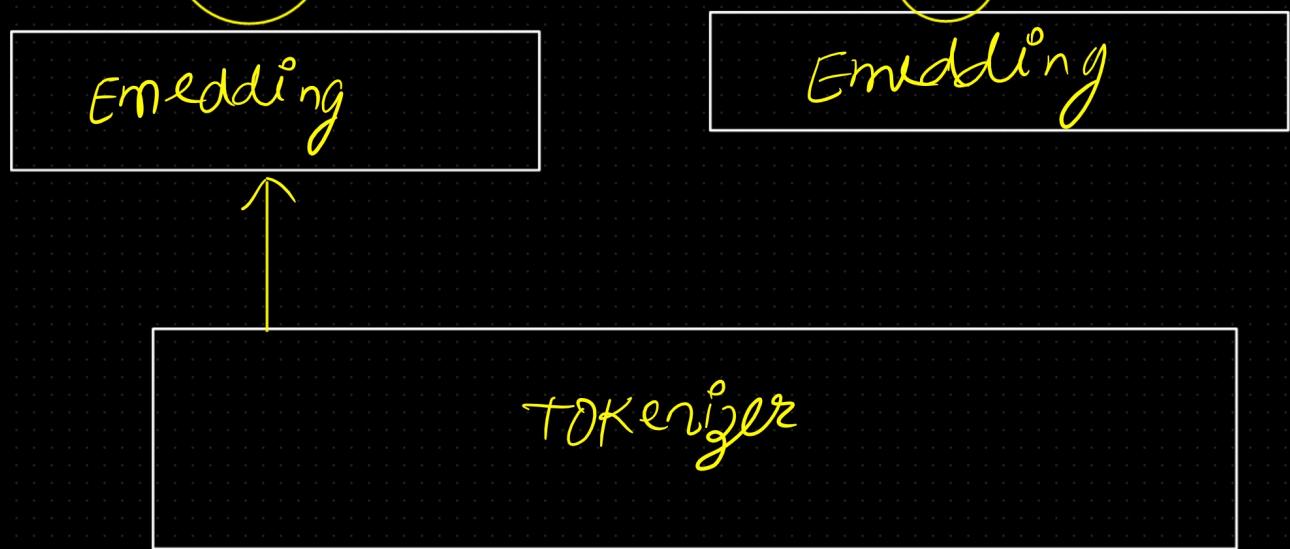
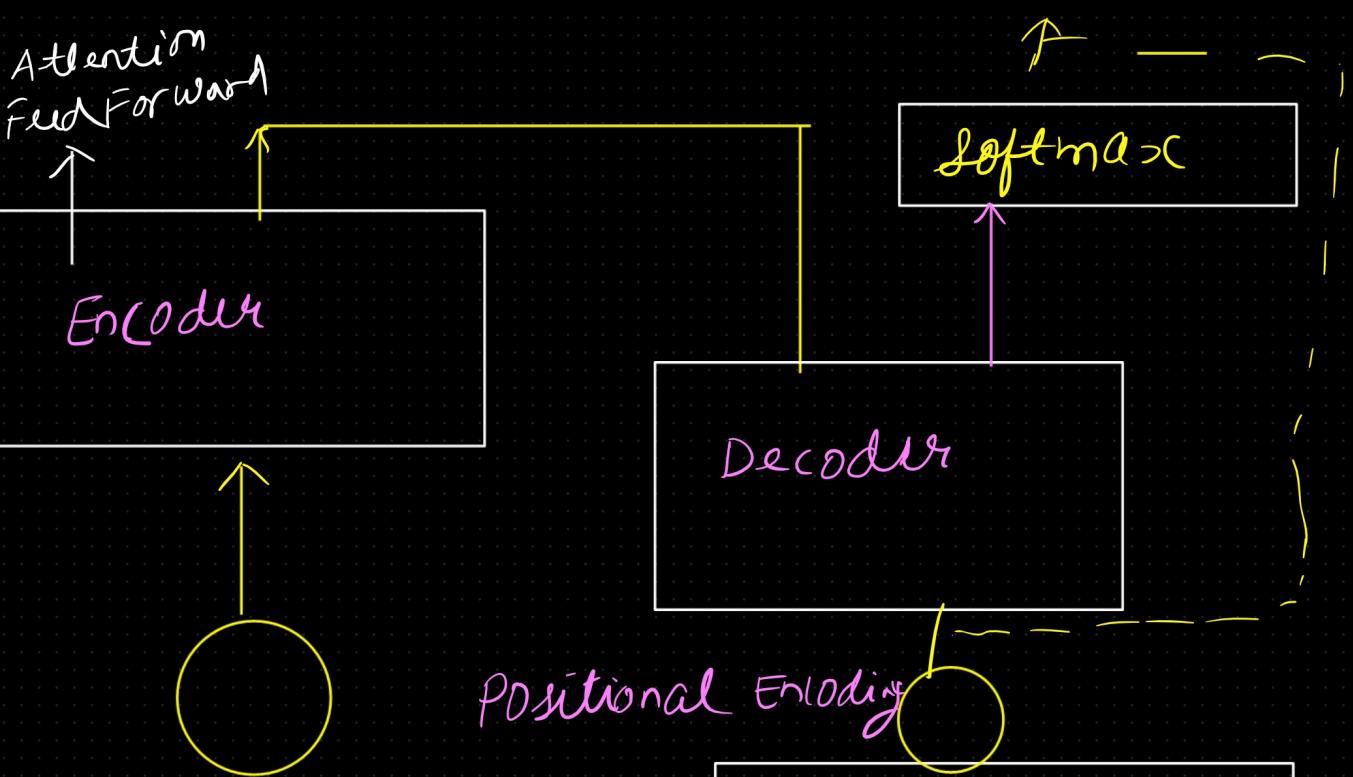
River Bank? Financial
Bank?

→ Transformers are able to pay
Attention to the meaning of the
words.



Multi-Head Attention

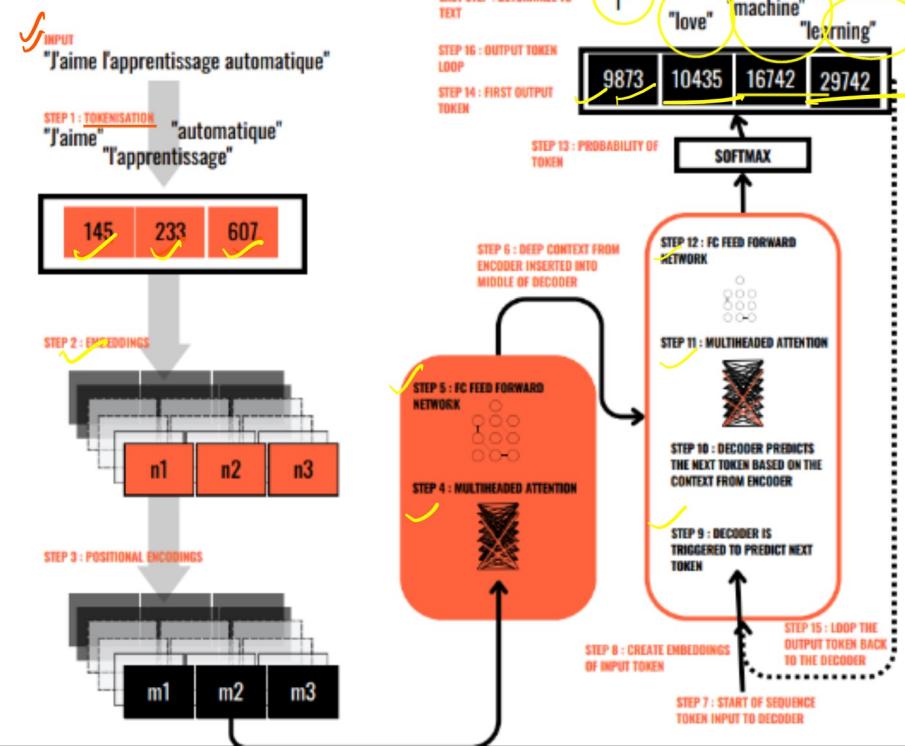
Output



↑ (words \rightarrow Token \rightarrow vector)
Input

How does a Transformer generate text?

The original objective of the transformers architecture was for **Language Translation** in form of a sequence-to-sequence task



What is prompt ?

→ The natural language instruction
in which we interact with LLM
is called prompt.

→ The construction of prompts is called
prompt engineering

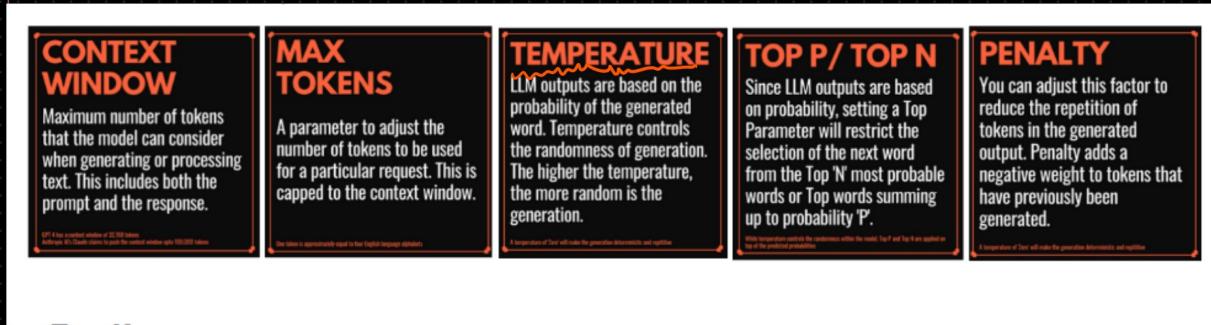
zero shot → Without giving an example

one shot → single example

Few shot → More than one example

The inferencing that an LLM does and completes the instruction given in the prompt is called 'in context learning'

↓ 0 - 1



Generative AI Project LifeCycle.

