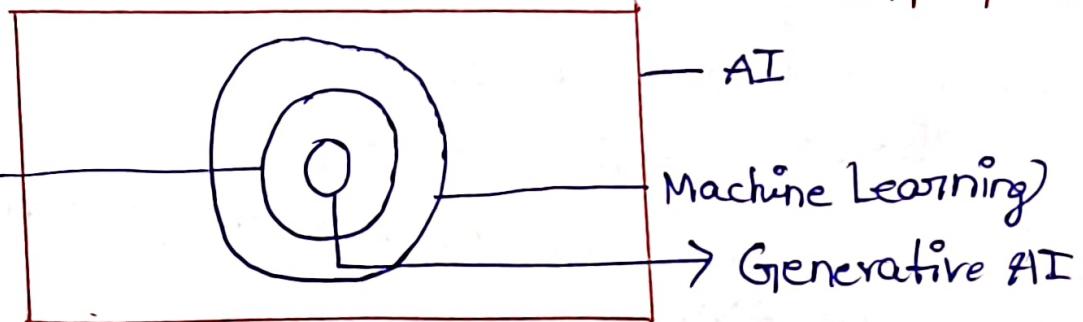


Generative AI

27/01/2025

Deep Learning)



AI → Build applications that can perform its own task with human intervention.

Ex: Netflix → Recommendation system
Self driving cars

AI Engineer → Integrating the model involves fine-tuning etc...
→ creating AI product.

ML → Machine Learning) is the subset of AI.

→ to provide statstools to perform:-

✳ Statistical Analysis, visualization, prediction and forecasting).

→ Data injection → Data transformation

→ Feature Engineering

→ understanding the data.

DL → Built to mimic human brain

→ Multi Layered Neural Network.

DL → ANN

→ CNN and object detection

→ RNN and its variants

CNN → Computer vision

RNN → Text Related usecases

↳ time series usecases.

ANN → We can train ML problem statement with help of ANN.

{transformers and Bert } Very Advanced.



{Generative AI }

Generative AI → Subset of Deep Learning

models :-

→ Discriminative Model

→ Generative AI or Model

Discriminative Model

Generative AI

Task → classification,
prediction

Task → Generates new data
trained on huge dataset

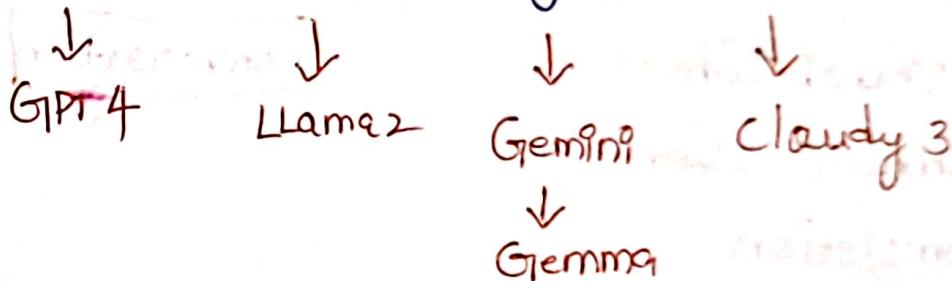
Dataset → Trained on
Labelled Dataset

In Gen. AI two types of models:-

① LLM → Large - Language Model \Rightarrow Text Data

② LIM → Large - Image Model \Rightarrow Images

LLM :- openAI, Meta, Google, Anthropic



(T42) \rightarrow Foundation Models \Rightarrow Trained on Huge Data
 \rightarrow pretrained Models

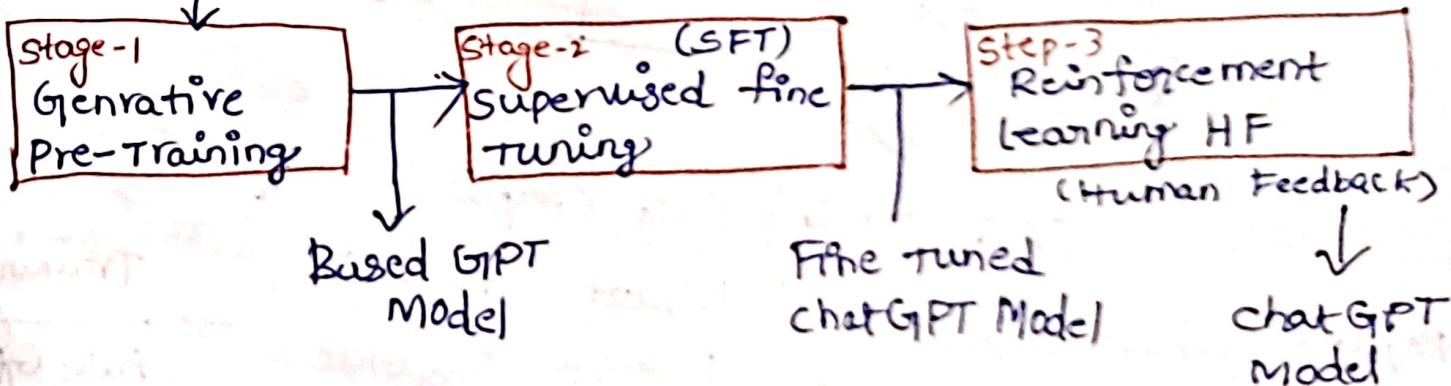
Langchain :- It is a frame work

RAGI Application :- **Retriever Argumentation Query**

↳ Retriever Argumented Generation. \Rightarrow develop chatbot

How, chat GPT is trained :-

Internet Data



① Stage-1 \Rightarrow [Generative Pre Training]

Internet Huge Input Text Data \rightarrow **Transformers** \rightarrow Base GPT Model

Tasks:-

What we want ?

① Language Translation

conversation \Rightarrow chatbot

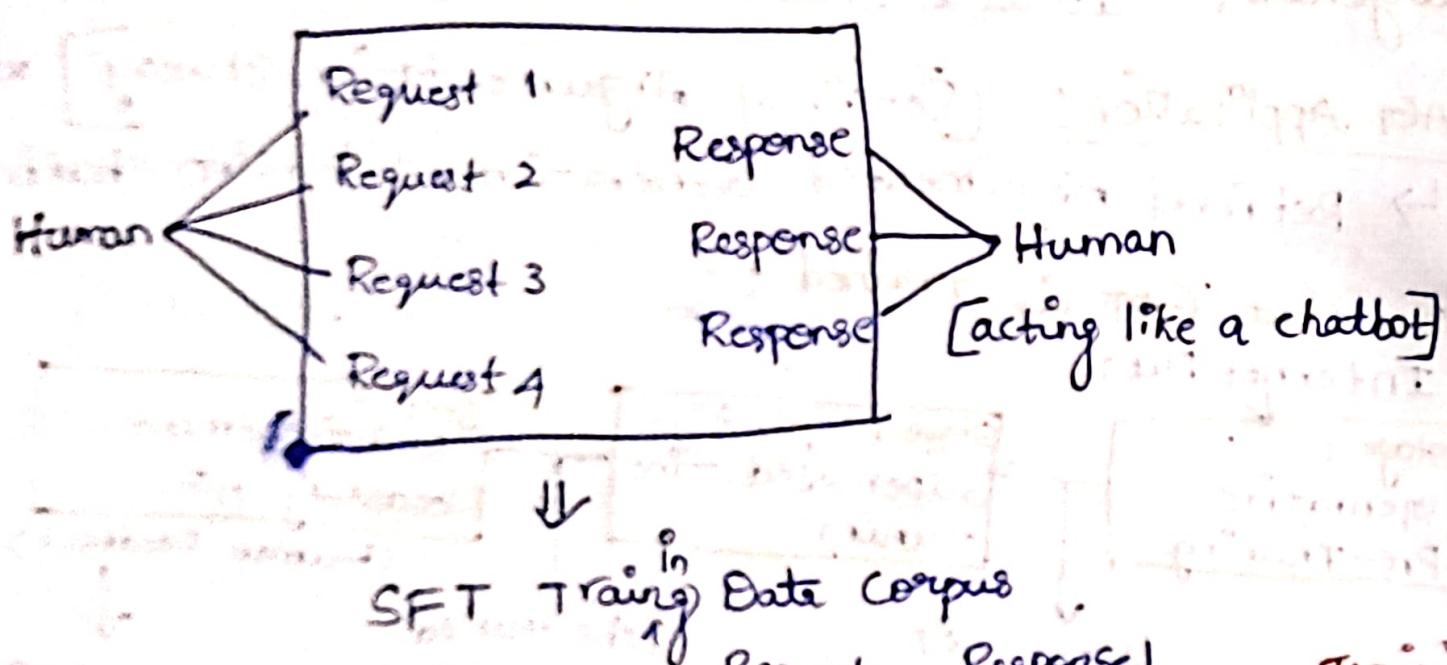
② Text Summarization

③ Text Completion

④ Sentiment Analysis

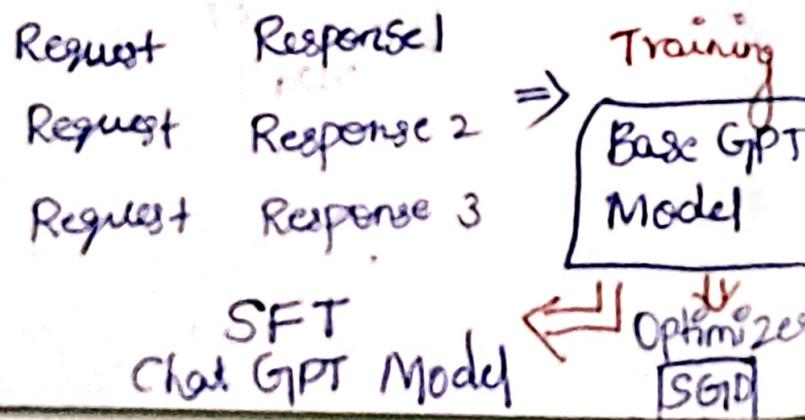
② Stage-2 [Supervised Fine Tuning] (SFT)

Real Conversation

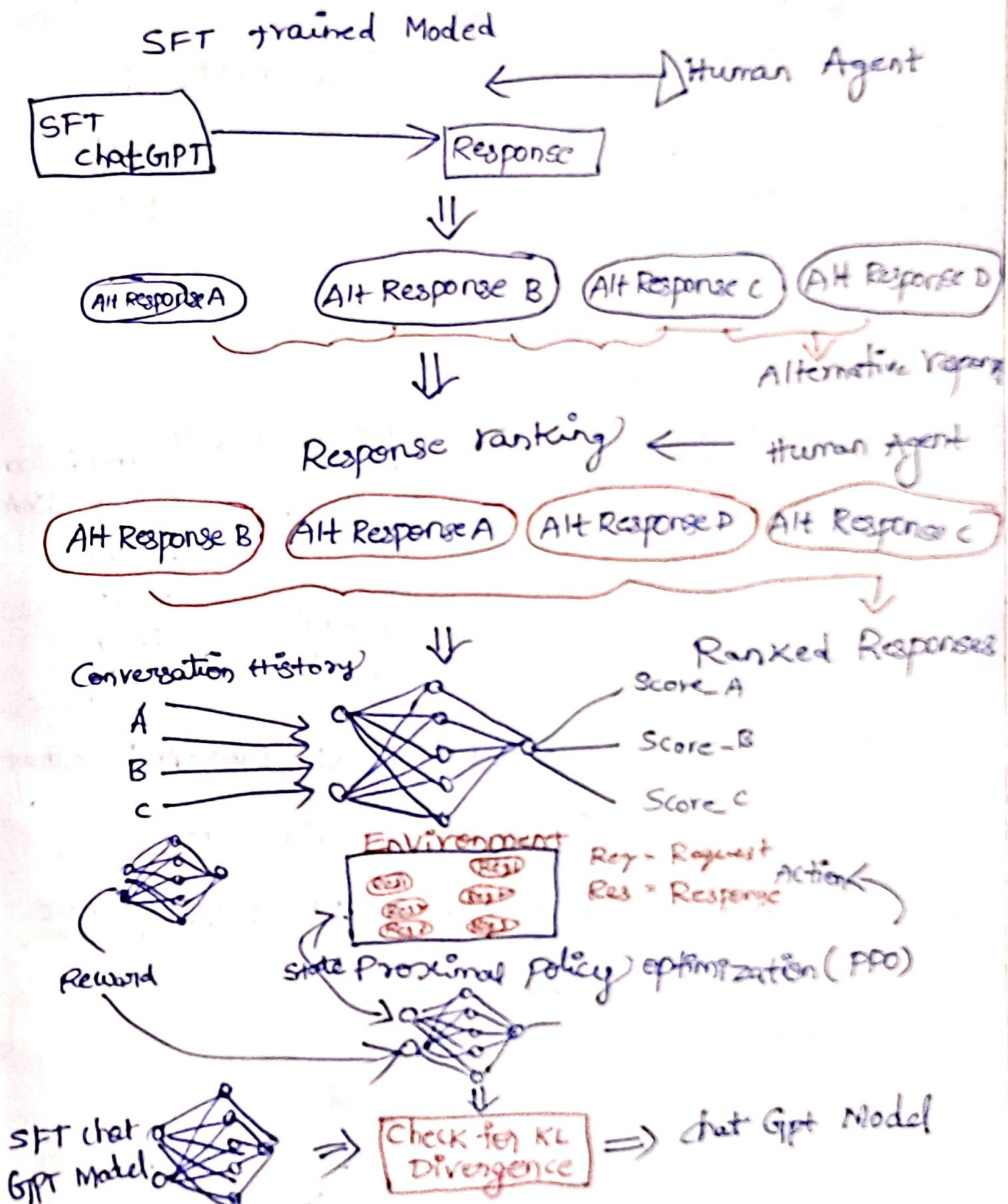


Request = conversation history

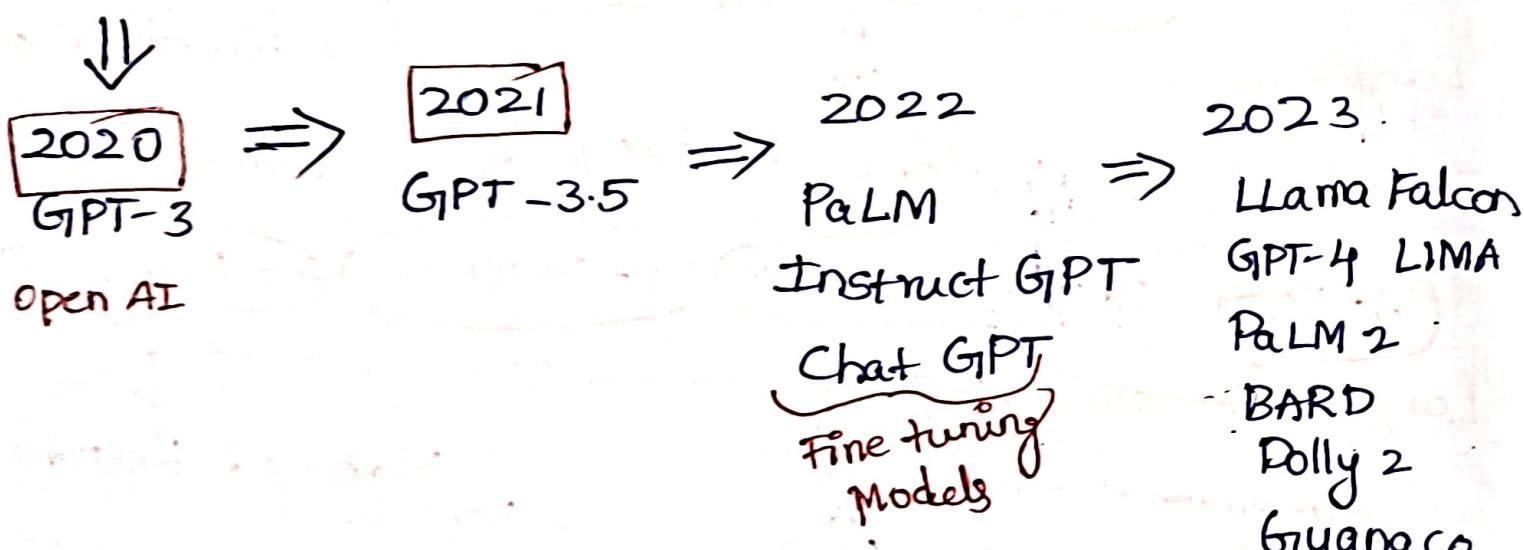
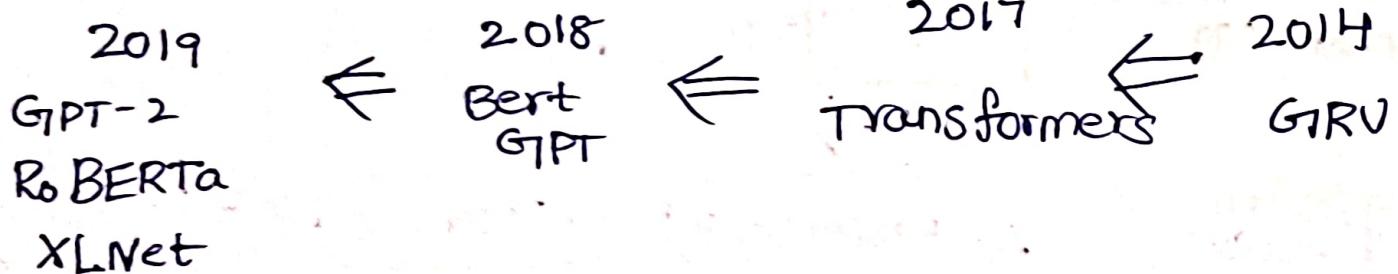
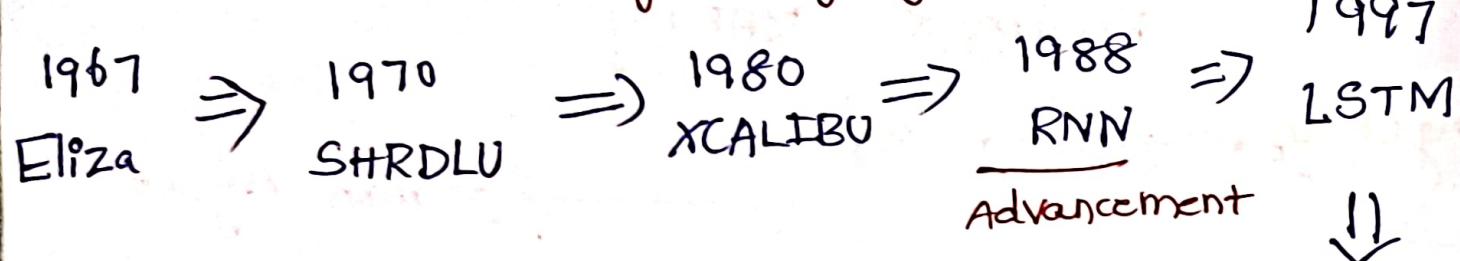
Response = Best Ideal Response



③ Stage - 3 = [Reinforcement Learning through Human Feedback (RHF)]



Evaluation of LLM (Large Language Models)



2024 → openAI → Multimodal → GPT 4.0 ← NLP → Images

2024 → Google Gemini Pro 1.5 } paid Models

2024 → Llama 3, Anthropic Claude opus, huggingface

⇒ Google Gemma 2

Langchain Ecosystem :-

→ why Langchain

→ Langchain Ecosystem ?

⇒ Langchain is the most common framework that is specifically used to build Generative AI application

Langchain :- open source

→ Frame work to develop Gen AI app.

Gen AI APP → LLM OPS [Debugging, Evolution, Monitoring]

Langsmith → Debugging → play round → Evolution
{observability} → Annotation Monitoring

LangServe → chaining as Rest API's
{Deployments}

Templates → Reference Applications

LangChain → chaining, Agents, Retrieval Strategies
{Cognitive Architecture}

Langchain - Community → Agent Tooling
{Integration components} Retrieval
Model I/O

→ model

→ prompt

→ Example selector

→ Output parser

→ Retriever

→ Document loader

→ Vector Store

→ Text Splitter

→ Embedding model

→ Tool

→ Toolkit

Langchain Core → { protocol } //

LCEL → Langchain Expression Language

→ Parallelization → Fallbacks → Tracing → Batching

→ Streaming → Async → Composition.

→ open 'cmd' with command prompt
→ cd 'folder path' & hit enter to go into folder
→ d:
→ code... (code space dot)

.env → OPENAI_API_KEY = "your api key" → LANGCHAIN_API_KEY = "" → LANGCHAIN_PROJECT = "your project name" → components of Langchain:-

① RAG → Retrieval Argumented Generation

loading the dataset → Diagram / PDF / Image / URL / JSON

[Data Ingestion]



Text or document

splitting the data into smaller, chunks

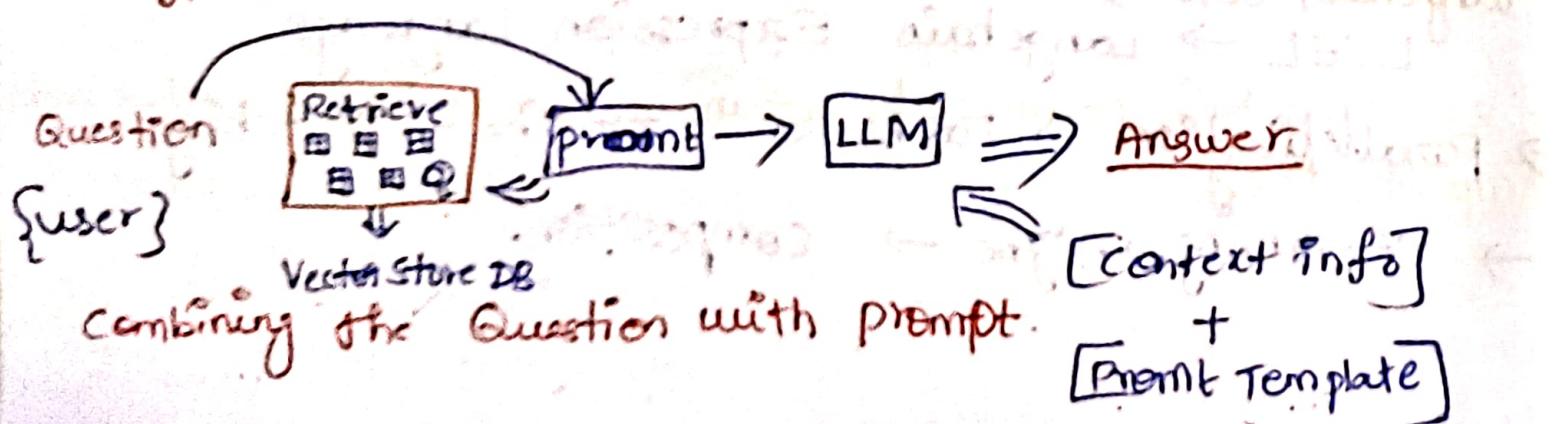


Embedding the text into vectors



storing the data into the vector store DB.

Types:- FAISS / CROMA DB / ASTRAD DB



→ Retrieval chain



Interface (responsible for querying the vectorstore DB)



Vector Store DB

Data - Ingestion :-

```
from langchain_community.document_loaders import TextLoader
```

```
Loader = TextLoader("speech.txt")
```

```
text_documents = Loader.load()
```

{Requirements.txt}

→ langchain

→ python-dotenv

→ ipy kernel

→ langchain_community

→ PyPDF

→ bs4

BeautifulSoup

```
from langchain_community.document_loaders import PyPDFLoader
```

```
Loader = PyPDFLoader('attention.PDF')
```

```
PDF_Documents = Loader.load()
```

```
from langchain_community.document_loaders import WebBaseLoader
```

```
Loader = WebBaseLoader('url')
```

```
webLoader = Loader.load()
```

BS4 => It's a python library used for web scraping and parsing HTML or XML documents.

```
from langchain_community.document_loaders import WebBaseLoader  
import BS4
```

```
Loader = WebBaseLoader(Web_Path=(url), )
```

```
bs_kwarg = dict(parse_only=BS4.SoupStrainer  
(class='post-title', post-content', post-header'))
```

WebBaseLoader = loader.load

Kwargs \Rightarrow keyword Arguments

BS4 Keyword arguments \rightarrow parse_only

filter \Rightarrow (SoupStrainer)

Loading Data Ingestion:-

\rightarrow Text

\rightarrow PDF

\rightarrow webURL

\rightarrow Arxiv

\rightarrow Wikipedia.

from langchain_community.document_loaders import
WikipediaLoader

WikiLoader = WikipediaLoader(query="Generative AI",
load_max_docs=2).load()

len(wiki_loader)

Text Splitter Techniques:-

- ① Recursive Text Splitter
- ② character Text Splitter
- ③ HTML Text Splitter
- ④ Recursive json splitter.

Types of Text Splitter Techniques.

EMBEDDING:-

① open AI — paid

② OLLAMA — open source

③ Hugging face — open source

Embedding Techniques

Embedding :-

Converting the text into vectors

- ```
=> { import os
 from dotenv import load_dotenv
 load_dotenv() # Load all environment variables}
```
- ```
=> os.environ["OPENAI_API_KEY"]  
    = os.getenv("OPENAI_API_KEY")
```
- ```
=> from langchain_openai import OpenAIEmbeddings
embeddings = OpenAI OpenAIEmbeddings(model="text-embedding-
-3-large")
embeddings.
```
- ```
=> text = "This is a tutorial on OPENAI Embedding"  
query_result = embeddings.embed_query(text)  
query_result => Convert Text into vectors
```
- ```
=> len(query_result) => 3072
```
- ```
embeddings_1024 = OpenAIEmbeddings(model="text_embedding-  
-3-large", dimensions=1024)
```
- ```
=> text = "This is a tutorial on OPENAI Embedding"
query_result = embeddings_1024.embed_query(text)
```
- dimensions =>
- Specifying the size of the Embedding vectors.

Langchain\_community). VectorStores

⇒ FAISS ⇒ Chroma ⇒ Pinecone ⇒ Weaviate

Ollama :- open source LLM models [platform]

→ Llama3, phi3, Mistral, Gemma

open source model ⇒ Langchain

In open source model Embedding techniques ⇒ ollama

Vector store DB :-

⇒ FAISS ⇒ Chroma DB ⇒ Astra DB

FAISS :-

```
from langchain_community.document_loaders import TextLoader
from langchain_text_splitters import characterTextSplitter
from langchain_community.Embeddings import OllamaEmbeddings
from langchain_community.VectorStores import FAISS
```

⇒ Loader = TextLoader('Speech.txt')

documents = Loader.load()

text\_splitter = characterTextSplitter(chunk\_size=200,  
 chunk\_overlap=30)

docs = text\_splitter.split\_documents(documents)

⇒ Embeddings = OllamaEmbeddings()

db = FAISS.from\_documents(docs, embeddings)

db

## # Querying in FAISS :-

```
query = " " " I am a man" query = db.similarity_search(query)
docs = db.similarity_search(query)
docs[0].page_content
```

For this we get same data.

① Query

② Retriever

## # Retriever :-

```
retriever = db.as_retriever()
docs = retriever.invoke(query)
docs[0].page_content
```

## # similarity\_Search\_with\_Score

```
docs_and_score = db.similarity_search_with_score(query)
docs_and_score
```

## # Embedding Vector :-

```
embedding_vector = Embeddings().embed_query(query)
embedding_vector
```

## => Components of Langchain :-

Import os

```
from dotenv import load_dotenv
load_dotenv()
```

## => Langsmith Tracking :-

```
os.environ["LANGCHAIN_API_KEY"]
```

```
= os.getenv("LANGCHAIN_API_KEY")
```

```
os.environ["LANGCHAIN_TRACKING_V2"] = "true"
```

```
os.environ["LANGCHAIN_PROJECT"] = os.getenv("LANGCHAIN_PROJECT")
```

```
=> from langchain_community.llms import Ollama
llm = Ollama(model="gen llama 2:7b")
print(llm)
=> response = llm.invoke("Explain Langchain in simple terms.")
=> response
```

# Chat prompt Template :-

```
from langchain_text_prompts import ChatPromptTemplate
prompt = ChatPromptTemplate.from_message([
 ("sys", "don"),
 ("user", "{ }")
])
Prompt.
```

# Chain :-

```
chain = prompt | llm
```

```
response = chain.invoke({"input": "you can tell me about
 longsmith?"})
print(response)
```

Simple App with Ollama - llama → Gen AI

```
import os
```

```
from dotenv import load_dotenv
```

```
load_dotenv
```

```
from langchain_community.llms import ollama
import streamlit as st
from langchain_core.prompts import ChatPromptTemplate
from langchain_core.input_parsers import StrOutputParser
```

# Langsmith Tracking)

# prompt template

# streamlit - frame work

# Ollama · Llama2 Model

⊕ LCEL => Langchain Expression Language

open source models => Llama2, Gemma2, mistral } platform  
} Groq

```
from langchain_groq import ChatGroq
model = ChatGroq(model="gemma2-9b-it", groq_api_key
 ="groq-api key")
```

from langchain\_core.messages import HumanMessage, System  
Message

messages = [

SystemMessage(content=" " ),

HumanMessage(content=" \*? ")

]

result = model.invoke(messages)

result //

```
From langchain_core.output_parsers import StructuredOutputParser
parsers = StructuredOutputParser()
parser.invoke(result)
```

## #langchain Server:- (serve.py)

```
import ① FASTAPI
 ② chatprompt template
 ③ StructuredOutputParser
 ④ flatten Chat-Groq
 ⑤ os
 ⑥ dotenv → load_dotenv
 ⑦ add_routes (langserve)
```

**Postman** → To check translation in langchain server.

chatbot → ipynb:-

→ Building Chatbot with message history using Langchain

```
import os
from dotenv import load_dotenv
load_dotenv()
```

```
groq_api_key = os.getenv("GROQ_API_KEY")
groq_api_key
```

```
⇒ from langchain.chat import ChatGrover
model = ChatGrover(model="Gemma2-9b-it", grover_api_key="some")
⇒ from langchain.core.messages import HumanMessage
model.invoke([HumanMessage(content=" ")])
⇒ from langchain.core.messages import AIMessage
model.invoke([
 HumanMessage(content=" "),
 AIMessage(content=" "),
 HumanMessage(content=" ")
])
```

## Message History :- (from langchain community)

```
from langchain_community.chat_message_histories import
 ChatMessageHistory
from langchain_core.chat_history import BaseChatMessageHistory
from langchain_core.runnables.history import RunnableWithMessageHistory
```

Prompt chain :-

```
from langchain_core.prompts import ChatPromptTemplate,
 MessagePlaceholder
```

Prompt = ChatPromptTemplate.from\_messages(

```
[("System", " "),
 MessagePlaceholder(variable_name="messages")]
```

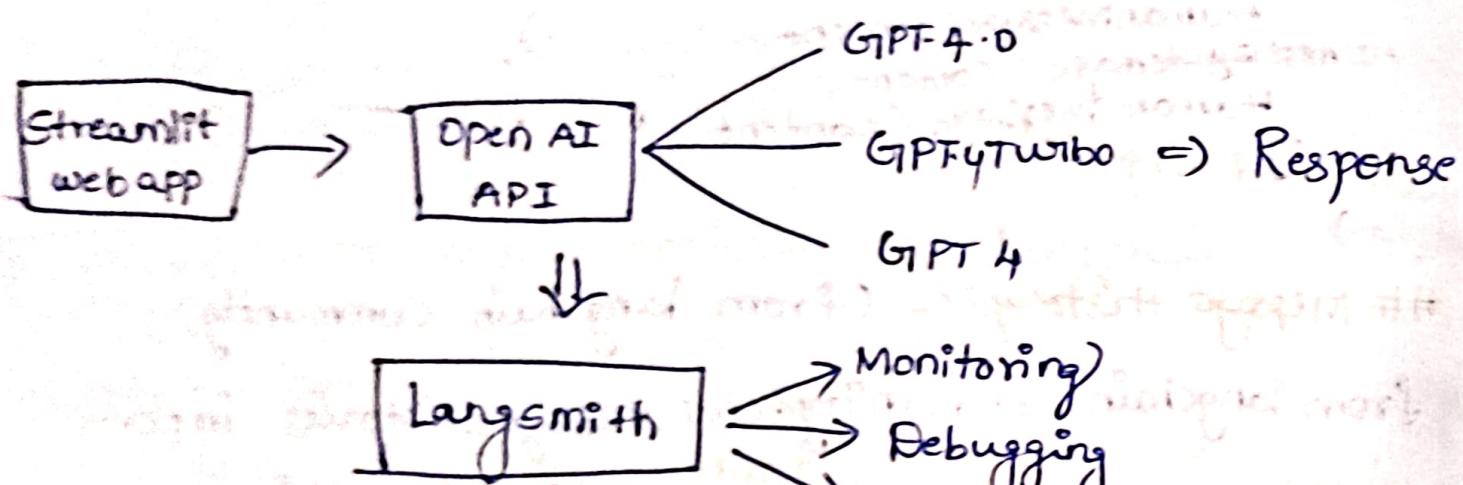
Trim\_message :-

Parameters → max\_tokens, strategy, token\_counter, include\_system, allow\_partial, start\_on

Q&A chatbot

① open AI LLM Models :- GPT-4.0, GPT-4 Turbo

② OLLAMA → open source → LLAMA3, GEMMA2, Mistral



- ① project
- ② Environment Variable
- ③ Required.txt
- ④ Streamlit web app
- ⑤ Deployment

Greg:-

Greg is a company known for developing ultra-fast AI accelerators and processors.

→ Their technology is designed to run large language models.

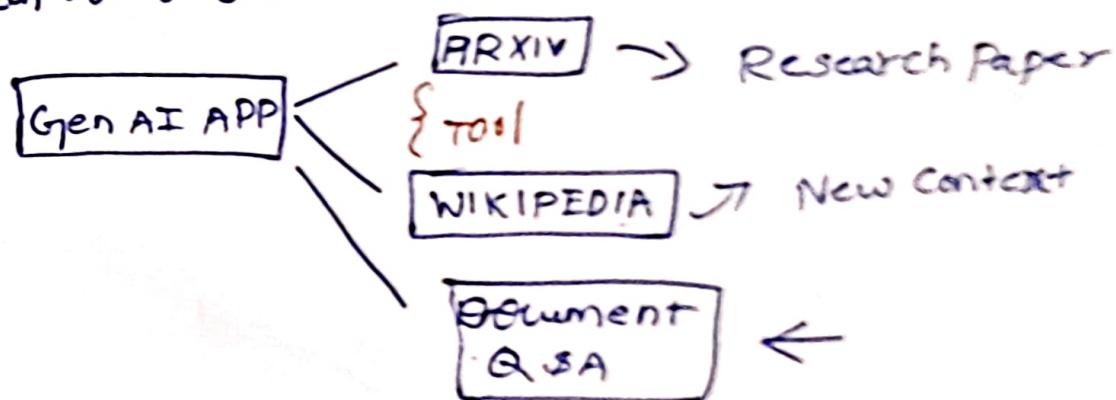
Rag:- [Retrieval-Augmented-Generation]

→ It is a smart AI technique that helps chatbots and AI Models give better, more accurate answers.

Tools and Agents:-

→ Tools are the interfaces that an agent, chain or ILM can use to interact with the world.

Eg:-



Pine Cone :-

- ① creating pinecone Account
- ② Sign in
- ③ Creating the project name
- ④ copy the API Key
- ⑤ creating the Index



configuration - Dimensions / Metric

Cosine  
Euclidean  
dot product

- ⑥ cloud provider
- ⑦ Region (aws - Virginia(open source))

# Azure Machine Learning:-

25/02/2025

Azure ML is a cloud based platform for

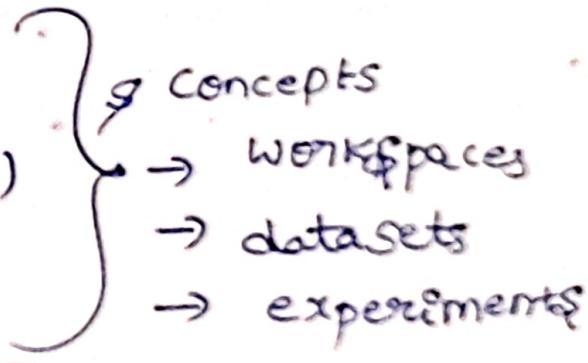
- Building
- Training
- Deploying
- Managing Machine Learning Models.

Key features of Azure Machine Learning:-

- ① AutoML → Automated ML / Model Selection / Hyperparameters / Tuning and Feature Engineering.
- ② ML Pipelines → ML workflows with reusable components.
- ③ Compute Options → Azure VMs, GPU clusters, Integration and Kubernetes for scalable F & D.
- ④ MLOps → CI/CD / Monitoring / Versioning.
- ⑤ Responsible AI → Interpretability

① Learn the Basics:-

- Azure ML Studio
- SDK (Python based)
- CLI



② Hands - on practice :-

- Azure ML studio ⇒ Auto ML Experiments
- Python - SDK ⇒ in Jupyter Notebook or Vs code
- Run notebooks on Azure ML compute Instances.

### ③ Build and Deploy models:-

- Train models using SkLearn, Tensorflow, Pytorch.
- Deploy Models as Azure ML Endpoints using Azure Functions or AKS (Azure Kubernetes Service).

### ④ Explore MLOps & Advanced Topics:-

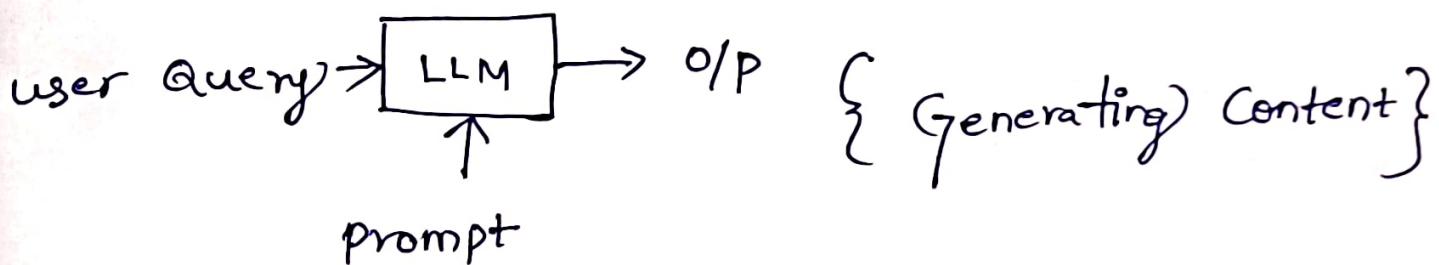
- Learn Azure ML Pipelines for model automation
- Use MLFlow
- Implement model monitoring.

Agentic-AI → AI Agents

- ① What is Agentic AI?
- ② How is it different from Generative AI?
- ③ Frameworks used in to Create Agentic AI?

Generative AI:-

Create Content → LLM Model



Agentic AI:- Autonomous AI System

→ Autonomous AI System that can be able to perform their own task.

⇒ To achieve a specific goal it is the AI System, will be working independently without any human intervention.

Langgraph — Framework

- It is a tool that helps to build AI workflows, where AI agents or steps work together in a structured way.
- It uses graph like system to manage decision-making, memory and interactions.

- It is useful for making smart chatbots, automation systems or AI assistants that needs to remember things and make decisions.
- In Generative AI the main aim is just to create the content.

Data → 100 log files or pdf's or text docs.

Step 0: ↴ reading the files from Azure Blob Storage.

Step 1: Data Extraction from log file or pdf

Step 2: chunking the data with the help of Recursive char splitter.  
(How will u decide the chunk size, how will u decide  
the chunk overlap?) .

Step 3: Embedding the data. (openai embedding).

Step 4: saving the Embedding in vector DB (Azure Cosmos DB) -

Prompt: Prompt Engineering completed.

LLM: for providing content to the LLM,

Based on the Prompt instruction LLM will give answer.

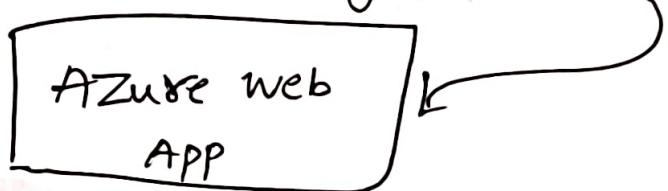
without  
RAG.

## Deployment in Azure

- CI/CD pipeline
- Git hub actions

① Web app → Docker image ⇒ private image

② Docker Image → container Registry (Azure)



⊕ Container Registry

Step 1 → Create Resource group - testdockervenkatesh  
Registry name - testdockervenkatesh

Step 2 → Access Key

Admin user → enabled

Step 3 → Web app

Select different Resource group

Name - testdockervenkatesh

Docker Container

→ Central us

Docker → Single container/ Azure Container Registry

Registry - testdockervenkatesh

Terminal :-

Creating the docker image → docker build -t testdockervenkatesh, azurecr.io/llmtest:late

Login into the docker → docker login testdockervenkatesh, azurecr.io

push the files into the docker → docker push testdockervenkatesh

username :- testdockervenkatesh

password :- from Azure Access Key

⇒ Create web app

---

① Container Registry — Store docker image

② Azure Web App

- ① What is Docker
- ② What are Containers
- ③ Containers Vs Virtual Machine
- ④ Docker Image Vs Container
- ⑤ Practical Implementation of Docker

Why Containers :-

Myself Venkatesh,

IT

- I do have exposure overall 3.1 years of IT Experience.
- Where I started my career as a data scientist.
- My previous company is Accenture, for 3.1 years I worked over there.
- When it comes to machine learning, DL and Gen AI I do have overall 3.1 years of experience
- So there I do have exposure to ML, DL, GenAI, NLP and at the same time connecting with the client
- taking the requirement from the client
- Doing the data collection -②
- understanding the problem statement -①
- Data preprocessing -③
- Data cleaning / Manipulation / Feature Engineering -④
- Model Building / Model Training / Model validating
- checking the accuracy
- Developing the production Ready code
- pushing the code to the Code cloud Repositories
- from there maintaining the repositories over there
- Maintaining the code in the feature Branches
- from their integrating the feature Branches with the develop branches.

- In Gen AI I do have exposure ~~of all 21 years~~ from the past 2 years, I have been working with Gen AI related projects.
- When it comes to my specialization NLP is my Strength and I am very good at gen AI dealing with Rag pipeline
- And finally when it comes to Cloud Experience
- I do have exposure of using Azure Data bricks
- Azure ML → Azure function → Azure Blob Storage
- Running the LLM Models from the Azure prospectus
- Connecting Azure data bricks with the Azure functions, like that I do have exposure till deployment of the models also.

⇒ Deployment type :-

- I do have Fast API and Flask API

Sometimes the frontend team will ask the predictions in the Tabular format and sometimes they will ask in JSON format

- Based on their requirement I will be providing the results to them.

When it comes to the Azure data bricks deployment so whatever the main function we are creating in the databricks we will be calling that function and we will be executing.

- once executed we will get job run id
- that job run id will be passed to the Azure func frontend program.
- creating Github actions
- Creating CI/CD pipelines
- Docker / Kubernetes that will be taken care by devops team.

project :-

①

- ① Bias and Variance
- ② cross validation and importance
- ③ overfitting and underfitting
- ④ How to handle missing data
- ⑤ techniques used in feature selection
- ⑥ what is Recursive feature Elimination
- ⑦ Bagging Vs Boosting
- ⑧ Model performing well but not in production  
how to handle.
- ⑨ classification metrics when to use for which type of dataset
- ⑩ . Difference between llm and NLP
- ⑪ What is LLM
- ⑫ Differences between fine-tuning and Prompt Engineering
- ⑬ NER in NLP ?
- ⑭ joins Concept in SQL .