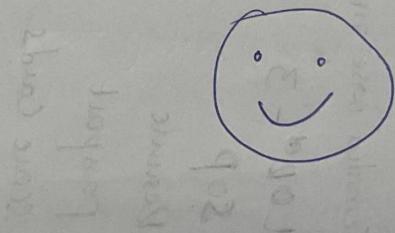


What is RAG (Retriever Augmented Generation)

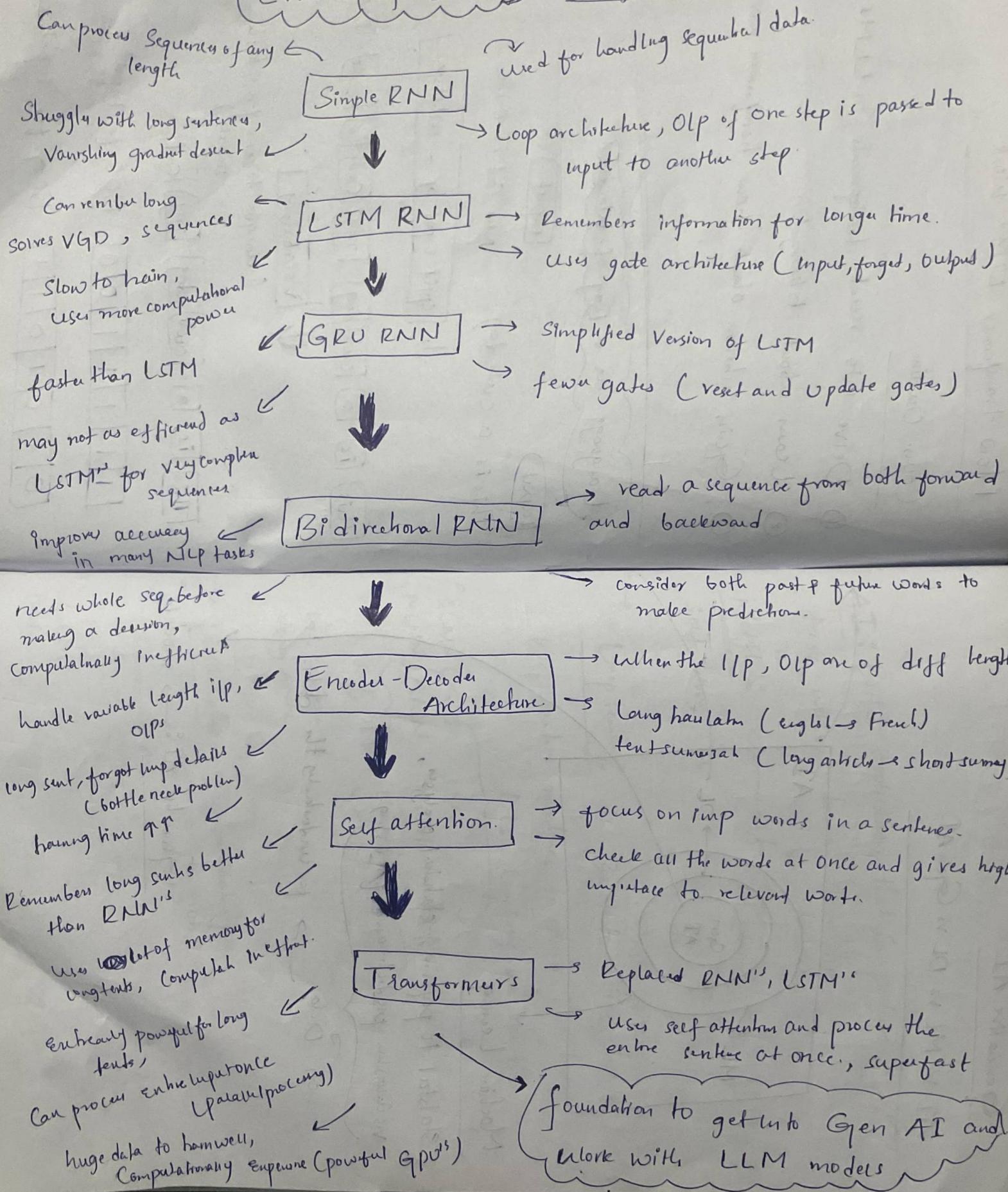
is a method where an AI model first finds relevant information from a data sources (like documents) and then uses that information to give a better and more accurate answer.

Lets Start Gen AI



Evolution of LLM's

Large Language Models :)



Generative AI

AI vs ML vs DL vs Gen AI



Machine Learning

Statistical tool to perform statistical analysis,

visualizations, prediction and forecasting.

on Data

for understanding the data.

Generative AI: → models

AI: Build applications that can perform its own tasks without human interruption.

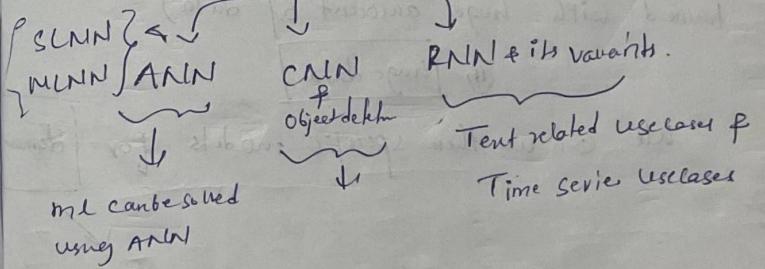
Netflix: recommendation system

self driving car

Deep learning: is built to mimic human brain.

We wanted machines to think how humans think

DL (majorly we have three)

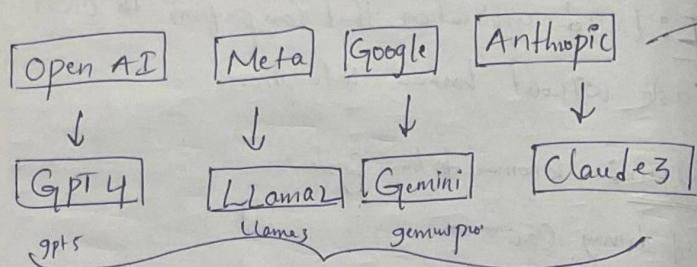


→ Discriminative → classification prediction.
Trained on label datasets

→ Generative AI → Generates new data, trained on huge data sets.

LLM (Text) LIM (Images)

Large language Models (LLMs)



foundation models

④ also called as

pretrained models (because these are

trained with huge amount of data)

↳ huge data

We can use these specific models for domain specific use cases also.

to use these pretrained models for domain specific use cases we should

Fine tune, it with custom data set

→ that is related to specific domain you are interested to work on.

Companies are in the race, to make the best LLM model

Companies and their LLMs

LangChain is the framework, which we can able to work with all LLM's models

Internet Data

↓ (Websites, articles, Wikipedia, books, etc..)

Stage 1: Generative pre-Training

Base GPT model

Stage 2: Supervised fine tuning (SFT)

Fine-tuned ChatGPT model

Stage 3: Reinforcement Learning HP

(human feedback)

ChatGPT model.

How ChatGPT is trained

Stage 1: Generative pre-training

Internet huge input text data

↓
Transformers

Tasks
↓

lang translation.

Text summarization.

text completion

Sentiment analysis

↓
Base GPT model

↓
what we want

↓
Conversation Chat Bot

req ⇒ Conversation history

res ⇒ Best ideal response

Reinforced learning through human feedback.

Stage 3

←

SFT ChatGPT → res req ← human agent

↳ alternate responses.

↳ response ranking

↳ Ranked responses

↳ Reward model →

↳ proximal policy optimization (PPO)

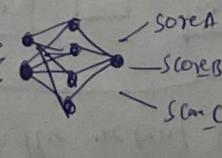
Conv hist

A

B

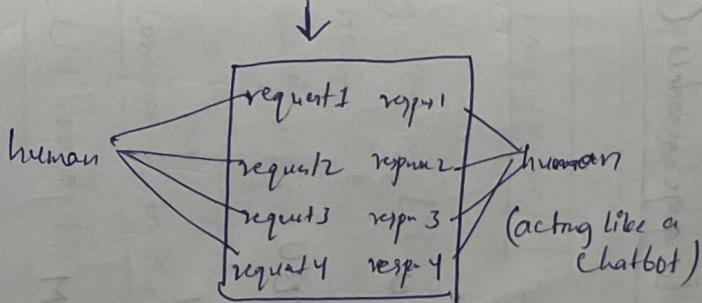
C

D



Stage 2: Supervised Fine tuning (SFT)

Real conversations.



↙ SFT training Data Corpus

req	res
req	res
req	res

⇒ millions of records

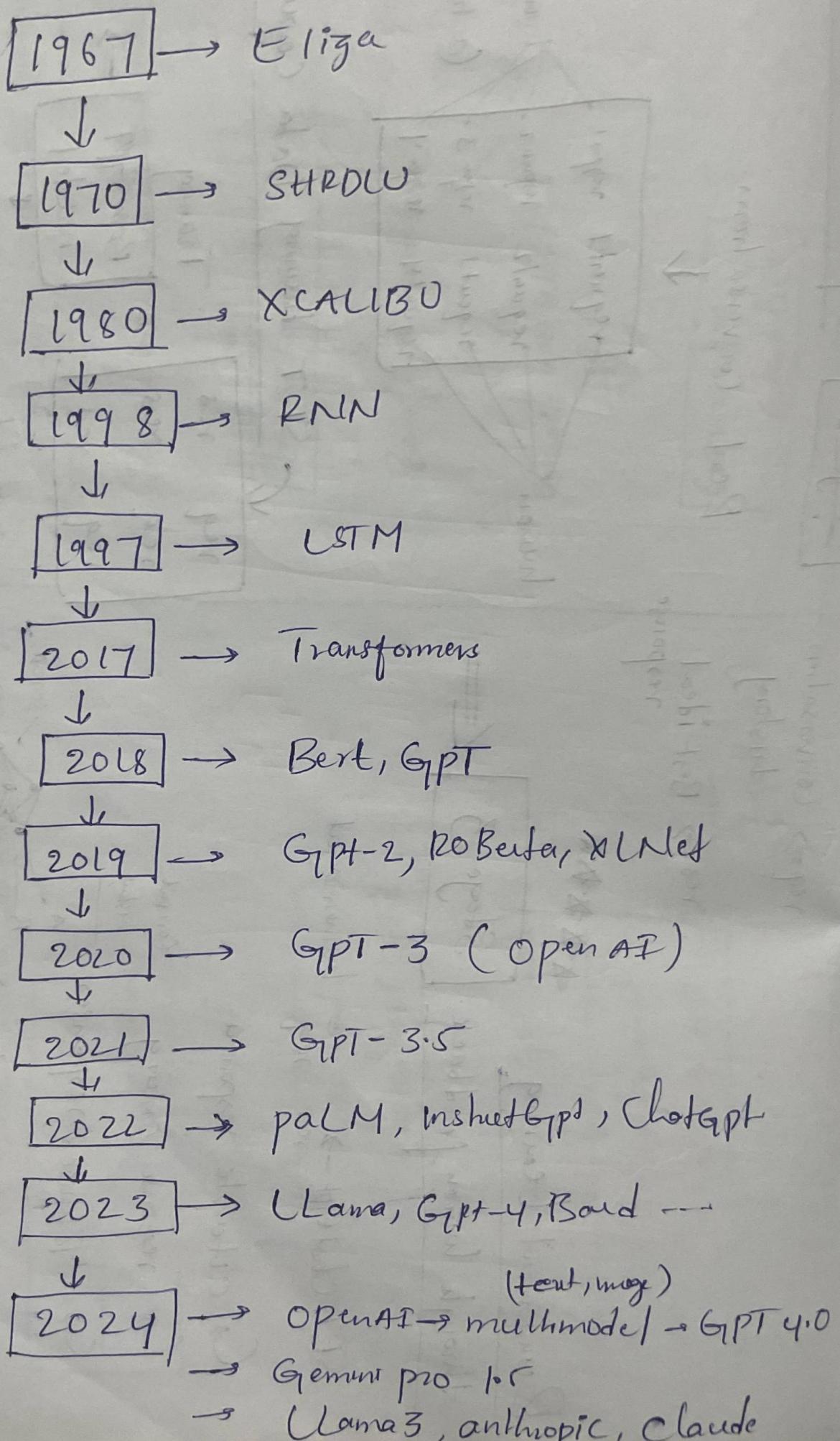
Training Base GPT model

⇒ chatGpt model

SGD optimizg

Evolution of Large language models (LLMs)

Lot of research happening in the recent years.



LLM model Analysis

artificialanalysrs.ai



Comparing various LLM's that are present right now.

(Quality, Speed, price)

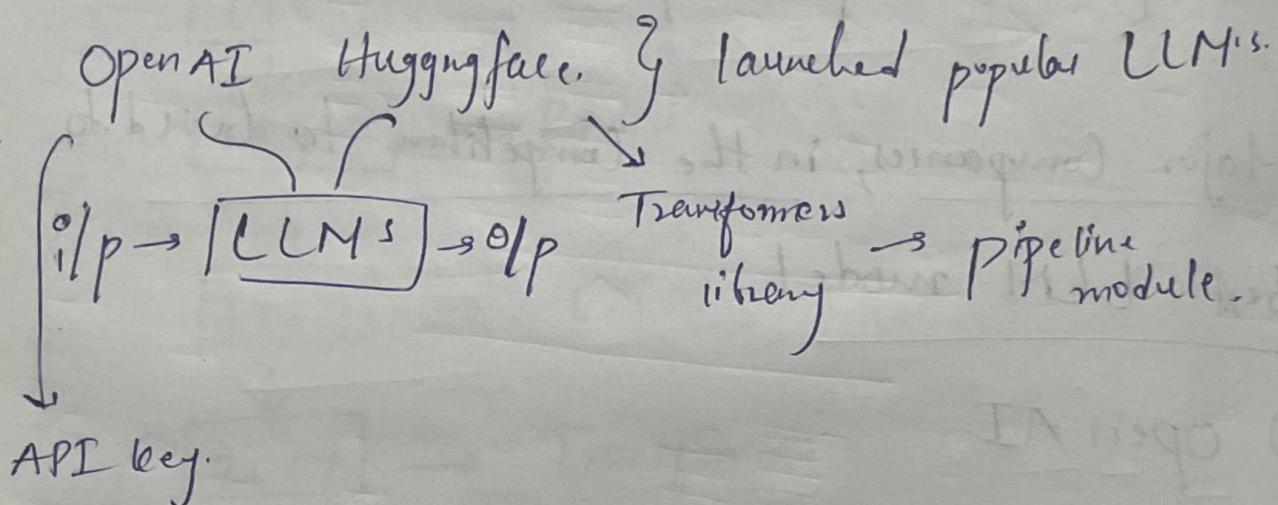
Major Companies, in the competition to build best LLM model

- ① open AI
- ② Google
- ③ Meta
- ④ MishaIAI
- ⑤ Anthropic
- ⑥ Cohere
- ⑦ Databricks
- ⑧ Reka
- ⑨ A121 labs
- ⑩ DeepSeek
- ⑪ snowflow
- ⑫ OLAi

Langchain Ecosystem.

- ① Why Langchain
- ② Langchain ecosystem.

Langchain right now is most common framework, that is specifically used to build GenAI applications.

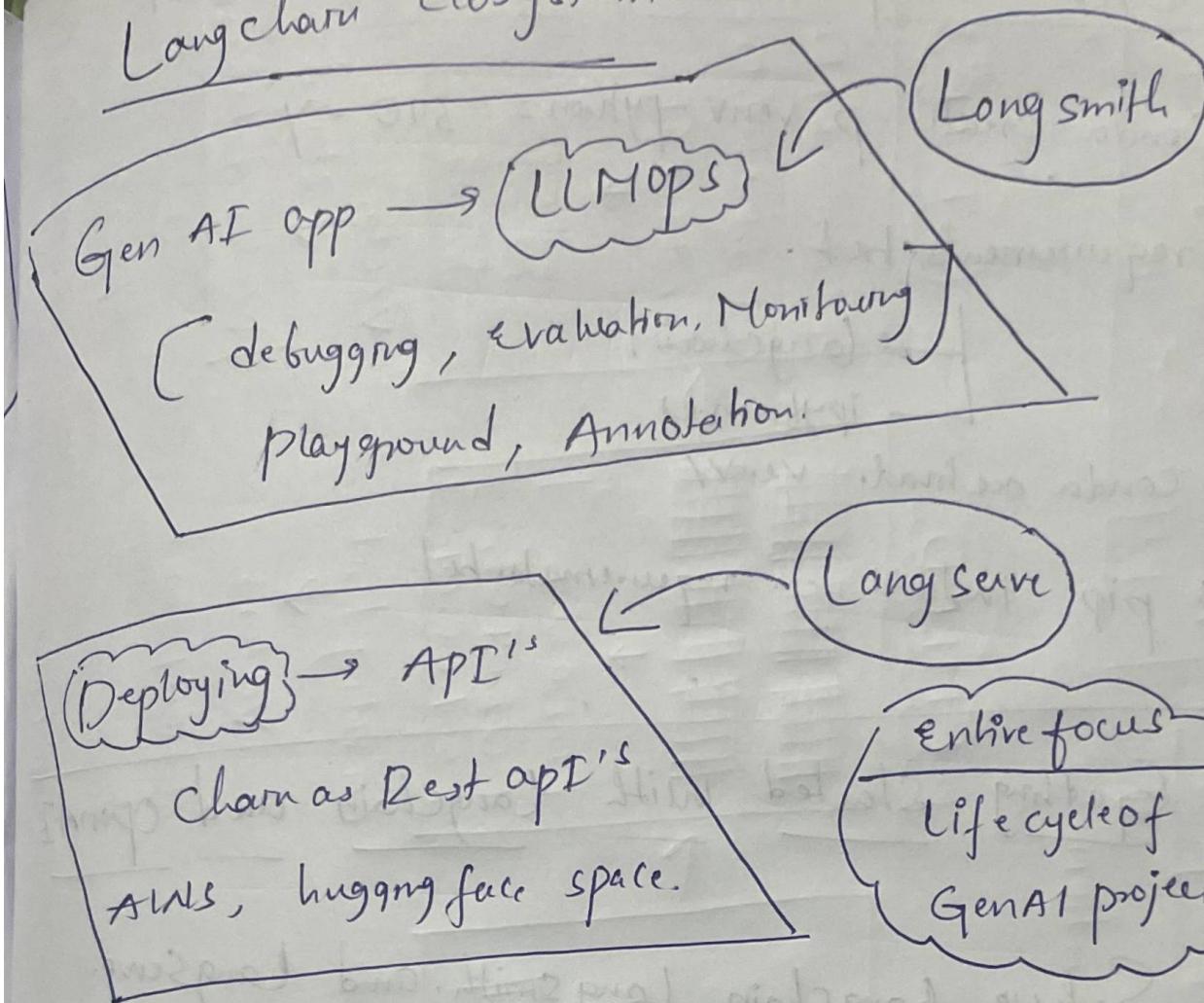


Similarly we have Meta, Google etc for accessing each companies model (there is going to be an different library).

Langchain (open source)

→ provided a common framework for accessing all LLM's models of different companies and for developing GenAI applications.

Langchain ecosystem.



Langchain → cognitive architectures

Chains agents Retrieval strategy.

Langchain Communities → Integration Components.

Model I/O
model, prompt,
example selector
output parser

Retrieval
Retriever,
DOC loader
vector store
text splitter
embedding model

agent Tooling
tool
toolkit

Langchain Core → protocols.

LCEL - Langchain Expression Language.

Parallelization
Fallbacks
Tracing
Batching
Streaming
Async Composition

CREATING VIRTUAL ENVIRONMENT

→ Conda create -p venv python==3.10 -y.

→ requirements.txt

```
└── langchain  
    └── ipykernel
```

→ Conda activate venv/

→ pip install → requirements.txt.

Getting Started with Langchain and OpenAI

- ① Set up Langchain, LangSmith, and LangServe.
- ② Use basic and common components of Langchain:
prompt templates, models and output parser
- ③ Build simple application with Langchain
- ④ Trace your application with Langsmith
- ⑤ Serve your app with LangServe.

Create langchain api key

openai api key

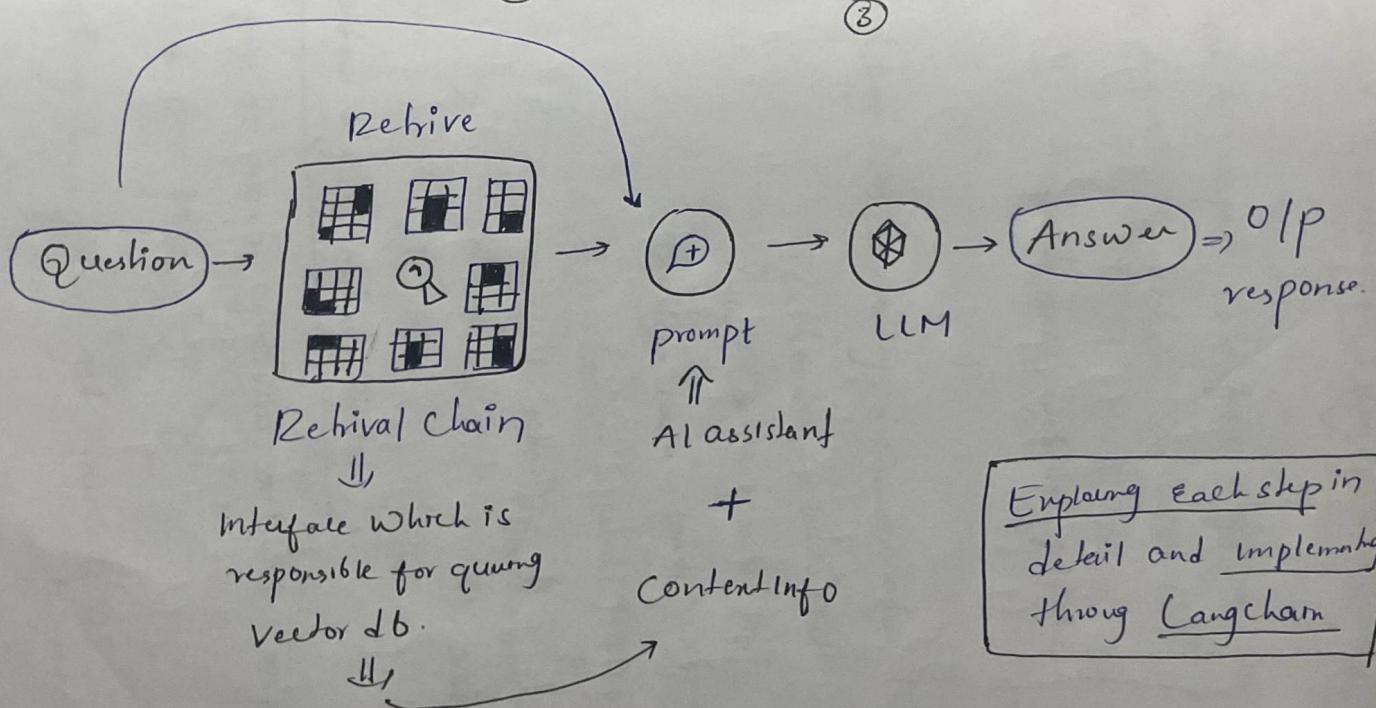
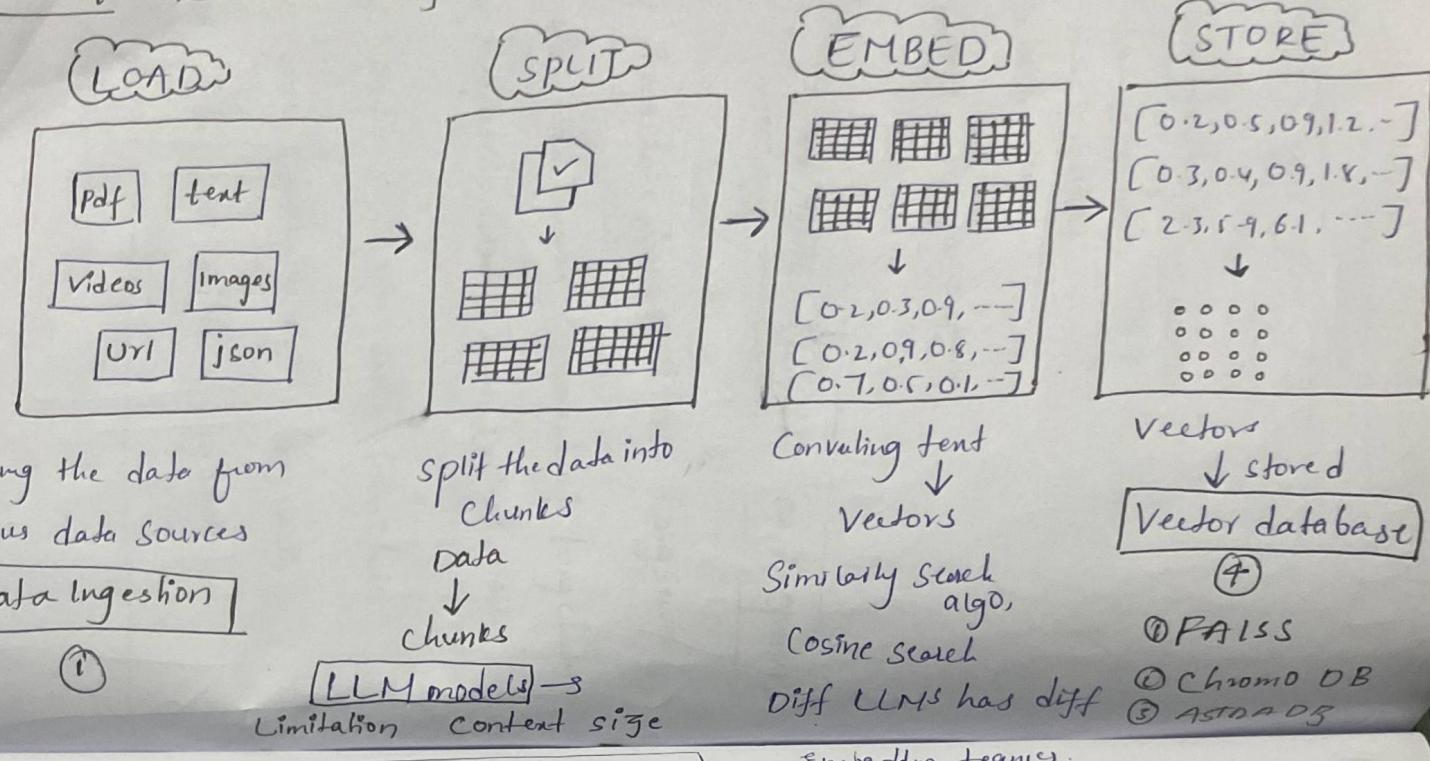
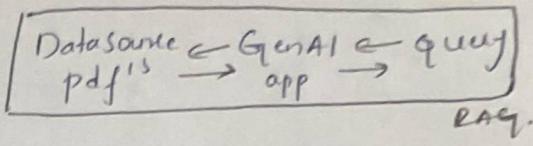
.env file

requirements.txt : python-dotenv

to import .env file
into code.

Important Components of Langchain

① RAG: Retrieval Augmented Generation.



① Document Loader in Langchain

↓
are used to import data from different sources
into a format Langchain can work with.
text, docs, websites, pdf etc.

from Langchain-community. document-loader import

TextLoader → loads "text"

Pypdf Loader → loads "pdf's"

WebBase Loader → loads "webpages"

Arxiv Loader → loads → research
papers

Wikipedia Loader → loads → wikipedia
articles

Content size (Content window):

→ refers to max no of tokens a LLM can process
at once in a single input

→ anything beyond that get truncated @)
forgotten

$$\boxed{\text{i/p prompt} + \text{response} = \text{Total tokens}}$$

↓
Content size.

② Splitting the documents

① Recursive Text Splitter:

Splits nicely first by paragraphs, then by sentences, then by words, until each piece is the right size.

② Character Text Splitter:

Splits based on fixed character count.

③ HTML Text Splitter:

Keeps to split the HTML documents into clean feed chunks while ignoring HTML tags and keeping meaningful content.

④ Recursive JSON Splitter:

Used to split large JSON documents into smaller meaningful chunks.

$$\text{Length of list} = \text{messages} + \text{tokens per message}$$

③ Langchain Embeddings

Embedding:

Converts text (chunks) to Vectors

↳ There are different LLM models which provides Embedding model like-

① OpenAI

② Ollama

③ Hugging face

④ Claude

⑤ anthropic

→ Here we are going to focus on

① OpenAI (paid)

② Ollama (open source)

③ hugging face (open source)



These three provide good accuracy which is more than sufficient.

① OpenAI Embeddings

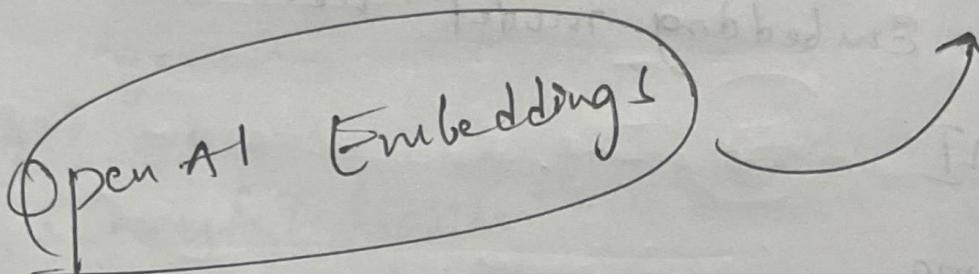
→ OpenAI API key (5 USD)

→ Billing

→ Create new Secret API key.

→ OpenAI offers three popular embedding models for converting text to vectors

- ① text-embedding-3-large most capable one
- ② text-embedding-3-small smaller embedding model
- ③ text-embedding-ada-002 the older one.



② QLLama Embeddings

QLlama is a platform where you run open-source models like Llama(Meta), Marshal (Marshal AI), Gemma (Google), Code LLaMA (Meta), Alvatelnet (Autell), and Phi (Microsoft) locally on your computer for privacy, offline access and easy AI experimentation.

QLlama embeddings

using open source models

Llama, Gemma, deepset-R,
etc.

using Embedding models

- mubai-embed-large
- nomic-embed-text
- all-nomic

⑧ Huggingface embedding

Sentence transformers (library)

from langchain-huggingface import
huggingFaceEmbeddings

① all-MiniLM-L6-V2

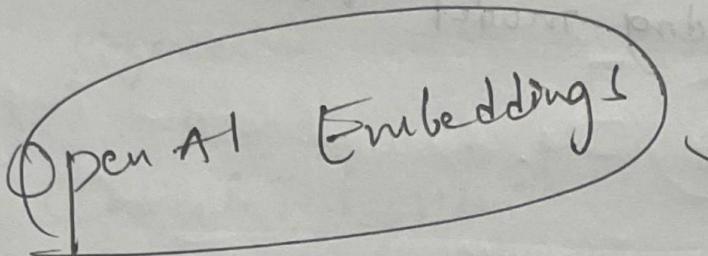
② all-mnert-base-V2

③ paraphrase-MiniLM-L6-V2

and many more - - -

→ OpenAI offers three popular embedding models for converting text to vectors

- ① text-embedding-3-large most capable one
- ② text-embedding-3-small smaller embedding model
- ③ text-embedding-ada-002 the older one.



② QLLama Embeddings

QLLama is a platform where you run open-source models like Llama (meta), Moshal (moshal AI), Gemma (Google), Code LLaMA (meta), Neuracraft (Antel), and Phi (microsoft) locally on your computer for privacy, offline access and easy AI experimentation.

QLLama embeddings

using opensource models

(Llama, Gemma, deepset-R,
etc.)

using Embedding models

→ mubai-Embed-large
→ nomic-Embed-text
→ all-nomic

⑧ Huggingface embedding

Sentence transformers (library)

from langchain-huggingface import
huggingfaceEmbeddings

① all-MiniLM-L6-V2

② all-mnpt-base-V2

③ paraphrase-MiniLM-L6-V2

and many more --

④ Vector Store DB's

in Langchain are tools used to store and search embeddings - especially useful in Retrieval-Augmented Generation (RAG) systems.

→ There are many Vector Store Db, here are few

- ① Faiss
- ② Chroma db
- ③ Ashadb
- ④ Pinecone

①

Faiss: Facebook AI Similarity Search is a

library for efficient similarity search and clustering of dense vectors, mainly used to quickly find similar items (like text or images) based on their vector embeddings library.

faiss-cpu

from Langchain-community. import

While working with LLM models, we cannot use vector store db directly instead we convert them into Retriever and give to LLM models.

$$\text{retriever} = \text{db}.as_retriever()$$

Similarly we have

↳ Similarly Search score db

↳ Similarly Search by-vector

↳ Storing vector store db in local

db.save_local('faiss-db')

↳ Loading Faiss locally

newdb = Faiss.load_local(1, 2, 3)

docs = newdb. similarity_search(query).

library

faiss-cpu

② Chromadb

Chromadb

library

(*)

→ open source vector database

langchain_chroma

chromadb = vector

met std function

db address function

vector-pd. db = push

load db from file

(db.load) (load.vae.js)

load db from pub

(load) (load.bach.js) = db with

(push) (push.vae.js) = db

(load)

Building Important components of LangChain

- ① Set up LangChain, LangSmith, LangServe
- ② Using basic common components like -
prompt templates,
models,
output parsers
- ③ Building a simple application with LangChain.
- ④ Trace your application with LangSmith.
- ⑤ Serve your application with LangServe.

Project 1

Building GenAI app using Langchain with OpenAI

↓
Querydocs AI

Answers your questions
from group of documents

model name is
"gpt-4o"

Steps involved in making project.

Step 1:

Setting up environment and API keys

↳ Load API keys from .env file.

Step 2:

Load Data from website

↳ WebBaseLoader

Step 3:

Load website data as a document

↳ loader.load()

Step 4:

split documents into smaller chunks

↳ RecursiveCharacterTextSplitter()

Step 5:

Convert text to vectors (Embeddings)

↳ OpenAIEmbeddings()

Step 6:

Store Vectors in a Vector Databases (Faiss,
Chroma -)

↳ FAISS, CHROMA, ..

Step 7:

Search relevant chunks from Vector store.

↳ SimilaritySearch()

Step 8:

Create the LLM model (GPT-4o)

↳ ChatOpenAI

Step 9:

Create the prompt format (Chat prompt template)

↳

ChatPromptTemplate.from_template ("")

Step 10:

Create the document chain

↳ create stuff-documents chain (llm, prompt)

→ Combines model & prompt into a chain

Step 11: → takes multiple docs + prompt → model → answer

Try chain with manual content

↳ manually giving the document and

Step 12:

Checking

Turn Vector Store into Retriever

→ vdb-as-retriever()

→ can auto pick the best matching documents

Step 13:

Create a Full Retrieval Chain.

Retriever → to find relevant docs

document-chain → takes all docs adds prompt and
Ask a Real Question

↳ check for new question

Step 14:

View the Retrieved Documents.

resp = retrieval_chain.invoke({ "input": ... })

resp['answer']

Project - 2

Building Gen AI App using Langchain with Ollama

↓
Ollama Q&A

model is

"deepseekr1"

Step 1: Imports - Setting Up the Environment.

↳ libraries, .env file

Step 2: Langchain and Ollama imports

→ Ollama, ChatPromptTemplate, StreamOutputParser

Step 3: prompt template creation

→ helps in controlling how AI should respond.

Step 4: Building the Streamlit UI

→ for building UI

Step 5: LLM setup and chain creation

→ prompt / LLM / output - parser

Step 6: Handling User Input and Showing O/P.

↳ asking Ques from user and answering back.

This run locally using deepseekr1 via Ollama.

Building Basic LLM application using LCEL

Langchain Expression Language

Getting started with Open Source Models Using
Groq API

① open AI API key → paid

② Open source models

① Llamas } through platform

② Gemma2 }

③ mistral }

Groq

Similar to

QLlama

Create Groq API
key from Groq
Cloud

→ low latency

→ Cloud host

→ CPU

language proxy

unit.

helps in calling
Open source models.

Why Groq

AI inference

process of using a learned AI model to make predictions (①) generate O/p's from new input data.

Training → model learns from data

Inference

When the model uses, what it learned ~~to calculate~~ to give answers

Ques → ChatGPT reply → Inference

Google translate converts sentence → Inference

AI models detect faces in photos → Inference

why fast inference matter.

① Real time Experience

slow AI → Frustrating

fast AI → Seamless Experience

② Better User engagement

③ Useful in Critical applications

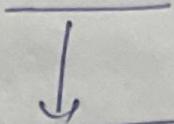
Healthcare, finance, security ---

④ Cost efficiency at scale

less compute time / require

project - 3

Building Basic LLM application. using LCEL



Using Grog
Gemmar-9b-It

Langchain Express
language

Steps involved in project

Step 1: Importing required libraries

fastapis , shoutputparser , OS , add_routes
chatpmptemplate chatGrog , dotenv

Step 2: load api key from .env file

→ load_dotenv()
→ Grog-api-key

Step 3: Initialize Grog LLM model

ChatGrog(Gemmar -9b-It)

Step 4: Define prompt Template.

↳ LLM understand its role & task clearly

System prompt , user prompt

Step 5: Define Output Parser

↳ converts LLM OLP into plain string format

Step 6: Build the Langchain Chain

Building chain → Pipeline

prompt template / model / parser.

Step 7: Create Fast API APP.

→ creates app

→ allows people to send (or) get
request response.

Step 8: Expose Langchain Chain as API Route

→ connects langchain model to API

Step 9: Run the APP with Uvicorn

→ Starts the web server that will run your
Fast api app on your computer (localhost)

Simple translation API that runs on your computer
using FastAPI. It sends your input text to GPT
model online and returns the translated sentence

Langchain ↗ helps to connect requests to AI smoothly.
Lang Serve

project - 4

Building Chatbots with Conversation history using Langchain.

Load Credentials from .env

os.getenv(''), apikeys)

Run inference using Groq UCM

model, ChatGroq

chat without Memory

stateless.chat, model.invoke()

keep history per session

Memory store, ChatMessageHistory

Add memory to model

Memory wrapper, RunnableWithMessageHistory

I identify each user's history

Session ID's, session-id = "chat1"

Define Structured Inputs.

prompt templates, ChatPromptTemplate

trim messages

Token Trimming, avoid exceeding token limits

Combine trimming, templating, inference

full chain, RunnablePassthrough → prompt → model.

Important functions in Langchain, that everyone should know before Creating Chatbots with Conversation history.

① ChatMessageHistory()

Stores all messages in a single chat

② get_session_history(session_id)

gets message history for session (or)

Creates new one if it doesn't exist

③ RunnableWithMessageHistory()

adds memory support to your chatbot

④ ChatPromptTemplate.from_messages([...])

Builds a prompt using system and user messages

⑤ MessagePlaceholder(varname="messages")

Reserves space in the prompt to insert past messages

⑥ trim_messages(...)

removes old messages (or) unimportant ones

⑦ itemgetter("messages")

Pulls out the messages key from dictionary

⑧ RunnablePassthrough(...)

prepares data (like trimming msgs) before sending to the model

Project - 5

Conversational Q&A chatbot with Message

History [Full steps]

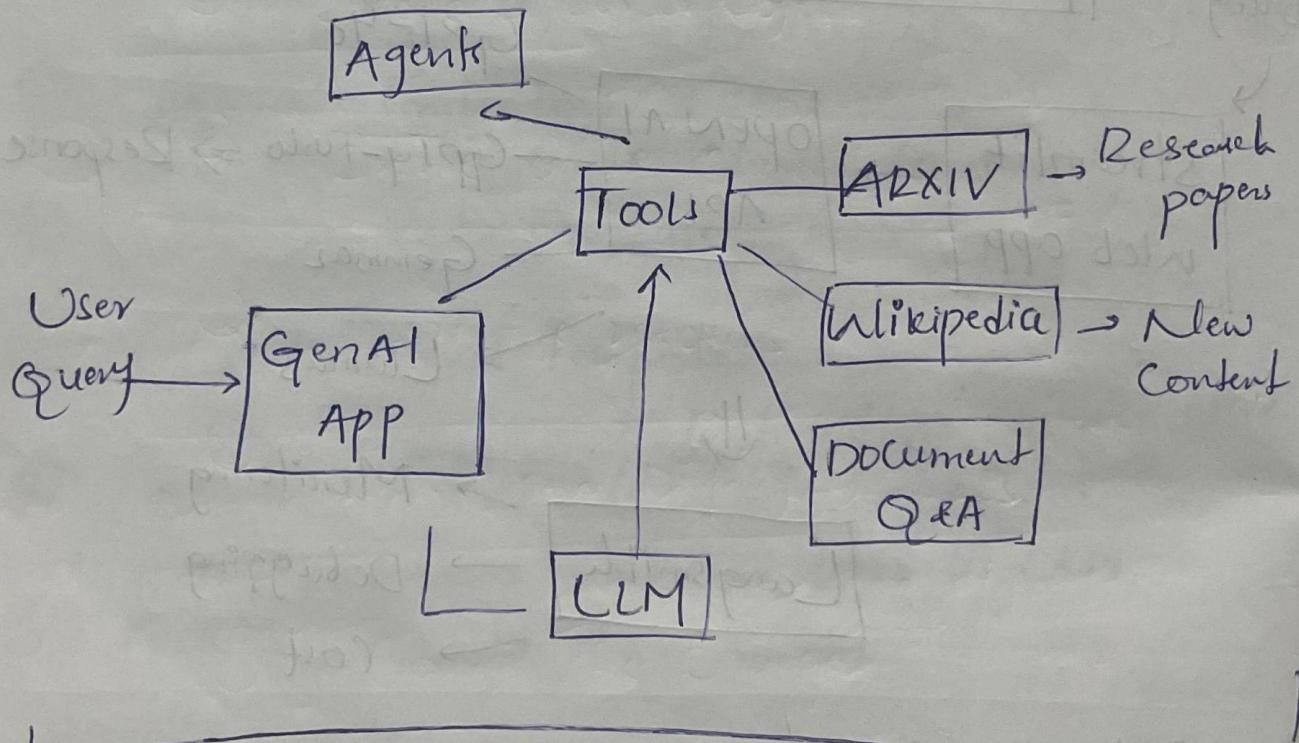
- ① Setup and Load environment
→ Loading API keys
- ② Initialize LLM
→ Loading Model
- ③ Load Data from Website
→ Scraping the data
- ④ Split Text into Chunks
→ Breaking long text → chunks
- ⑤ Generate Embeddings
→ Chunks to dense numeric vectors
- ⑥ Store in Vector Database
→ Store embeddings
- ⑦ Create prompt templates
→ Define how System & User interact.
- ⑧ Create Retrieval Chain
→ Retrieve relevant chunks and pass to LLM
- ⑨ Ask Questions Using the Chain
→ ask ques, get content based answers
- ⑩ Enable Chat History Awareness
Add memory to prompts, make retrieval aware of flow
- ⑪ Manage Chat History
Keep track of human and AI messages
- ⑫ Ask follow-up Questions.
→ play with chatbot.

Search Engine with Tools And Agents

Tools: tool is a function that does a specific task, like doing math, searching the web or running code. It helps the chatbot get things done.

Agents: is a smart controller that decides which tool to use when, in order to answer the user question.

OpenAI LLM model → Gpt-40 ⇒ December 2023 → Data.



End to End Search Engine

Wrappes: Gets data from website.

QueryRun: Uses the Wrappes to answer your question.

Steps to create an Agent

① Creating Wrappers

→ Wikipedia API Wrappes, ArxivAPI Wrappes etc...

② Creating tools using those Wrappers

→ ArxivQueryRun, WikipediaQueryRun.

③ Combining all the tools

→ add all tools

④ Creating LLM

→ Ollama, openai, Gptq.

⑤ Create prompt Template.

→ langchain.hub.pull(* *)

⑥ Create Agent

→ langchain.agents.CreateOpenaiToolsAgent.

⑦ Create Agent Executor

→ Langchain.agents.AgentExecutor

⑧ Invoke Agent Executor with different questions

→ agent.executor.invoke("input": "What is AI")

project 7: Chat with SQL DB with langchain

SQL Toolkit and Agent type

① preparing the data for SQLite3 database

① Import the SQLite Library

→ import sqlite3 (python library)

② Connect to the Database

→ Creates a new database ←

③ Create a cursor object

→ used to execute SQL commands

④ Create a Table

→ Define table

⑤ Insert the records

→ add records to table

⑥ Fetch and Display Records

→ Retrieves and displays all the

⑦ Commit the changes records

→ Save all changes made to db

⑧ Close the Connection.

→ Close to free resources.

② Steps involved in project

- ① Import libraries
Streamlit, pathlib, langchain, SQLAlchemy, Sqlite3, ChatGPTQ
- ② Page setup
Filters, Readings, styling, sidebar, footer.
- ③ Database Selection
Sqlite3 local, MySQL external.
- ④ Collect MySQL details if selected
Username, password, DB name.
- ⑤ Enter GPTQ API key.
GPTQ api key
- ⑥ Set up the LLM
Initialize LLM
- ⑦ Configure the database.
Sqlite @) MySQL (create_engine)
- ⑧ Connect to Selected DB
User choice.
- ⑨ Prepare SQL toolkit
Creating tool
- ⑩ Create LangChain SQL agent.
Creating agent
- ⑪ Session history initialization
→ Starts new chat history
- ⑫ Show Chat history
→ Display messages
- ⑬ User sends a New Question
→ User input.
- ⑭ Run Query and Show Response.

project 8: Text Summarization with Langchain

① When you want to summarize small number of tokens: Fits in LLM context window

→ Manual prompting

→ prompt template with LLM chain.

① Manual prompting

↳ using AIMessage
SystemMessage } = msgs
HumanMessage }

→ llm(msg) simple ✓

② Prompt template with LLM chain

modules → LLMChain, PromptTemplate

Prompt template

→ Input variables = []

→ template =
generic
prompt template

LLMChain

→ llm = llm

→ prompt = prompt

② Steps involved in project

- ① Import libraries
Streamlit, pathlib, langchain, SQLAlchemy, Sqlite3, ChatGPTQ
- ② Page setup
HTML, filters, headings, styling, sidebar, footer.
- ③ Database Selection
Sqlite3 local, MySQL external.
- ④ Collect MySQL details if selected
Username, password, DB name.
- ⑤ Enter GPT API key.
GPT api key
- ⑥ Set up the LLM
Initialize LLM
- ⑦ Configure the database.
Sqlite @) MySQL (Create_engine)
- ⑧ Connect to Selected DB
User choice.
- ⑨ Prepare SQL toolkit
Creating tool
- ⑩ Create LangChain SQL agent.
Creating agent
- ⑪ Session history initialization
→ Starts new chat history
- ⑫ Show Chat history
→ Display messages
- ⑬ User sends a New Question
→ User input.
- ⑭ Run Query and Show Response.

Project 8: Tent Summarization with Langchain

①

When you want to summarize small number
of tokens: Fits in LLM context window

→ Manual prompting

→ prompt template with LLM chain.

②

Manual prompting

↳ using AIMessage

SystemMessage } = msgs

HumanMessage }

→ llm(msg) simple ✓

②

Prompt template with LLM chain

modules ↳ LLMChain, PromptTemplate

Prompt template

→ Input variables = []

→ template =

generic
prompt template

LLMChain

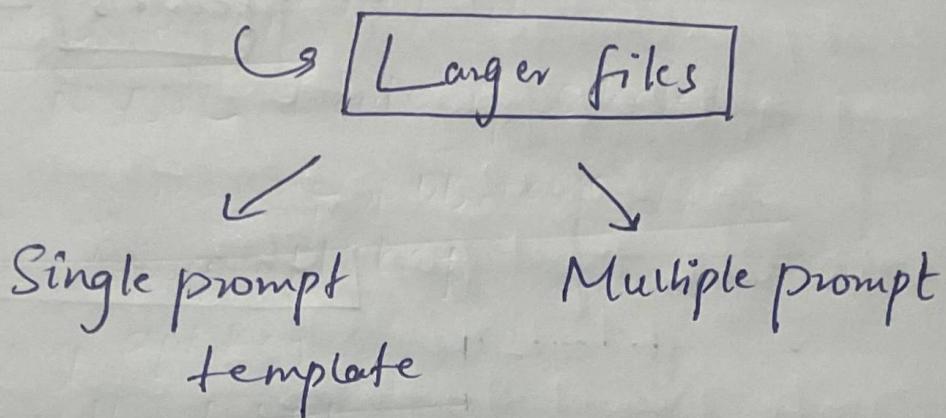
→ llm = llm

→ prompt = prompt

② When you want to summarize large number of tokens : does not fit in LLM Content window

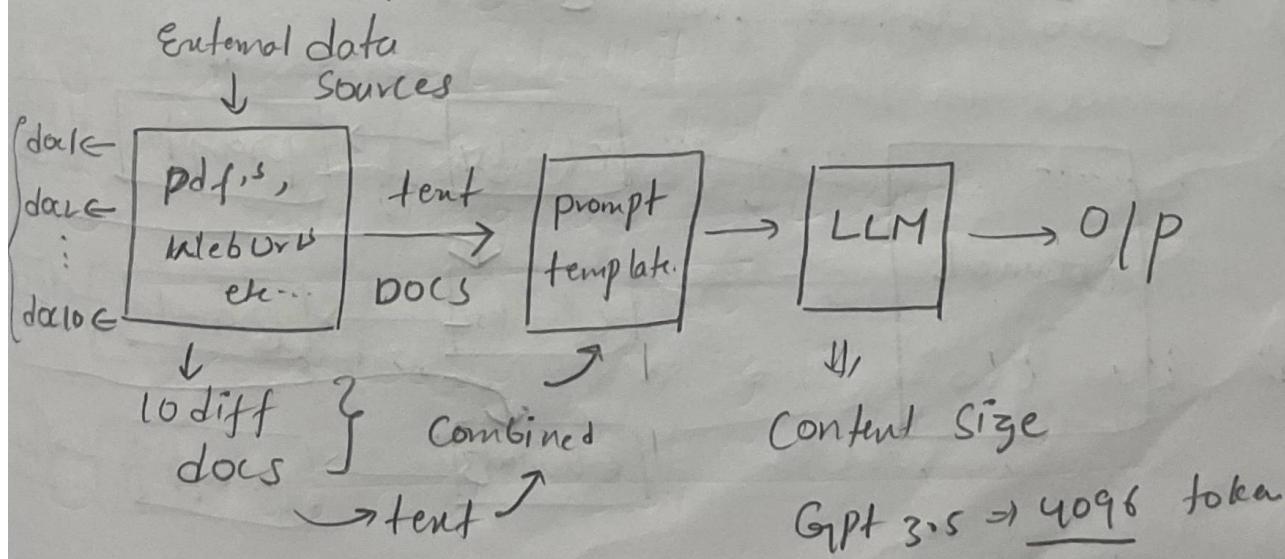
① Stuff Document Chain Summarization.

② Map Reduce Summarization Technique



③ Refine Chain Summarization.

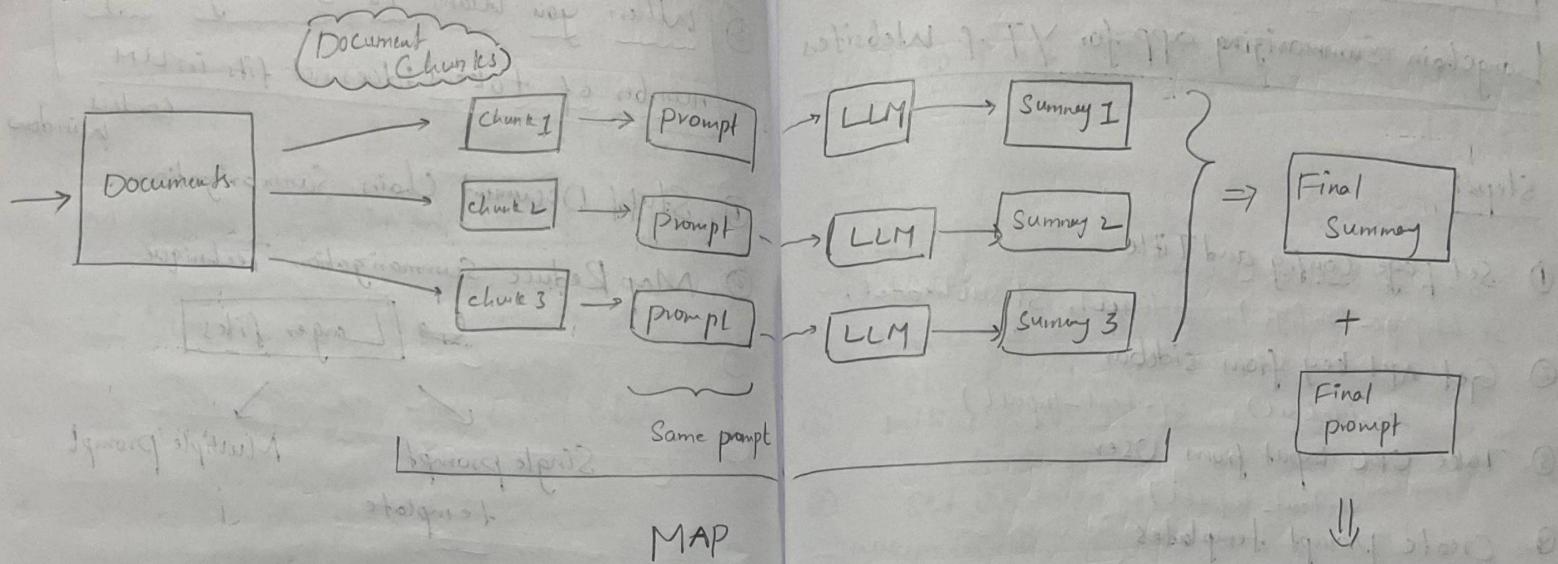
① Stuff Document Chain Summarization.



challenges :

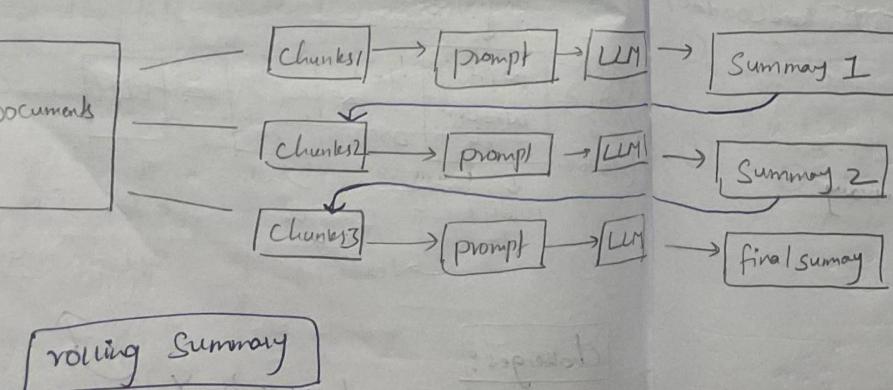
Larger documents X

② Map Reduce Summarization Technique



③ Refine Chain Summarization

→ Update a rolling summary by iterating over the documents in a sequence.



Project 10:

Langchain Summarizing app for YT & Website

Steps:

- ① Set page Config and Title
st.set_page_config(), st.title(), st.subheader()
- ② Get API key from sidebar
st.sidebar(), st.text_input()
- ③ Take URL input from User
st.text_input()
- ④ Create prompt templates
PromptTemplate()
- ⑤ Validate inputs
validators.uri(), st.error()
- ⑥ Initialize LLM
chatGPT()
- ⑦ Load Content from URL
YoutubeLoader.from_youtube_url(), unstructuredLoader
- ⑧ Create Summarization chain
load_summarize_chain()
- ⑨ Run chain and display summary
chain.run(), st.success()
- ⑩ Handle errors.
try-except(), st.error()

Project 11: Text to Math problem Solver Using Langchain

Text to Math problem Solver Using Langchain

Steps:

- ① Set up page config and title
st.set-page-config(), st.title()
- ② Enter Groq API key
st.sidebar.text_input()
- ③ Initialize LLM
chatGroql
- ④ Set up tools
wikipedia tool, calculator tool, Reasoning tool
- ⑤ Initialize agent with tools
initialize_agent(tools=[])
- ⑥ Set up chat history
st.session_state(), st.chat_message()
- ⑦ User input section
st.text-area()
- ⑧ Solve Button logic
Input check, append user messages,
agent response,
- ⑨ Display final answer
st.chat_message("assistant").write(response)
st.success()

85% of Course Covers:

- General and LCM models
- LangChain for GenAI
- LangChain with OpenAI
- Components and Modules in LangChain
- OpenAI + Ollama
- LCBL
- more than 12+ projects covered

↓
Database, how to deploy
Astra dB genAI app w/m

Shuttle

Project 12:

PDF Query RAG with Langchain & AstraDB

AshaDB is a ready to use version of Cassandra
made by DataStax, that runs in the cloud

Steps:

1. Install and Import required Libraries
pypdf2, Cassandra, VectorStoreWrapper, OpenAI, Cassio, openAIEmbedding
etc...
2. Load api keys and connect to AshaDB
openai.apikey,
ASTRA-DB-APPLICATION-TOKEN, ASTRA-DB-ID, cassio.init()
3. Read text from the pdf file
Read the data
4. Split the text into chunks
CharacterTextSplitter()
5. Set up LLM and embeddings
OpenAI(), openAIEmbedding()
6. Create vector store backed by Ashadb and insert
chunks
Cassandra(), vectorstore.add_texts()
7. Wrap the vector store for easy querying
VectorstoreIndexWrapper()
8. Ask queries
asha_vector_index.query(question="")

How to deploy GenAI app on Streamlit Cloud

Prerequisites:

- ① Github account
- ② Streamlit Cloud account (linked to Github)
- ③ A working Streamlit python app.

Steps:

- ① Create your Streamlit app.
- ② Create requirements.txt
- ③ push your project to Github
- ④ Deploy app on Streamlit Cloud.
 - Visit Streamlit Cloud
 - Sign in with Github
 - Click + New app
 - branch → main, filepath → app.py
 - Click deploy
- ⑤ Optional
 - add configurations, add secret for api keys
- ⑥ Share the app.

Deployment of Gen AI app in Hugging face spaces :

Prerequisite:

- ① Github account
- ② huggingface.co account.

Steps:

- ① Create new file in the git hub repository
.github/workflows/main.yaml
- ② Copy the code from "huggingface spaces"
"github actions" and past in "main.yaml"
- ③ Now go back to hugging face spaces,

Create new space :

- spacename ↗
- license.
- Streamlit
- cpuBasic 2vCPU free
- public
- Create space.

④ Create access token

- Go to huggingface.co
- access token
- new token
- Generate token (Write mode)

⑤ replace your Username in main.yaml to

your hugging face username.

- Similarly replace spacename with the spaces name you gave before

⑥ Now go to github settings

- Secrets and Variables
- Actions
- New repository Secret
 - ↳ Name HF-Token
 - ↳ Secret
 - ↑
paste the hugging face access token here.
- add secret

NVIDIA NIM

NVIDIA Inference microservice is a prebuilt, optimized container that delivers fast and scalable AI models inference with minimal setup.

Project 13:

Building RAG Document Q&A with Nvidia NIM And Langchain.

Steps:

① Import Required libraries:

st, os, dotenv, Nlver embeddings, ChatNvidia ,
PyPDFDirectoryLoader , RecursiveChatHistorySplitter , False ,
ChatPromptTemplate , CreateRetrievalChain , stuff doc
chain .

② Load API key:

NVIDIA-API-KEY = os.getenv('')

3. Initialize Nvidia LLM

chatNvidia() \Rightarrow Hafas - 706 - Instruct

4. Create Vector Embeddings from pdfs

NvidiaEmbeddings(), pypdf Document Reader(),
RecursiveChatTextSplitter().

5. Streamlit User Interface

title(), write(), buttonClicked() etc.

6. Prompt template:

Chat prompt template - from template(q1 , q2) equal

7. Take User question and Get answer.

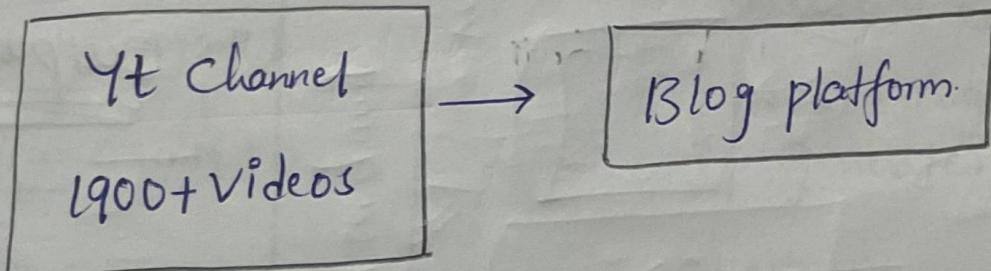
User input \rightarrow retrieve relevant documents \rightarrow

Combine document with prompt \rightarrow Send to Nvidia LLM
 \rightarrow generate the answer.

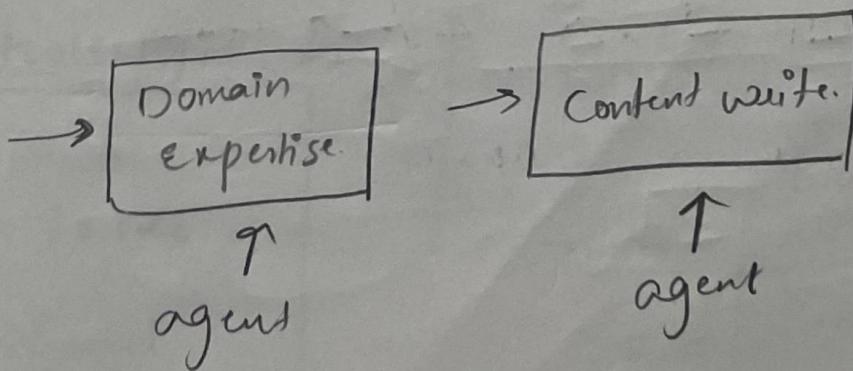
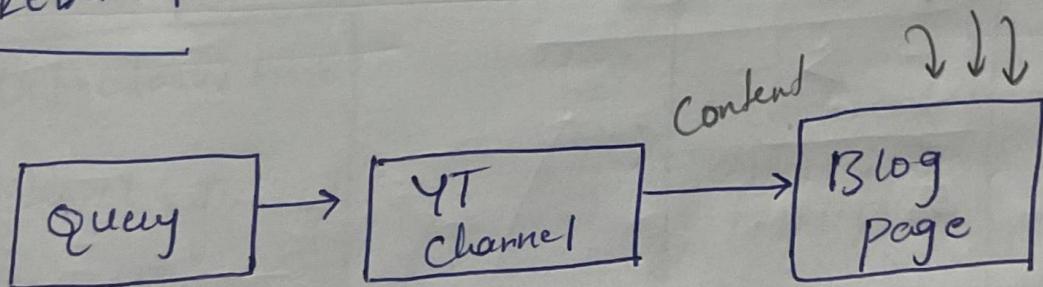
CREW AI

CREW AI is a framework that allows multiple AI agents to work together as a team to complete complex tasks efficiently through collaboration and coordination.

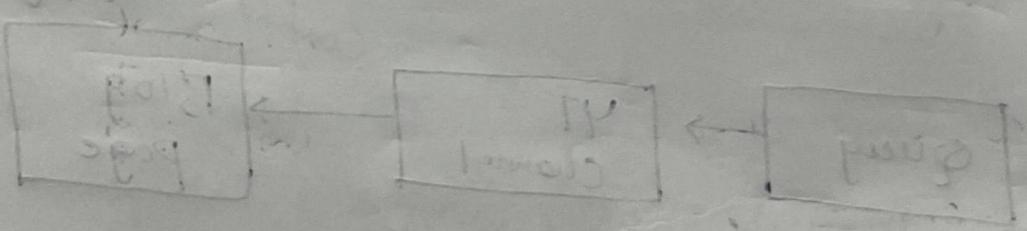
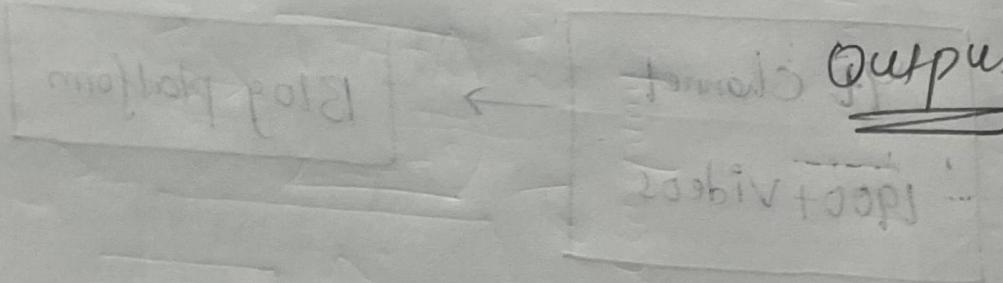
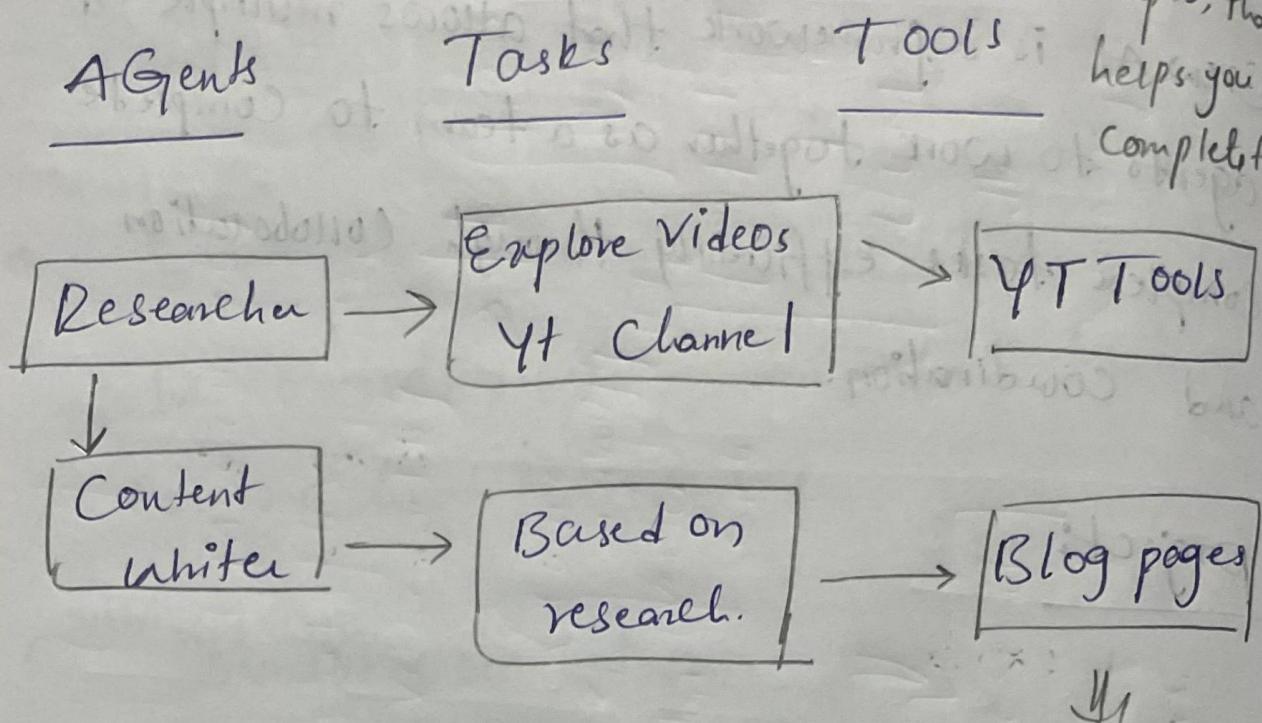
project:



CREW AI



CREW AI



Project 14:

YouTube videos to Blog page Using CREW

AI Agents

Project flowchart:

agents.py

```
Define agents here:  
researcher=Agent()  
writer=Agent()
```

tasks.py

```
Define tasks here  
research=Task()  
writer=Task()
```

tools.py

Define tools to run tasks

YT channel
ScrapeTool()

-pyCache folder
db folder

While crew.py is running it generates

crew.py

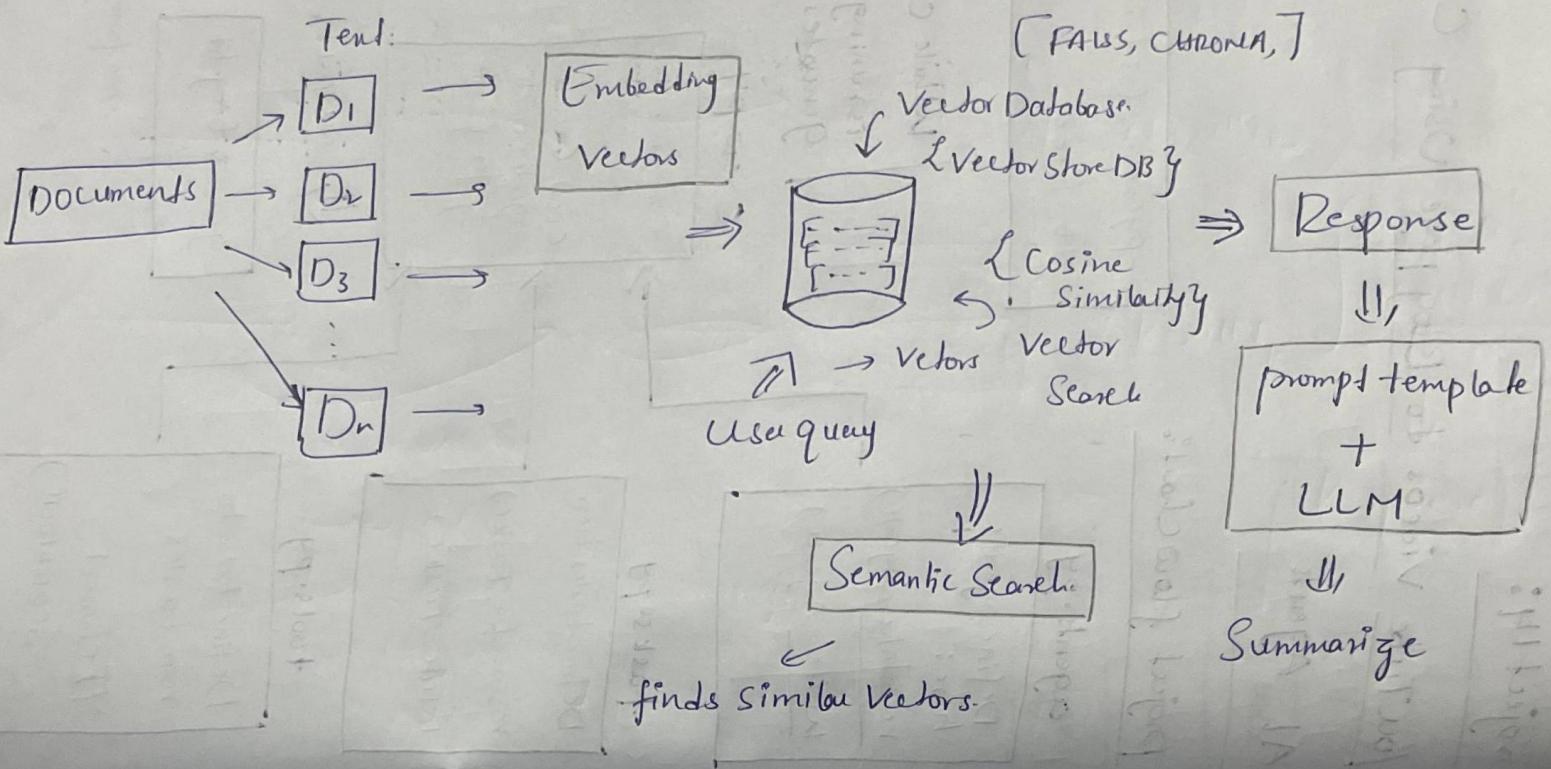
```
Combine all  
crew=Crew()  
agents=[...]  
tasks=[...]
```

||, output

generate blog file

.md format

Hybrid Search RAG



Hybrid Search Mechanism: here we are not just focusing on Vector dense Search.

Combines Multiple Search techniques.

- ① Semantic Search → Dense vector Search → Similar.
- ② Syntactic Search → Exact Search → keyword Search

Tent → Sparse matrix

By using OHE, BOW, TF-IDF

↳ Sparse Vector Search.

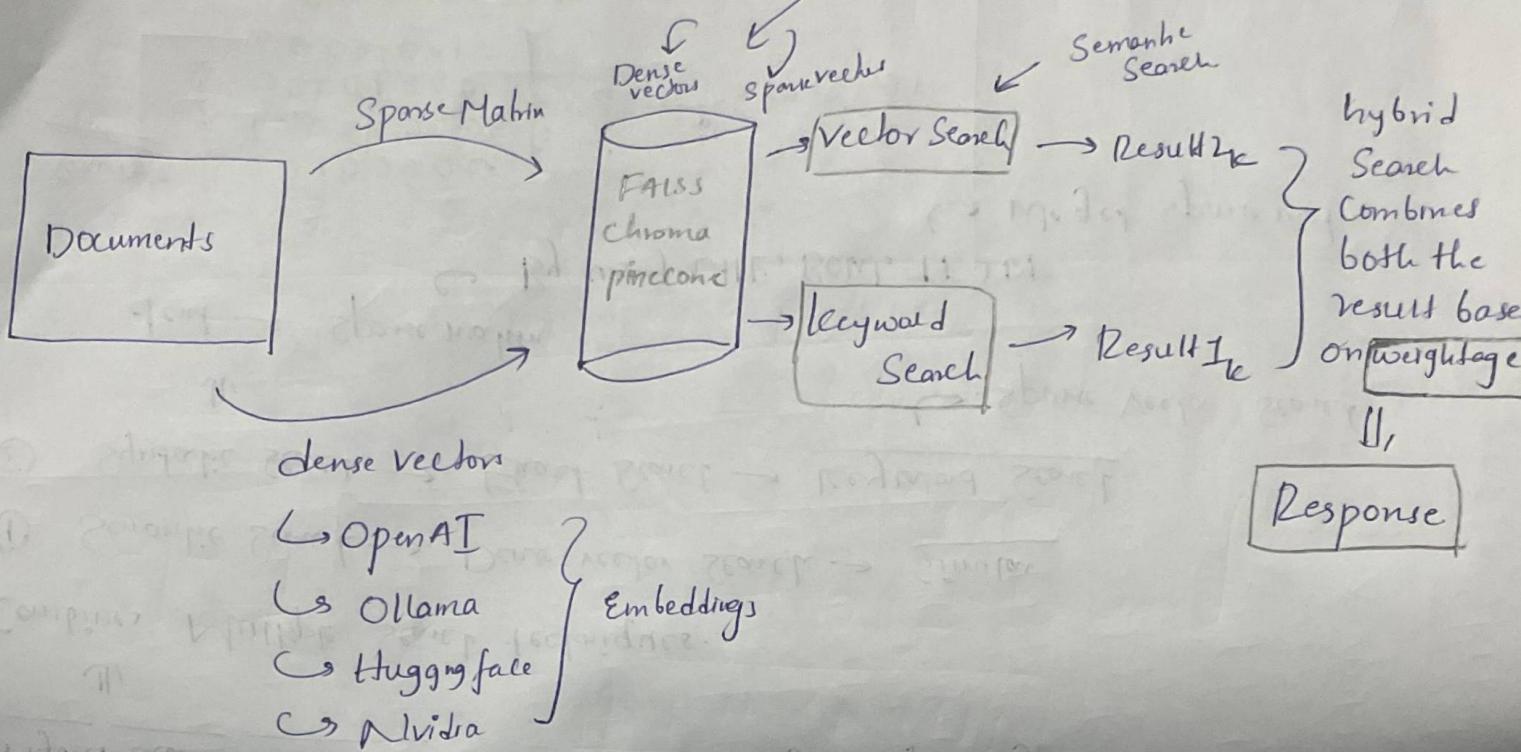
↳ We get Sparse matrix.

Example

$$\begin{bmatrix} 0 & 0 & 1 & 0 & 0 & 0 & 0 \\ 1 & 0 & 0 & 0 & 0 & 0 & 0 \\ 0 & 0 & 0 & 1 & 0 & 0 & 0 \\ 0 & 0 & 0 & 0 & 0 & 0 & 1 \end{bmatrix}$$

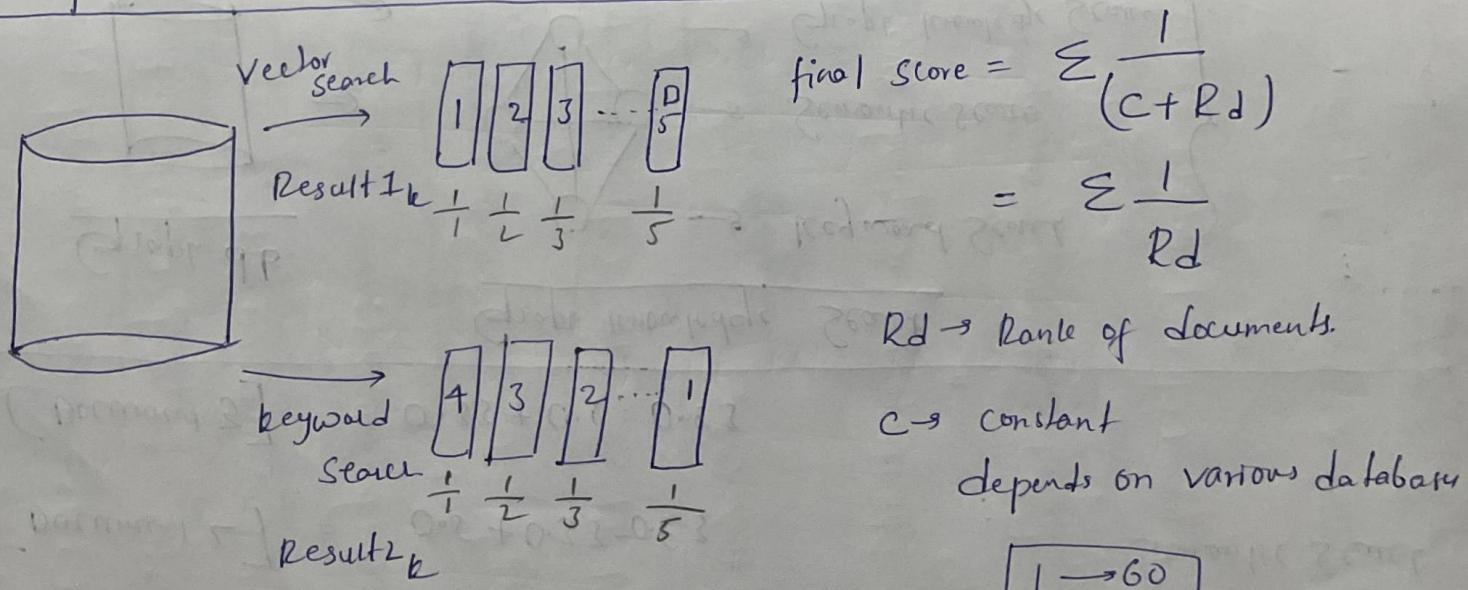
Sparse matrix

How Does Hybrid Search Work?



How weightage is given : Based on the technique

Reciprocal Rank Fusion in Hybrid Search



$$[Document 1] = 1 + 0.2 = 1.2$$

50%

\Rightarrow keyword search

$$[Document 2] = 0.5 + 0.33 = 0.83$$

Semantic Search

$$[Document 3] = 0.33 + 0.5 = 0.83$$

Graph db

Graph knowledge Search

\rightarrow Keyword Search

\rightarrow Semantic Search

\rightarrow Graph knowledge Search

use
query

Graph db [Neo4j]

Graph DB with Langchain

- ① Knowledge Graphs
- ② Neo4j Graph Databases
 - ↳ Cypher Queries
- ③ RAG application with Graph Databases

LangChain

↓ amazing module

Lang graph

↓
Where we can Create

Multi AI agents

↓ particularly uses

knowledge Graphs

① knowledge Graph:

also known as Semantic Network



Network of real world Entities

Eg: Events, situation, Concepts



illustrate the relationship b/w them

NLP Use Case:

Rohit Sharma is the Captain of India"

entity

entity

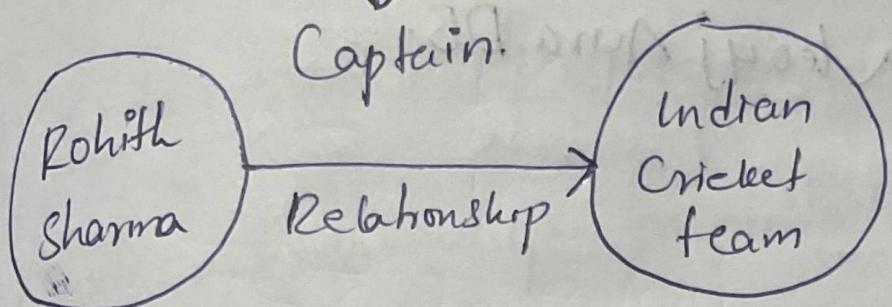
knowledge Graphs:

Three main Components

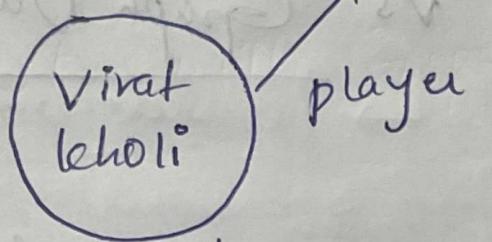
① nodes → any object, place, person ② Edges

③ Labels

Label → specifies relationship



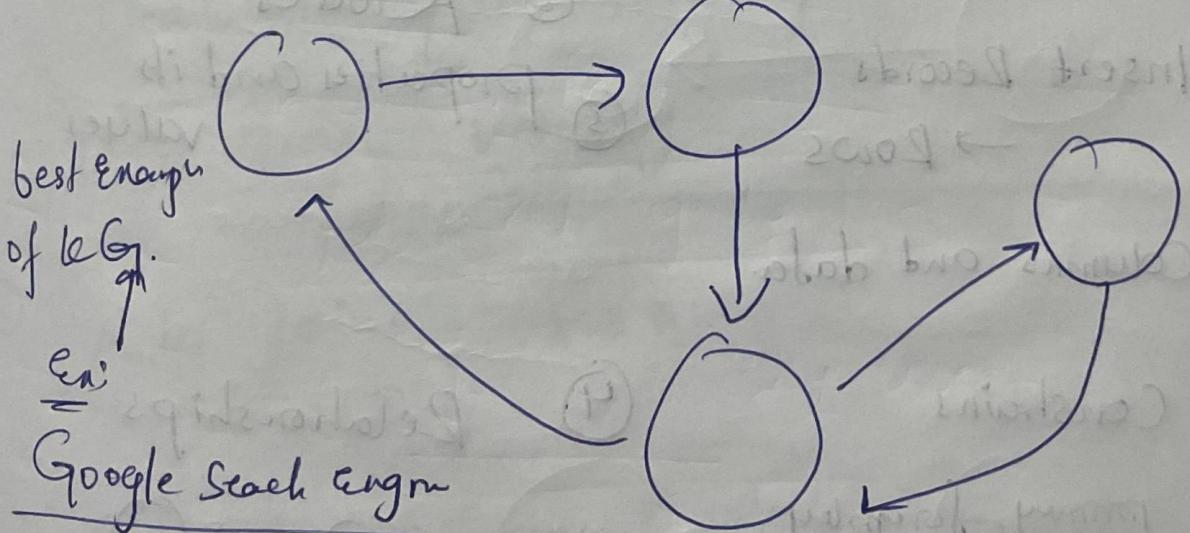
node



node

"Virat Kohli is the player of Indian Cricket team"

knowledge Graph



Search: virat kohli

↓
Suggestion → Other players

→ Create Neo4j database from
Neo4j Aura DB.

RDBMS vs Graph Databases



Relational DBMS → MySQL, PostgreSQL, Sq / Server

RDBMS

- ① Tables
- ② Insert Records
→ Rows
- ③ Columns and data.
- ④ Constraints
primary, foreignkey
etc--
- ⑤ Joins Queries

Graph Database

- ① Graphs
- ② Nodes
- ③ properties and its values
- ④ Relationships
- ⑤ Traversal.

Advantages of Neo4j

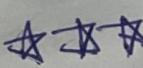
① Data model : Responsible for creating
Nodes, Relationship, properties and values

② Real time insights:

③ Easy Retrieval : Cypher query like
SQL query

④ Cypher query language :

Declarative query language to represent
the graph visually.

⑤ No joins 

⑥ Supports ACID properties:

atomicity, consistency, isolation, durability

⑦ No fixed Schema, instead supports
flexible Schema.

Neo4j property Graph Data model

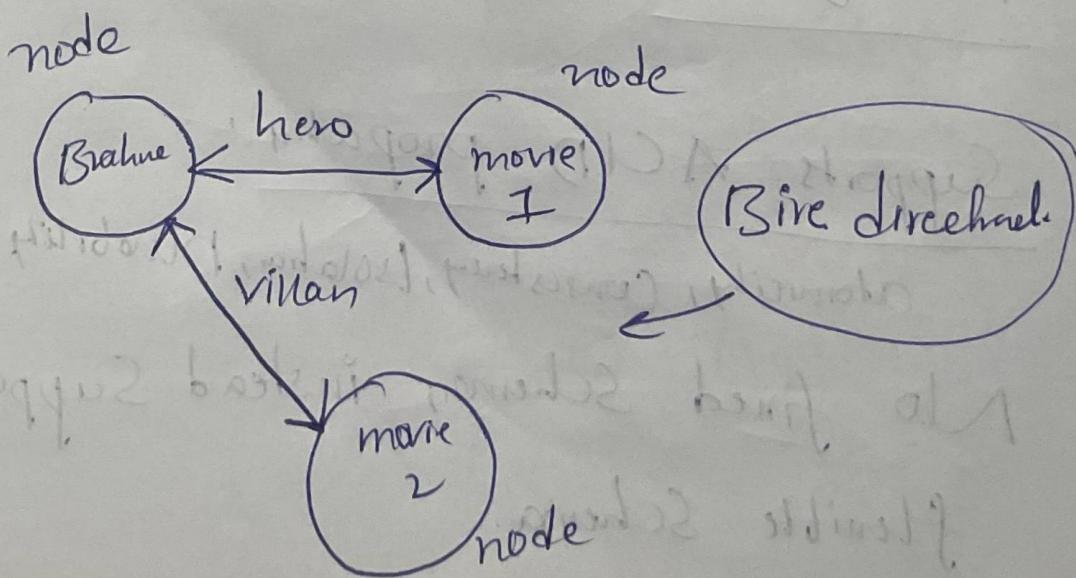
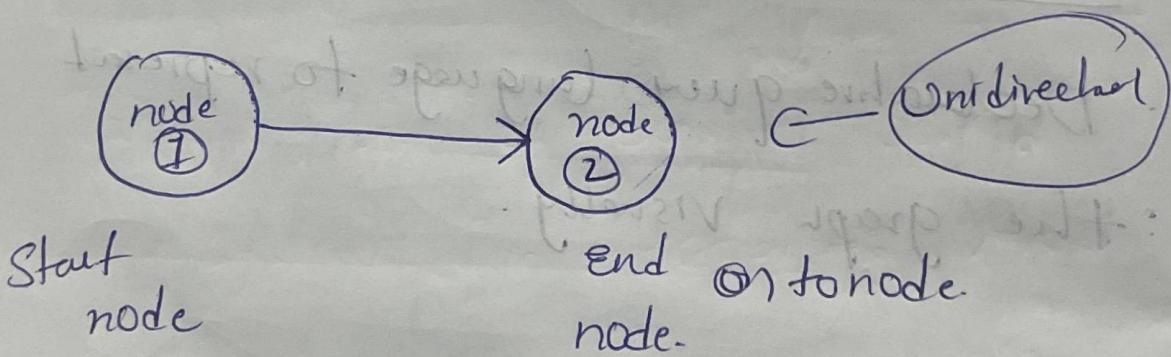
Neo4j Graph Database

↓
has property Graph Model

↓ ↳ to store and manage data.

Data = { Nodes, Relationships, properties }

↓ ↓
Unidirectional, key, value
Bidirectional pairs



Getting Started with Cypher Query language

Create(Nodel)

① Create(Elon:CEO)

② ↓ ↓
node value

③ Create(Elon:CEO:Employee)

④ Create(Elon : CEO name : "Elon musk",

YOB : 1979, POB : "SA" })

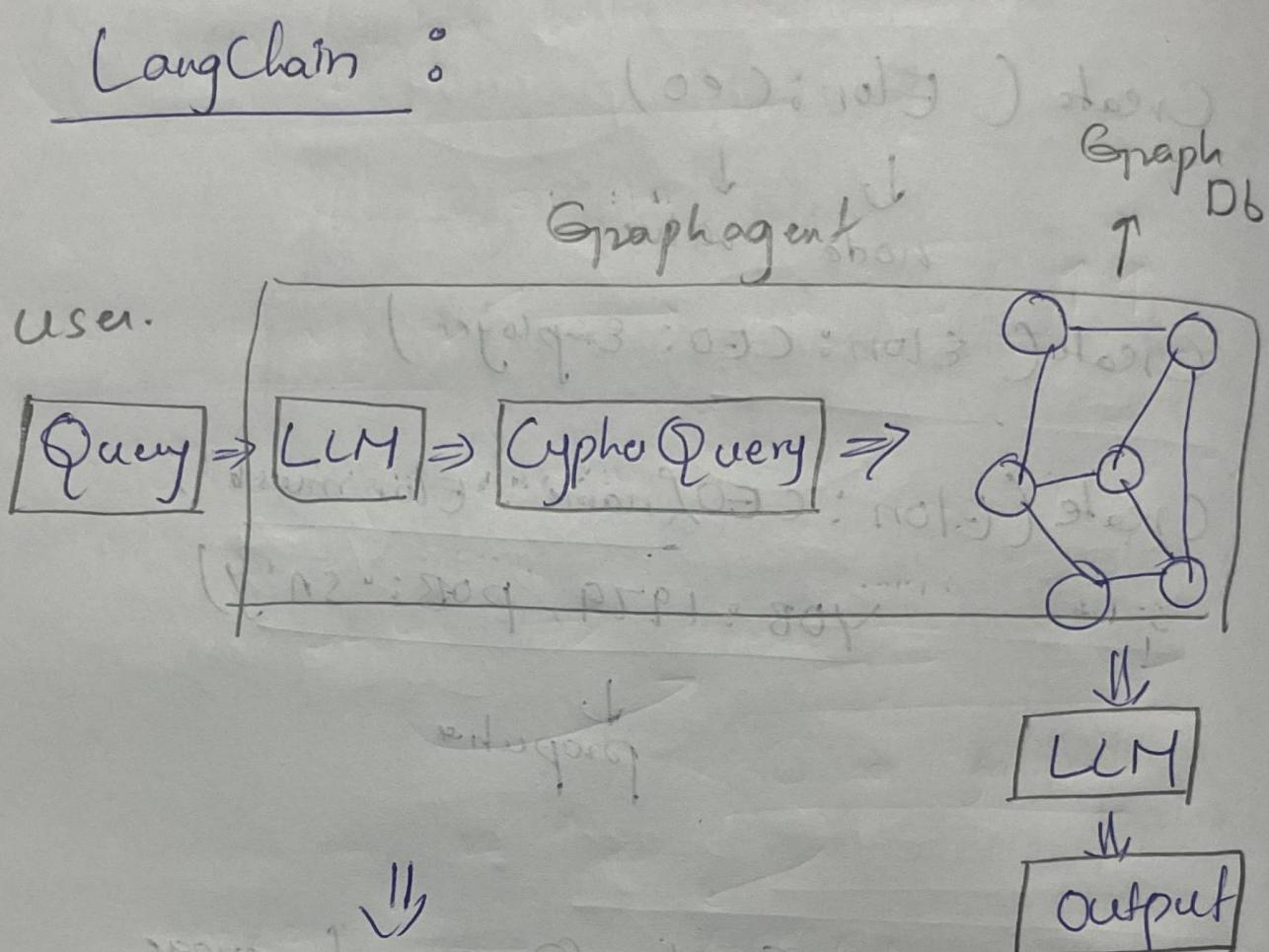
↓
property

Learn about Cypher Query language

that really helps you to create your
graph data bases.

Project 15^o

Practical Implementation of Graph DB with LangChain



Tried my best, but does not work in my system, due to changed Version and Compatibility issues.

What is Quantization In depth intuition

Quantization :

→ Data → weight & parameters.

① Full precision / half precision

② Calibration → problems

③ Modes of Quantization

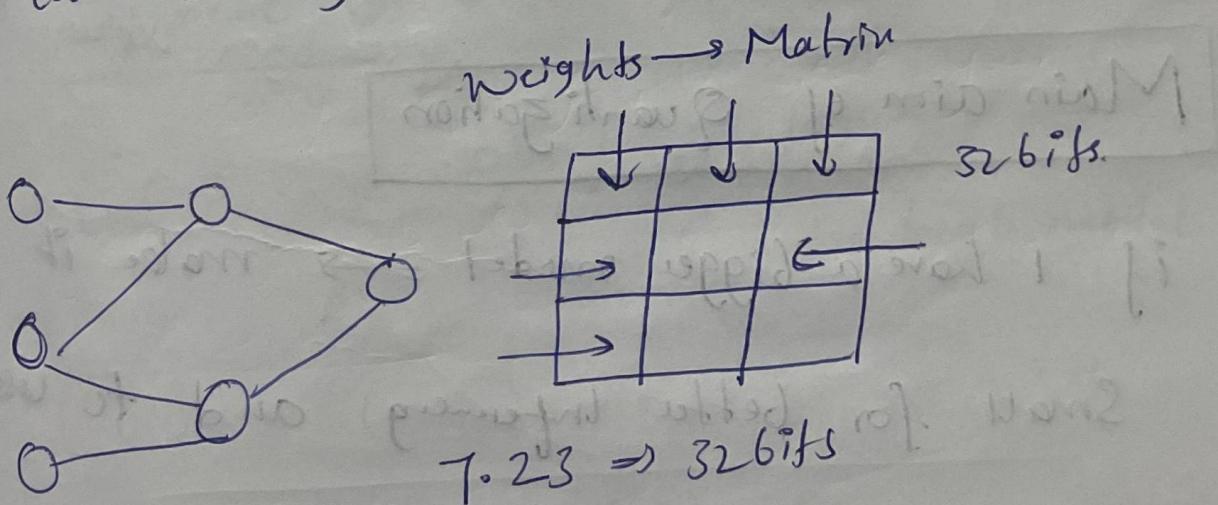
post learning
quantization

Quantization aware

Training.

Quantization

Conversion from higher memory format to a
lower memory format.



FP → Full precision / single precision

floating point 32 bit

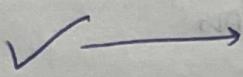
if you want to load

[Clama 2 706]

it is not supported by
your System with 32gb Ram.

↓ to do that

[Ans]



but [Cost] ↑↑

So that we convert

32 bit ⇒ int 8



⇒ [Quantization]

Very important for Inference

[Main aim of Quantization.]

If I have a bigger model → make it

Small for better inference and to use

in mobile phones, edge devices.

Smart watches etc..

We can also fine tune the model after quantization but there is loss of accuracy
↓
due to loss of memory

there are techniques to overcome it

How to perform Quantization.

- ① Symmetric Quantization \Rightarrow Batch normalization.
- ② Assymmetric Quantization.

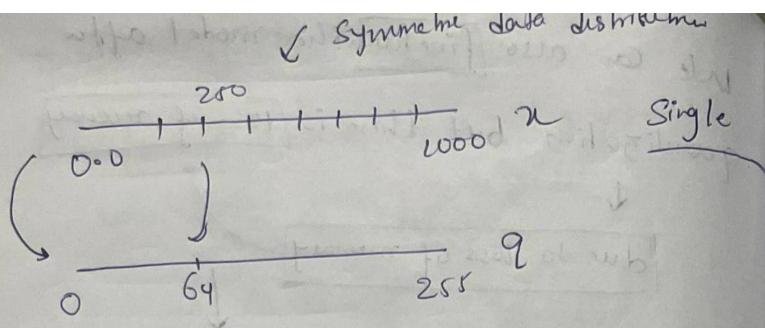
Symmetric Unsigned Quantization

$[0.0, \dots, 1000] \rightarrow$ Integers \times
for larger model.

$[0.0 \dots 1000] \rightarrow$ number \rightarrow 32 bits

Quantize
uint8 \Rightarrow 8 bit $\Rightarrow 2^8$

$0-255$



Mm mow Slabu

$$Scale = \frac{u_{max} - u_{min}}{q_{max} - q_{min}}$$

$$= \frac{1000 - 0}{255 - 0} = \boxed{3.92}$$

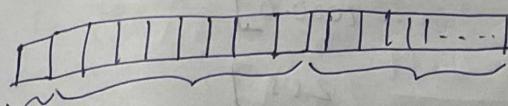
$$\text{round } \left(\frac{250}{3.92} \right) = 64 \leftarrow \begin{array}{l} \text{Scale factor} \\ \downarrow \\ \text{quantized value} \end{array}$$

Symmetric quantization

$\rightarrow 0$

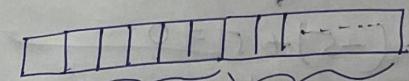
Precision Floatg point 32

how number is stored



↓ exponent Mantissa. + 1.5
 sign
 7th bit next 8 bits remain 6 bits 23
 ↓ ↓
 ↓ 1 0001 32 00101

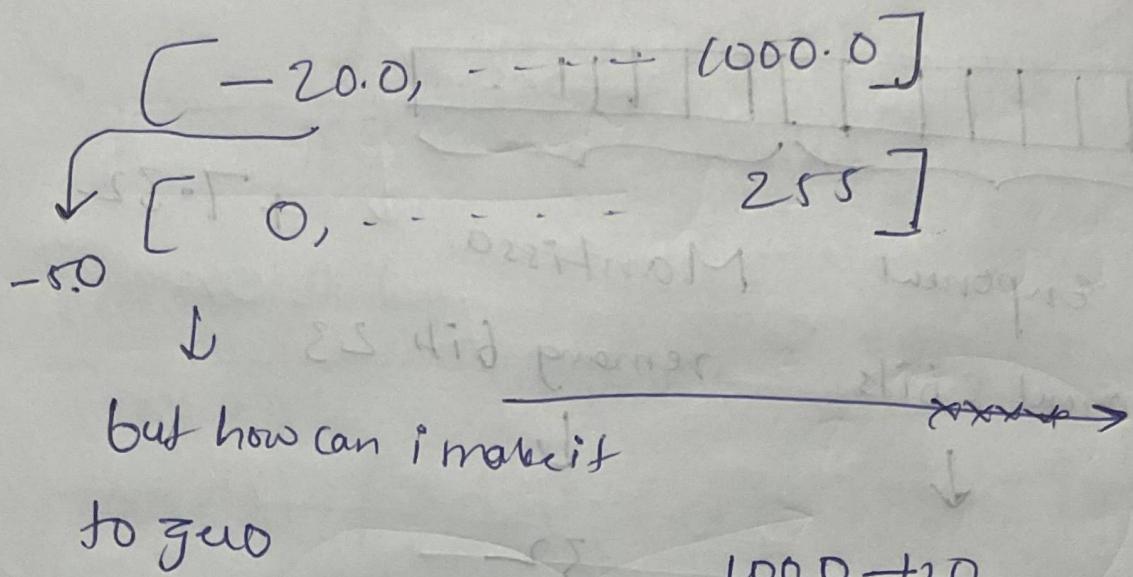
For FP 16



sign 5678 106.11
 Eng point mantissa

(ii) Assymmetric Quantization:

Assymmetric uint8 quantization



but how can i make it
to zero

$$\frac{1000 - 20}{255} = 4.0$$

Convert -20

$$\text{round}\left(\frac{-20}{4}\right) = (-5) + 5 = 0$$

1/1001

zero point

which is not on scale

that why we are adding zero point
of same number.

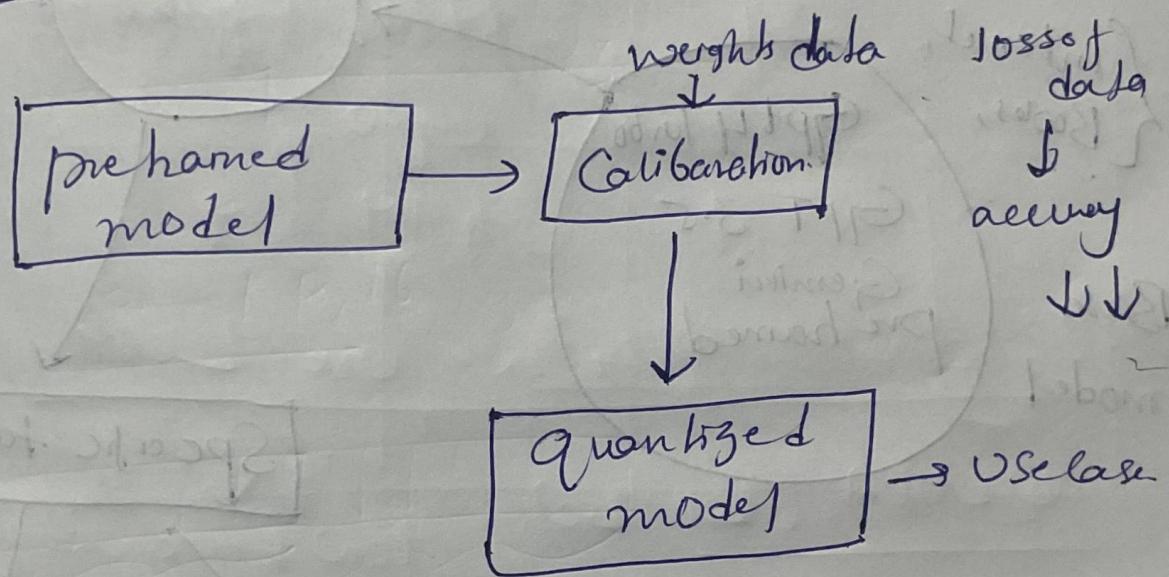
② Calibration:

Squeezing the values from higher to lower

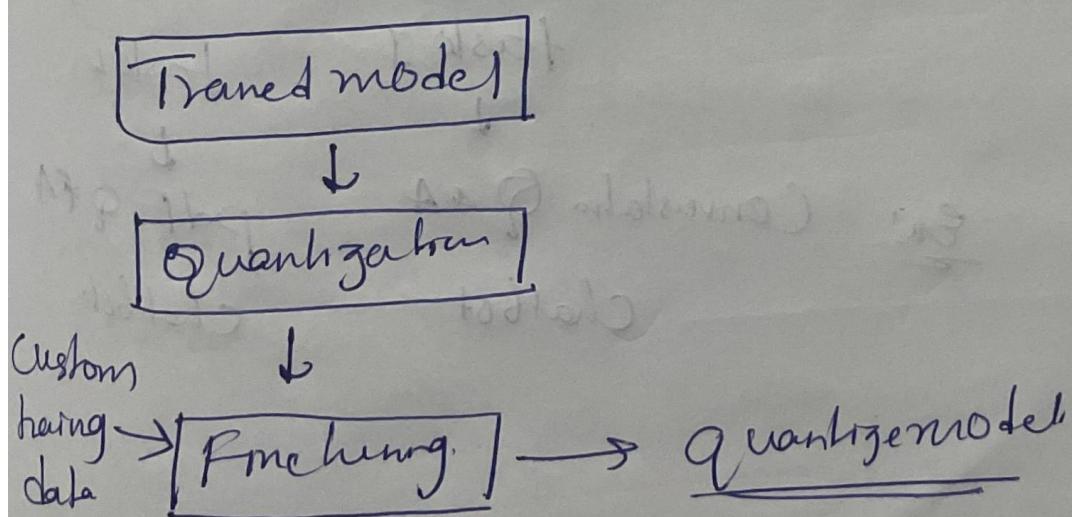
is called calibration.

③ Modes of Quantization:

① Post training Quantization (PTQ)

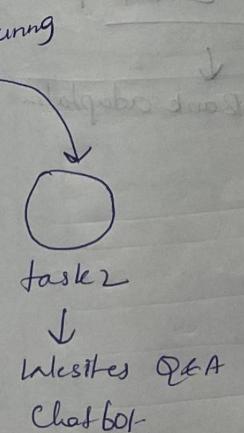
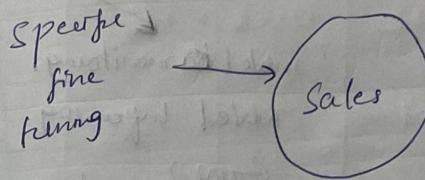
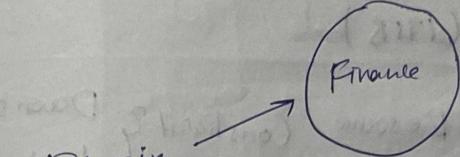
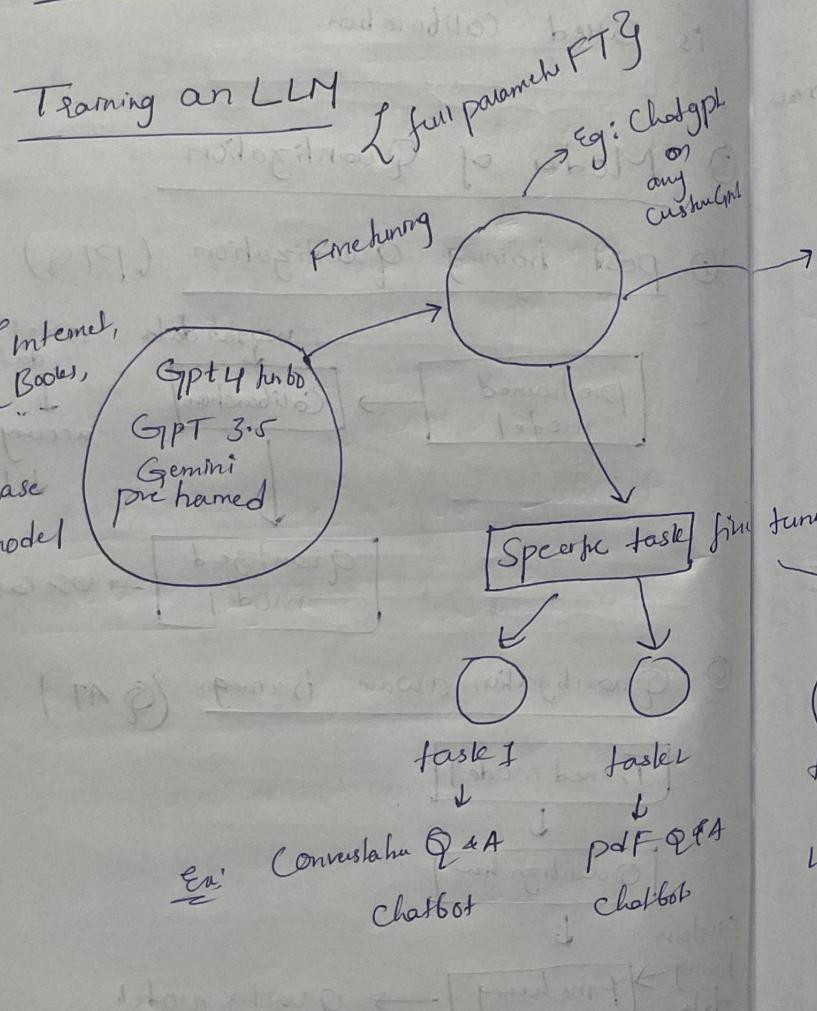


② Quantization aware training (QAT)



LORA + QLORA under mathematical

intuition: \downarrow and relevant example



Full Parameter Fine Tuning [challenges]

① Update all model weights

(175B)

② hardware Resource

Constraint } {

Downstream
task

model monitoring,

model inferencing

GPU } constraints
RAM }

To overcome these we use LORA &

QLORA [KORALO]

↓
low Rank adaptation

↓
Quantization low rank adaptation.

What does LORA DO? Instead of updating weights, it adds changes.

pretrained weights

3×3

LORA hacked wts

1	2	3
4	5	6
7	8	9

3×3

=

3×3

Fine tuned wts.

W_0

+

ΔW_0

\uparrow

$W_0 + \Delta W_0$

B

1
2
3

3×1

A

4	5	6
---	---	---

1×3

Rank = 1

No of learnable parameters

Rank	7B	13B	70B	180B
1	167k	228k	529k	849k
2	334k	456k	1M	2M
8	1M	2M	4M	7M
16	3M	4M	8M	14M
512	86M	117M	270M	434M

when to use high rank

When your model wants to learn
complex things

End to End Fine Tuning LLM Models

With Laini Platform: Project 16

Laini platform:

is a developer friendly platform that enables fast, efficient fine tuning and deployment of custom large language models using your own data on any hardware.

Project 16: End to End Fine tuning LLM

models with Lamini platform

Steps

① Define the data

- list of dictionaries
- input, output
- this is your training data.

② Import lamini library

```
import lamini
```

③ Set your api key

```
Lamini.api-key = " - - - "
```

④ Create a Lamini LLM object:

```
llm = Lamini(model="Lam3-8B-1shot")
```

⑤ Load the data:

```
data = getdata()
```

⑥ Fine tune the model:

```
llm.tune(data_or_datatxt=data,  
          finetune_args={"learning-rate": 1.0e-4})
```