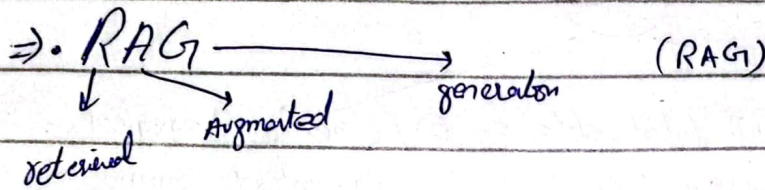


⇒ Lecture #36: Introduction to RAG System

- new technology
- Retrieval Augmented Generation (RAG)
- Quick earning potential



• LangChain

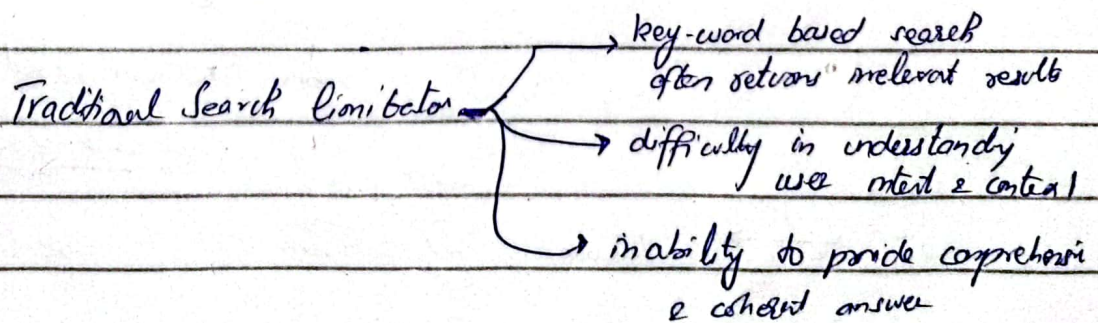
! & (S) retrieval

Bridge b/w your data source & ML/model/application

- sort of an orchestrator
- python framework helps in implementing RAG.

The Challenge of information

- digital age led to explosion of information
- making it difficult to find relevant & accurate data



So we need solution for better relevant and accurate retrieval of data.

Problem with LLM (providing information in desired output way is better)

- Not access to up-to-date information
- not access to real time & latest domain specific information

So we need LLM with our data

Solution

① Fine tuning with latest data

ایک ماڈل جو کسی specific data پر ٹرین کیا گیا ہو۔
 تو وہ ماڈل میرے ڈیٹا پر بھی اچھے نتائج دے کرے گا۔

Problems :- LLMs are very large and their training is very costly & time consuming

② Retrieval-Augmented Generation / in-context learning

ہم LLM کے پاس اپنی data دے دیتے ہیں۔
 اس سے LLM کو اس data کا analysis کرنا پڑے گا۔

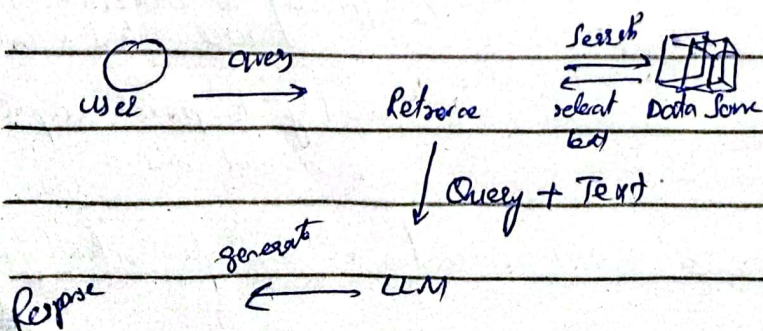
" RAG first retrieves relevant document/data from a knowledge base & then use that info

to generate response "

- giving our own context
- Can also fix to only our data, not the LLM.

Combine

retrieval-based & generation-based

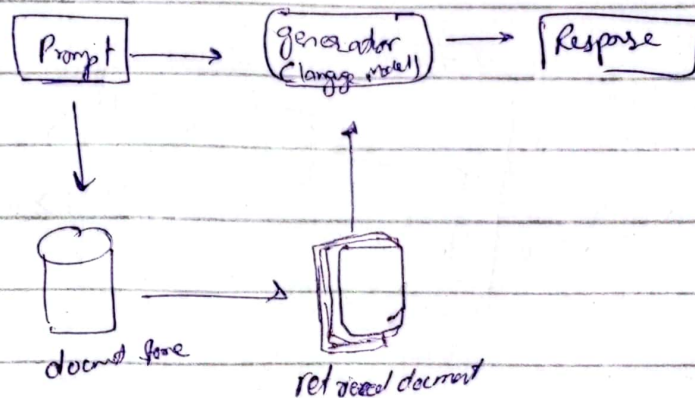


End goal - At the end of series you will be able to build any application

- 1) How data is attached to our system.
- 2) LLM also work similar to like this. (difference chatGPT is trained on short data)

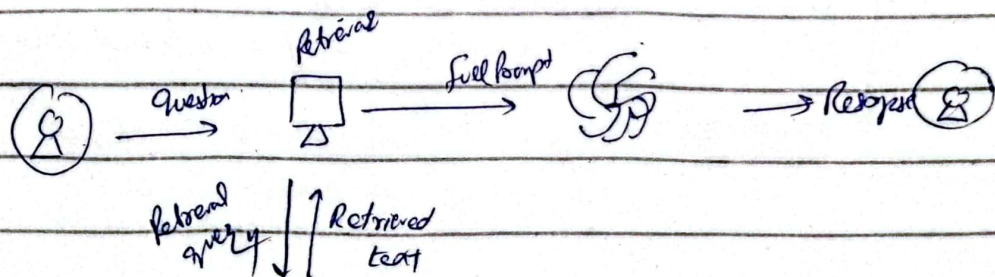
⇒ How RAG Works

- ① Query Input: we ask a question
- ② Document Retrieval: identify relevant document
- ③ Response Generator: generator crafts a response using retrieved information
- ④ Final Output: a factually accurate & contextually relevant answer.



→ Response on our data

→ only want to find generate response on our data.

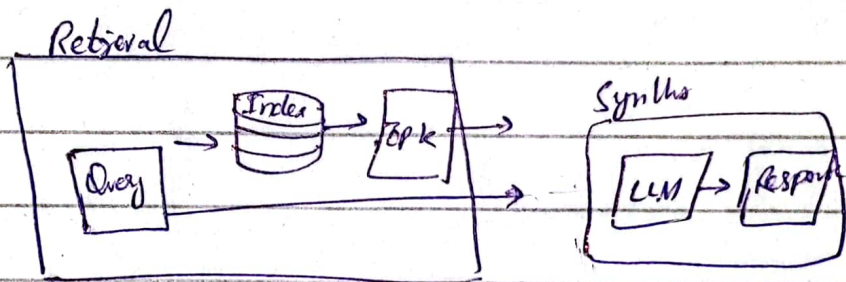
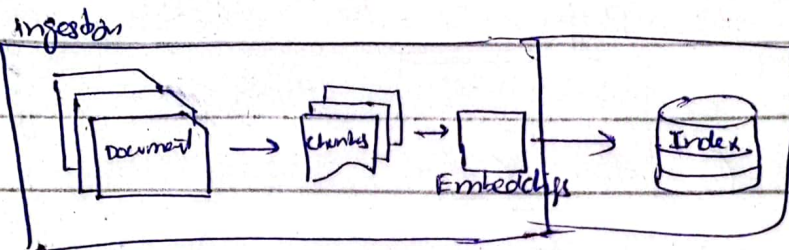


RAG Pipeline

(i) Ingestion: feeding raw data - collecting, ingesting & processing data
(give the data in easy words)

(ii) Retrieval: correctly retrieval only desired information

(iii) Synthesis: So modified readable response.



Why RAG

- ① Combine both world (retrieval based + generation based)
- ② Response are accurate & enriched (more control)
- ③ Contextual Awareness
(maintain query & retrieved document)
- ④ Improve the performance of baselines
(offers all abilities of LLM)
- question answering, summarization & more

intro to LangChain

→ Refer to the documents

- framework for building applications using LLM
- components :- include tools - like document transformers
 - connects to various data sources
 - chaining multiple queries

Ingestion ✓

- ↳ Document Loading :- process of ingesting document into system
- supported multiple format plaintext, HTML, PDF

Practical Examples (learn & Experiment)

- performance of RAG depend on ingestion
- installing Langchain library, langchain community

Loaders

components designed to ingest & pre-process data from various source

example, pdf (.pdf)

• text (.txt)

• document (.docx)

• urls

pypdf

- different loaders are with different nature.
- can add and edit meta-data