

DAY-11

Introducing Open source LLM models;

Open source LLM;

- i) not hosted anywhere; need to download the model & load the model manually
- ii) Advantage - no charge
- iii) Disadvantage - You need good configuration system
 - i3, i5 processor
 - 8GB of RAM
 - GPU would be an add on

Some very popular & powerful open source LLMs

- Meta Llama2
- Google PaLM 2 { Backend of Google Bard }
- Falcon

[Google Search open lms github]

↳ This repo has details of all
open source LLM

Llama2 is better than gpt3.5 turbo → meta claims

Search for =) llama2.ai

Quantized Model -

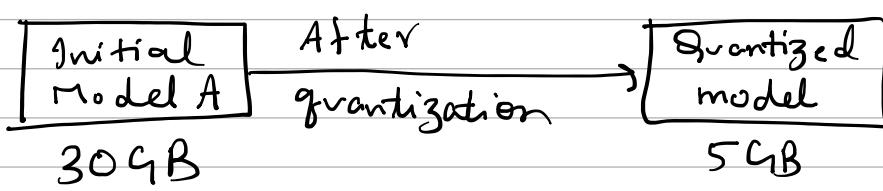
→ For a simple neural network there are multiples of parameters which are in float that gets trained /adjusted during Backpropagation (BP).

Data Type	32 bit	64 bit
short	1 byte	1 byte
String	2 byte	2 byte
int	4 byte	4 byte
long	4 byte	4 byte
float	8 byte	8 byte

} out of all datatype float consumes lot of memory / space during BP.

Therefore the process of converting float to integer is called QUANTIZATION.

$$\begin{aligned} w_1 &= 1.23 \xrightarrow{\text{After quantization}} 1 \\ w_2 &= 2.45 \xrightarrow{\text{After quantization}} 4 \end{aligned} \quad \left. \begin{array}{l} \text{ofcourse there is a} \\ \text{data loss but this is} \\ \text{done to quantize the} \\ \text{model} \end{array} \right\}$$



} The performance of quantized model is less compared to raw model

Input text → Pipeline → Response

{
① Apply preprocessing - Auto Tokenizer

② Text to numbers conversion

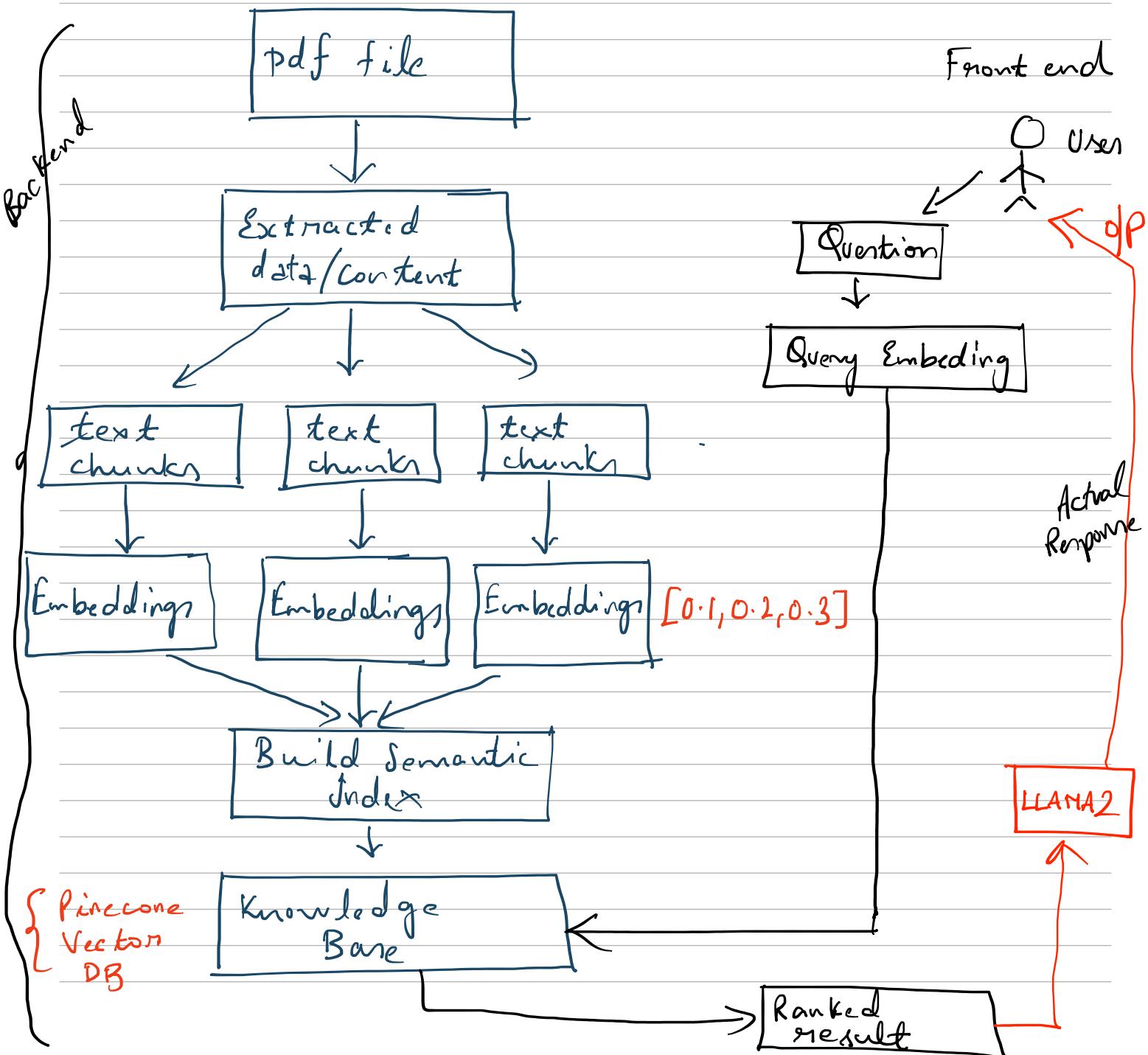
③ Model → prediction

④ Response }

02/01/2024

DAY-12

Implementing a Medical chart bot with custom data ;



Tech Stack used:

- ① Python
- ② LangChain / LlamaIndex → Gen AI framework
- ③ Frontend / web app → Flask
- ④ LLM → Meta Llama2
- ⑤ Vector DB → Pinecone