

Exploring LangChain:

Revolutionizing the Future of AI and Language Models

🔍 What is LangChain?

An open-source framework that simplifies the development and deployment of applications powered by LLMs.

- Chain multiple LLMs together for multi-step reasoning
- Integrate external data sources like APIs and databases (RAG – Retrieval Augmented Generation)
- Build dynamic workflows that are contextually aware
- Automate tasks and decision-making with agents

💡 A Real-World Example:

Building a customer support chatbot using LangChain:

- Chain models together for multi-layered processing (e.g. NLP and sentiment analysis)
- Use memory to remember customer interactions, creating a personalized experience
- Let indexes to fetch relevant data from your knowledge base

✚ 6 Key Components of LangChain



Models

The backbone of LangChain



Prompts

are instructions or queries provided to models



Chains

a sequence of operations where each step



Indexes

help optimize retrieval of relevant information



Memory

Allows LangChain applications to "remember" prior interactions or states



Agents

autonomous decision-making systems (e.g., creating ticket, pulling knowledge articles, etc.)

LangChain Alternatives?
LlamaIndex HayStack

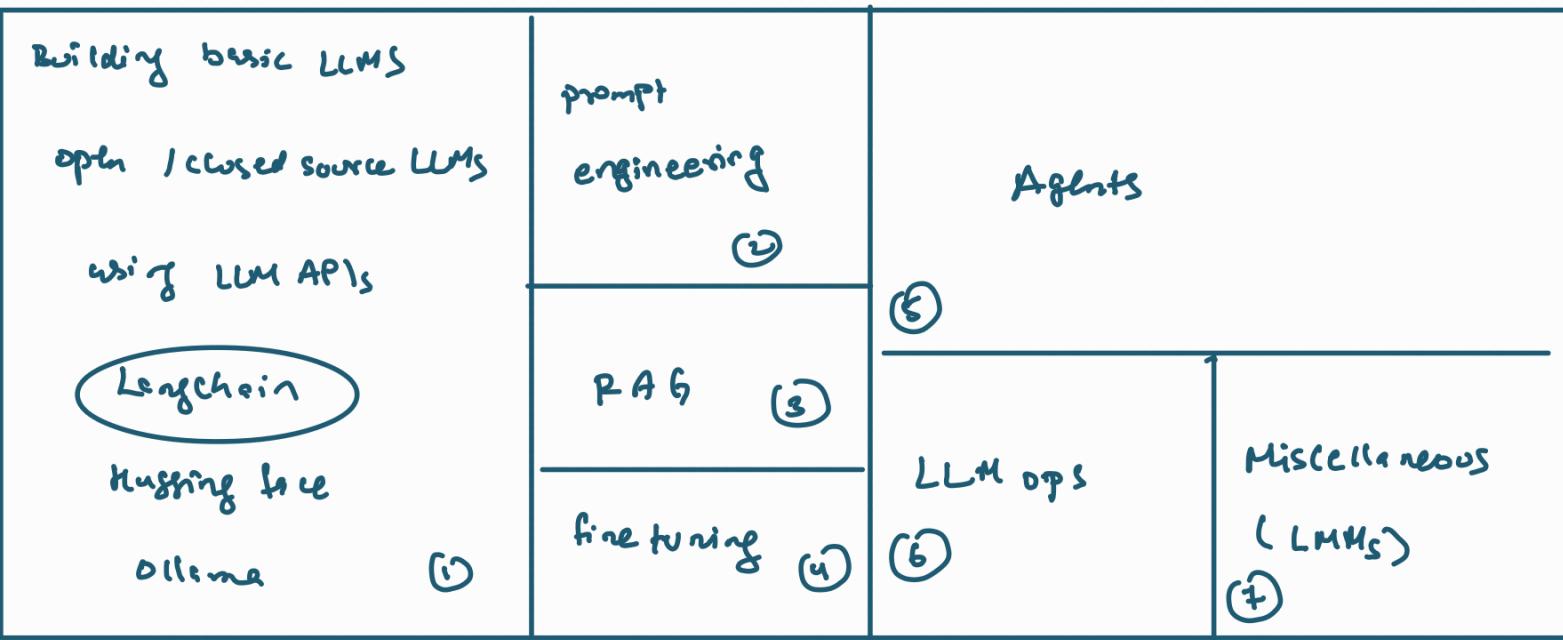
What is Langchain?

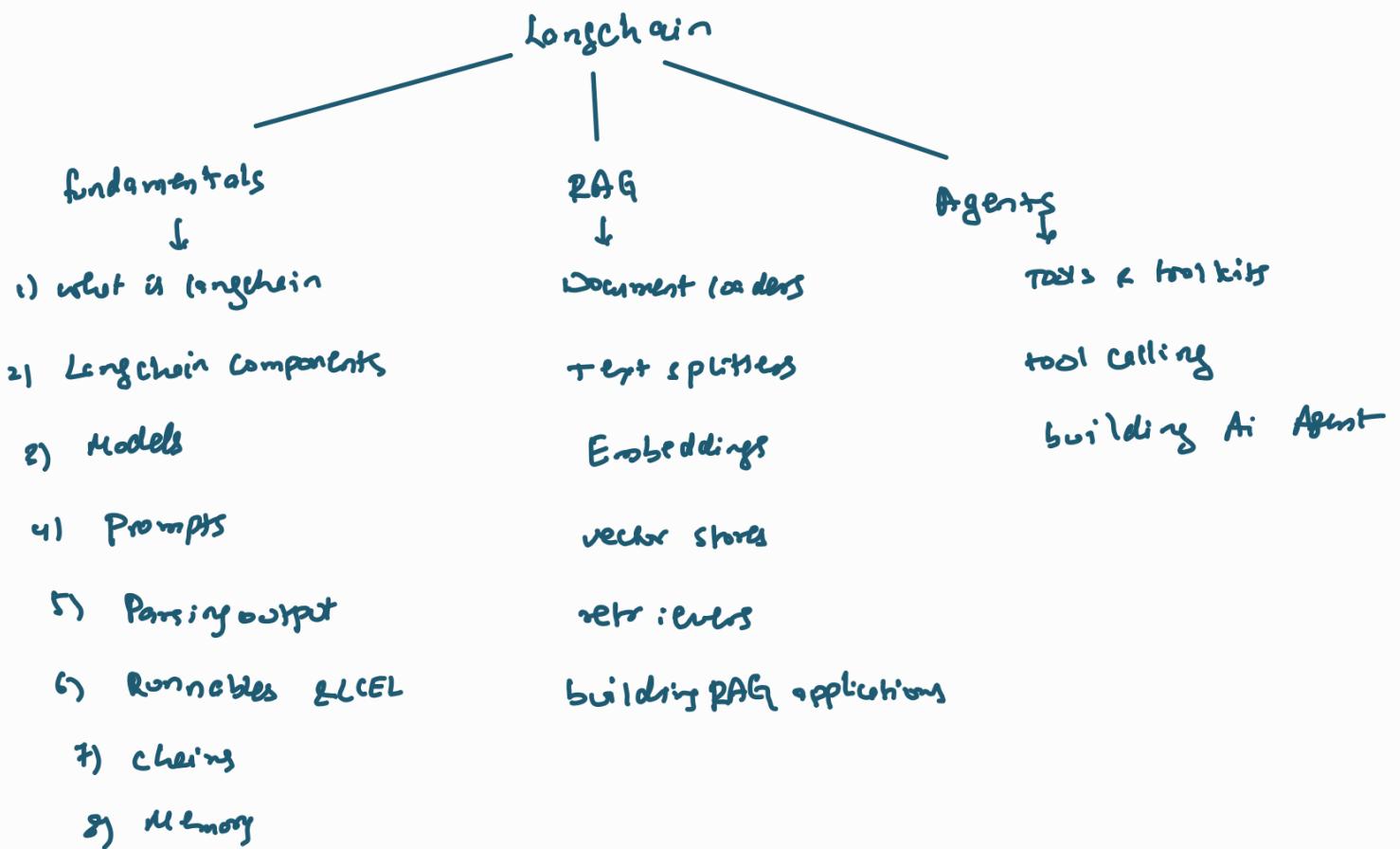
Langchain is an open source framework that helps in building LLM based applications.

It provides modular components and end-to-end tools that help developers build complex AI applications such as Chatbots, question answering systems, retrieval Augmented generation (RAG), autonomous agents and more.

- i) Supports all major LLMs
- ii) Simplifies developing LLM based Applications
- iii) Integrations available for all major tools
- iv) Supports all major GenAI use cases.
- v) Open source / free / Actively developed.

GenAI with Langchain Roadmap



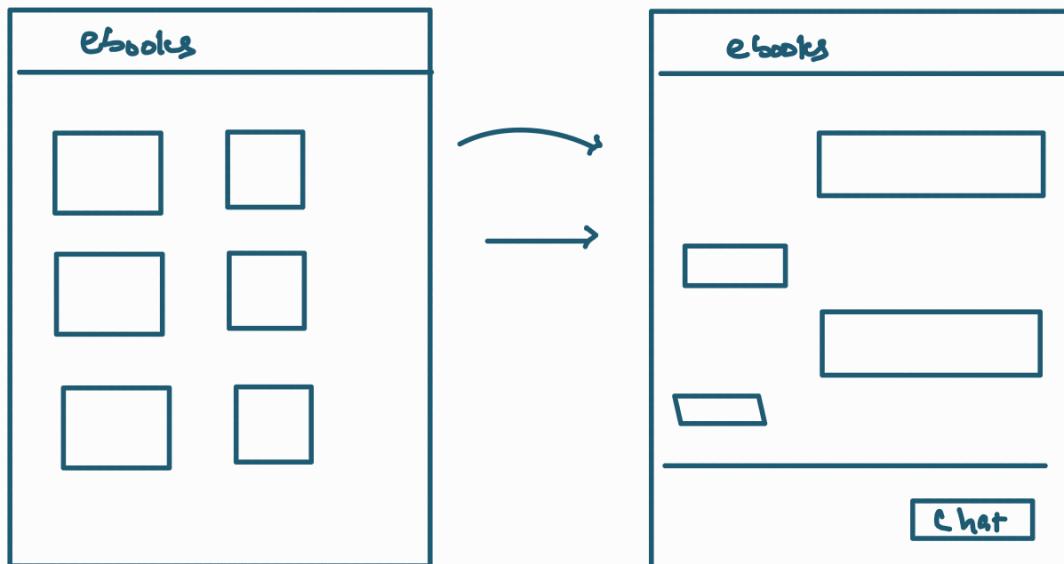


Intro to Langchain :

→ Langchain is an open source framework for developing applications powered by LLMs

Why Langchain? To make AI applications faster, scalable, model agnostic

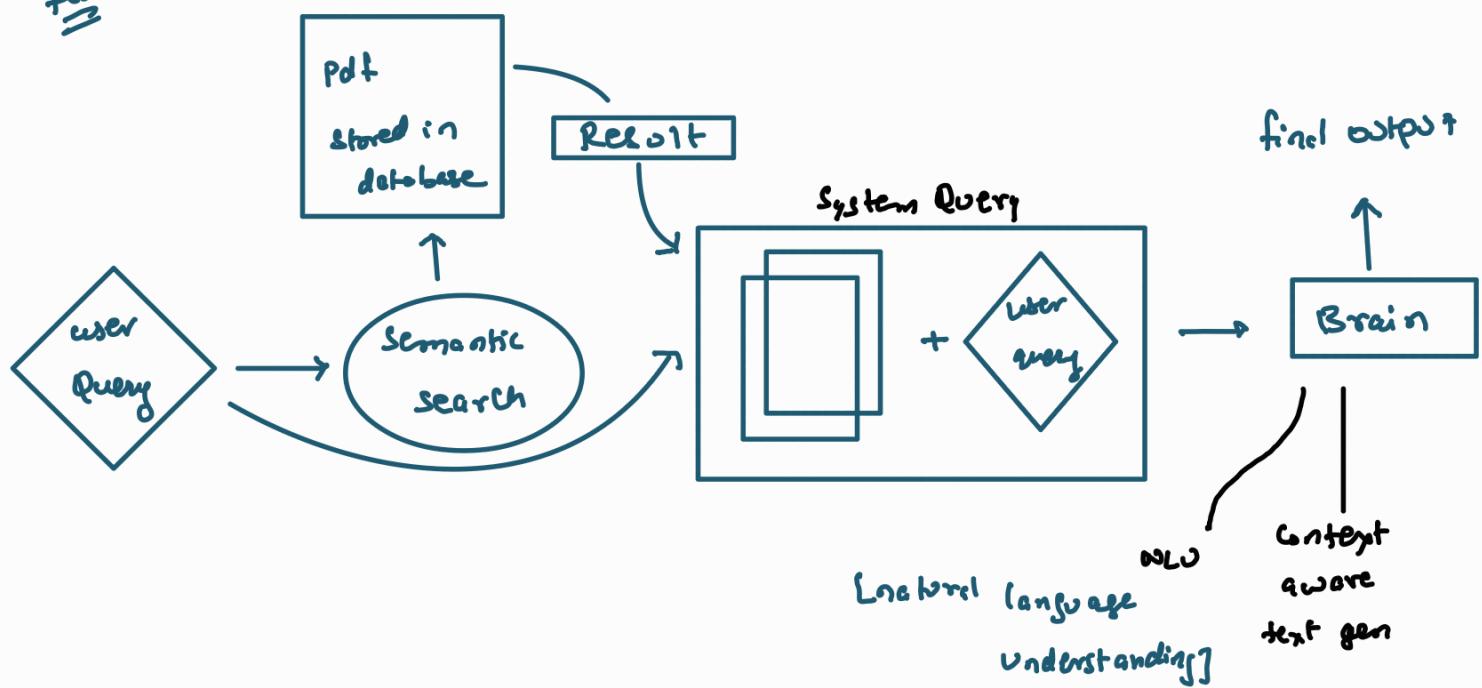
e.g.



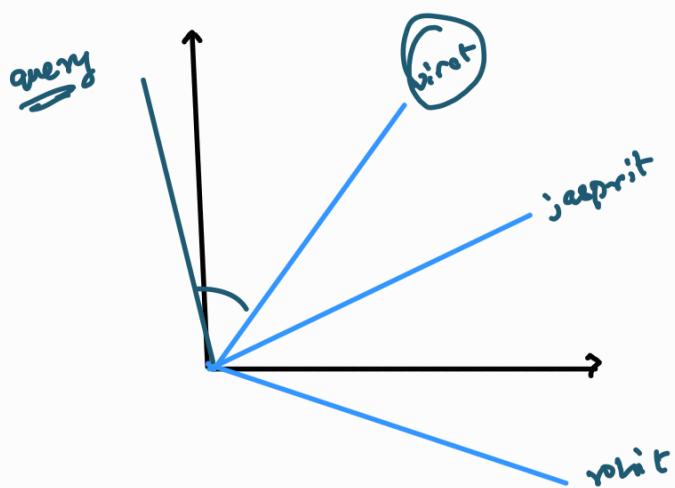
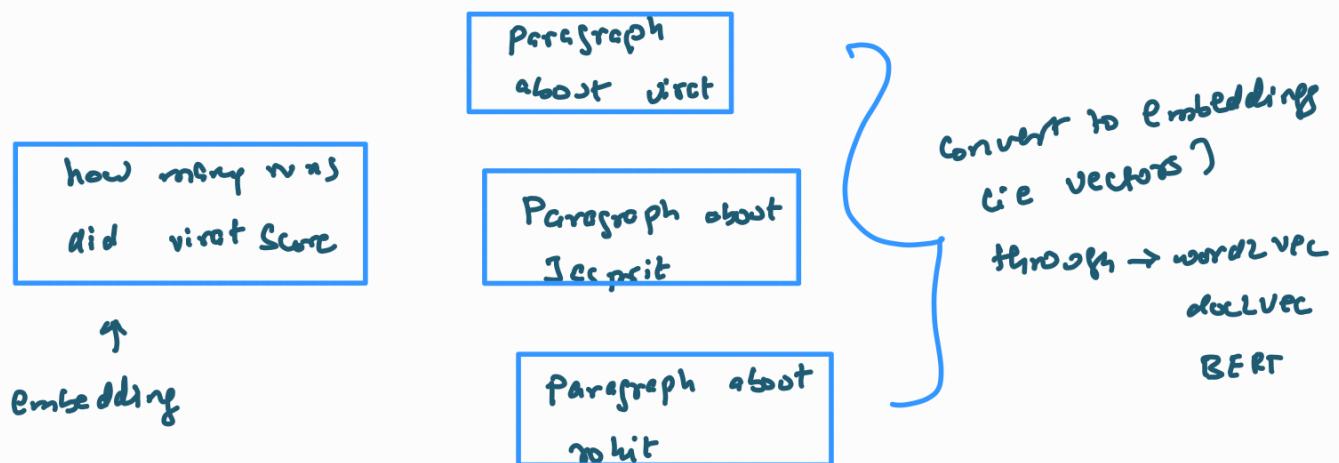
chat with ebooks

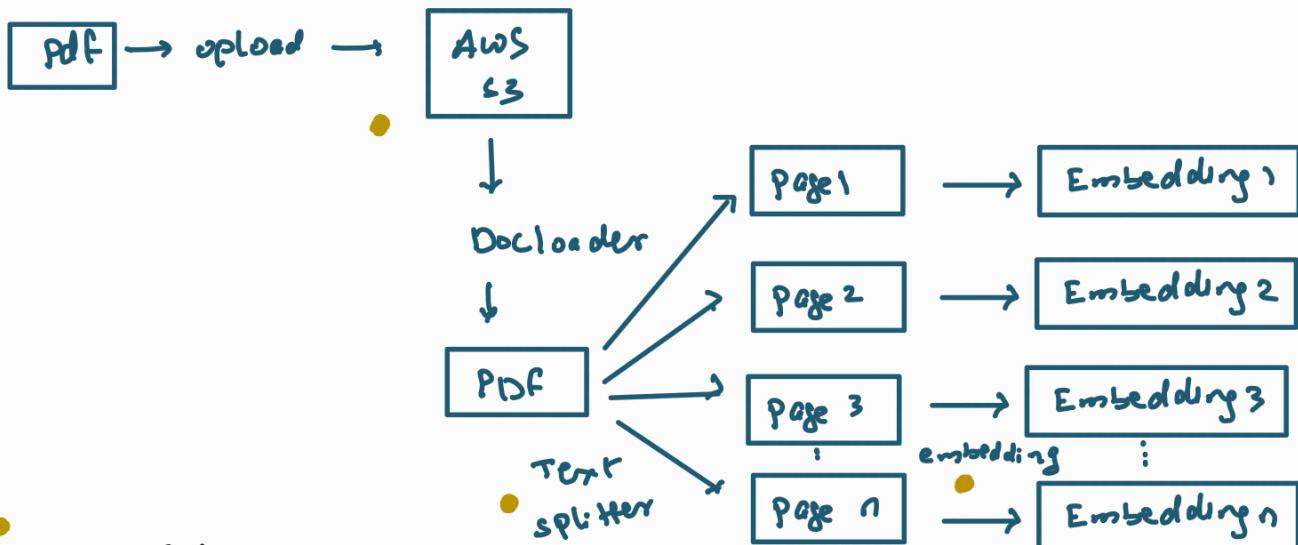
to create an application to chat with ebooks directly.

flow

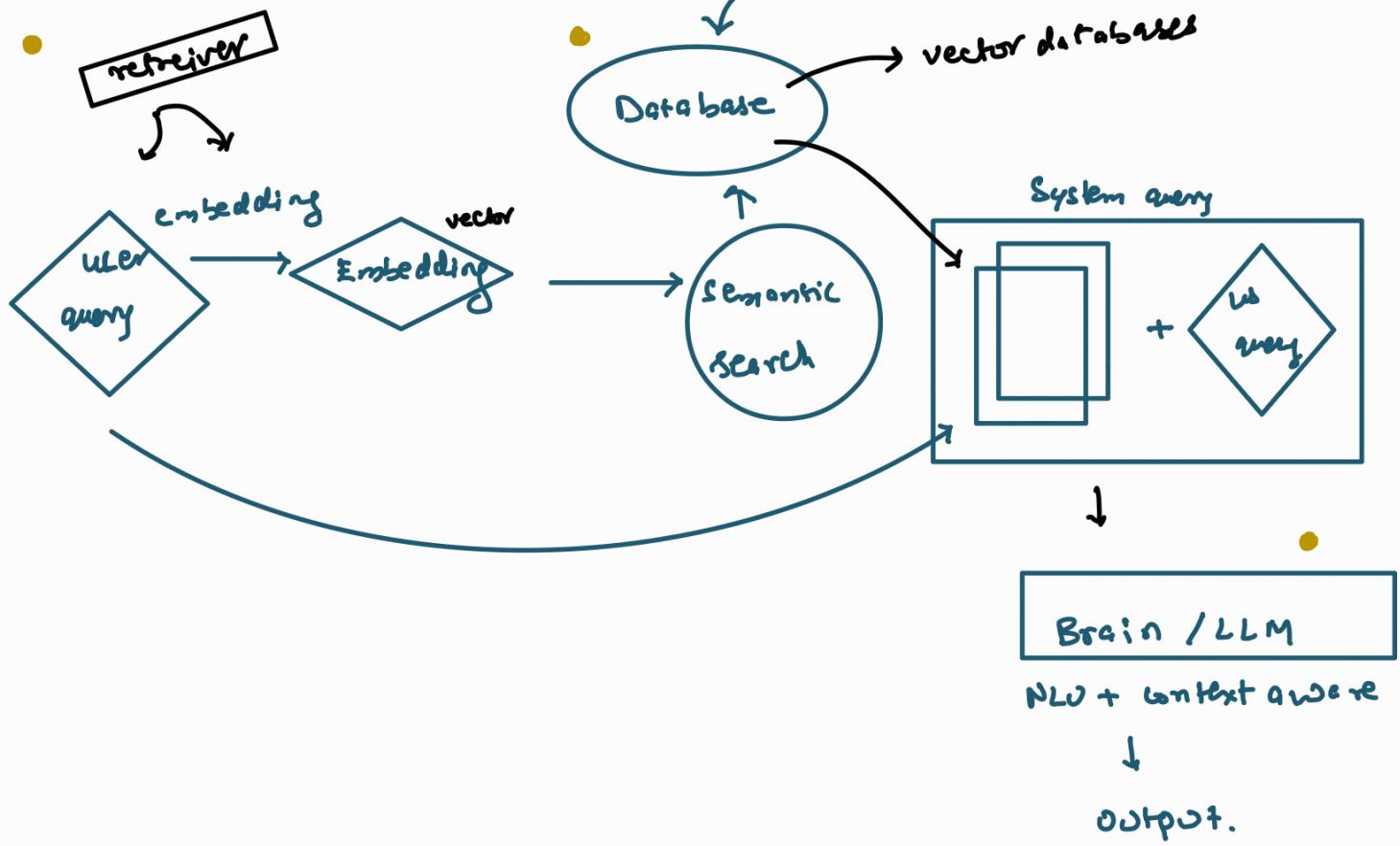


Semantic search → search based on meaning of text





↳ components



challenges:

i) building the brain. (NLU + text gen)

↳ use LLMs

ii) LLMs models are heavy (computation heavy)

↳ we can use the wrapper API

iii) Orchestration

↳ use LangChain

Benefits

- i) Concept of chains
- ii) model agnostic development
- iii) complex ecosystems (orchestration)
- iv) memory & state handling

What can be build?

Conversational chatbot

AI knowledge Assistants

AI agents

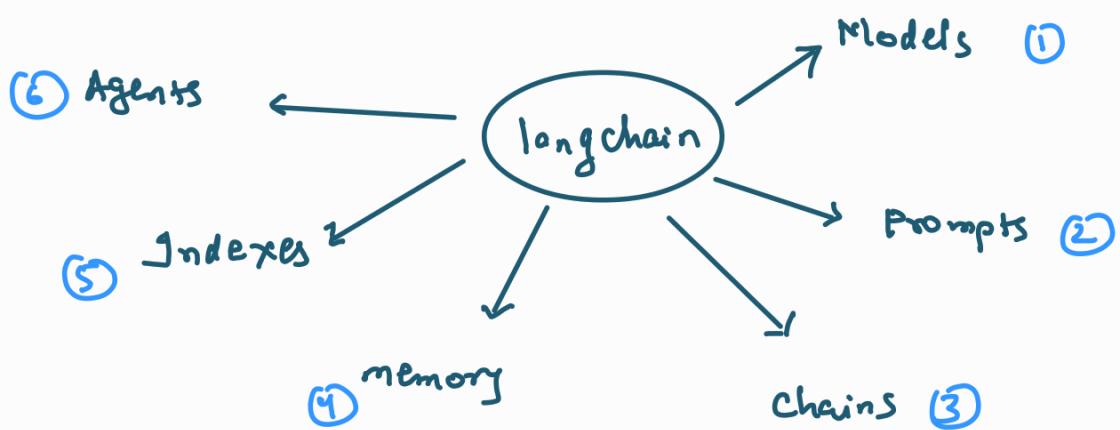
Workflow Automation

Sommernization / research helpers.

Alternatives

- i) Home index
- ii) haystack

LangChain Components:



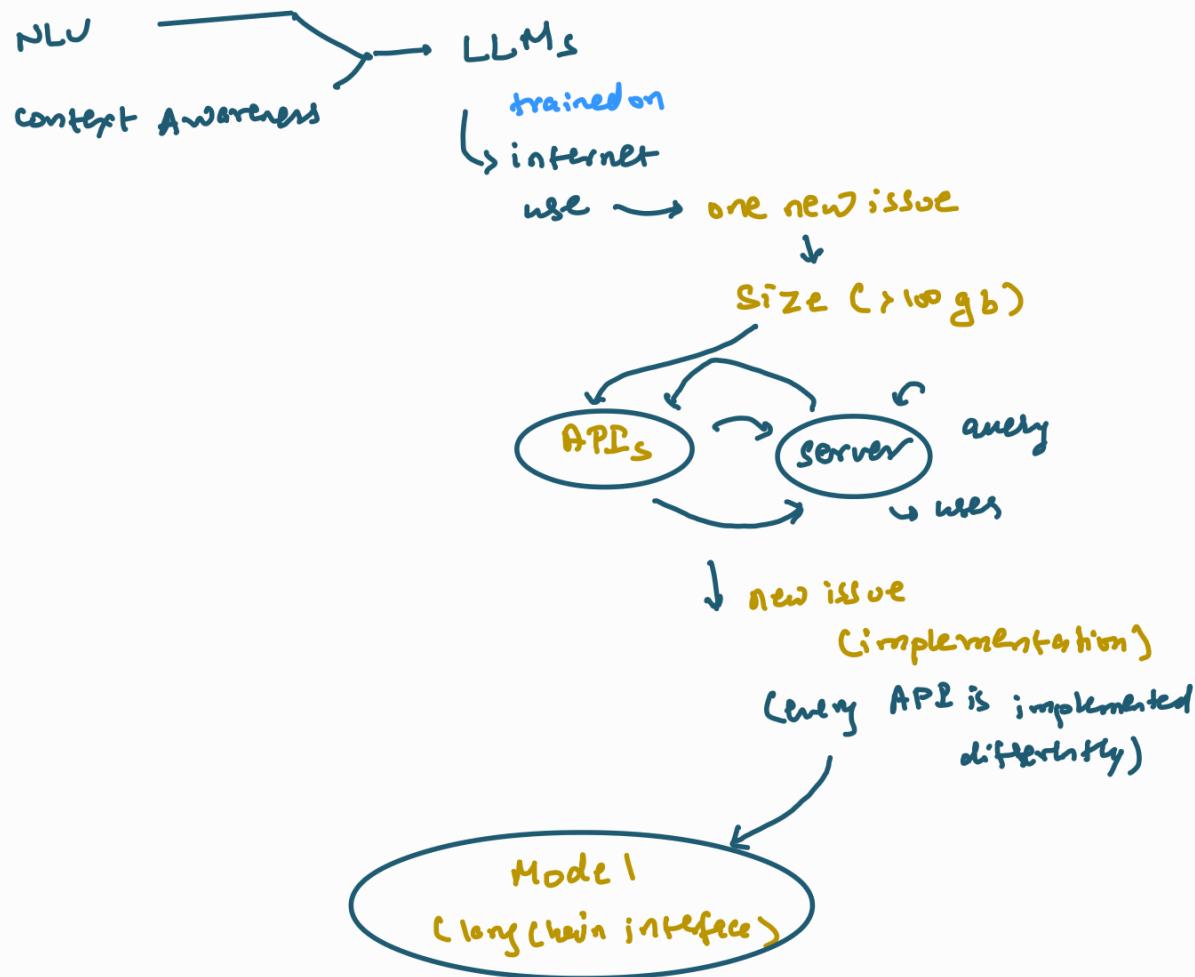
1) Models:

→ In Langchain, "models" are core interfaces through which we interact with AI models

problems before:

- i) making the NLP chatbot understand (NLU) for eg: make chat bot understand user query ("check my milk"),
- ii) if ever chatbot understand, it should provide relevant responses (i.e.) Context aware text generation

These problems:
Solved by LLM



ef

openAi

```

① from langchain-openai import ChatOpenAI
from dotenv import load_dotenv

load_dotenv()

model = ChatOpenAI(model='gpt-4', temperature=0)

result = model.invoke("hi who are you?")

print(result.content)

```

Anthropic AI

↑ (simple standardized model)

②

```

from langchain-anthropic import ChatAnthropic
from dotenv import load_dotenv

load_dotenv()

model = ChatAnthropic(model='claude-3', temperature=0)

result = model.invoke("hi who are you?")

print(result.content)

```

Langchain

Language model

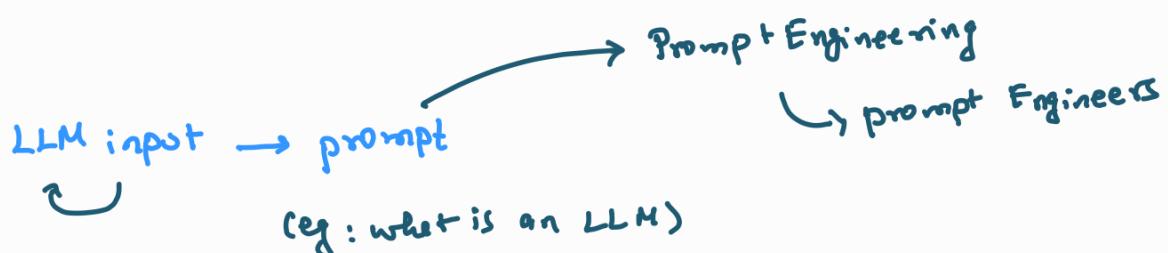
Embedding models

inp text → LMS
 LMS → o/p text

text → vector
 (i/p) (o/p)

↳ Semantic search

2) Prompts:



i) tone (academic, fun)

ii) Dynamic & Reusable Prompts.

```
from langchain_core.prompts import PromptTemplate
prompt = PromptTemplate.from_template ('summarize {topic} in {emotion}'
                                         ' tone')
print (prompt.format (topic='wicket', emotion='fun'))
```

iii) Role-based Prompt

```
chatPrompt = ChatPromptTemplate . from_template ["
    ("system", "Hi you are an experienced {profession}"),
    ("user", "tell me about {topic}"),
]
```

format the prompt with variables

```
formatted_message = chatPrompt.format_message (profession="doctor",
                                               topic="viral fever")
```

iv) few-shot prompt (based on examples)

① examples = [

```
{ "input": "I was charged twice", "output": "Billing issue" },
{ "input": "The app is crashing", "output": "Technical issue" }
```

]

② step 2 create a template

example-template = """

Ticket: {input}

Category: {output}

"""

③ step 3 : build the few shot prompt template

```
few-shot-prompt = FewShotPromptTemplate(  
    examples=examples,  
    example_prompt=PromptTemplate(  
        input_variables=["input", "output"],  
        template=example-template), prefix="Classify the following  
        customer support tickets into one of the following categories  
        'Billing issue', 'Technical issue' or 'General Inquiry'.\n\n",  
        suffix="\nTicket: {user-input}\nCategory:",  
        input_variables=["user-input"]  
)
```

3) Chains

Chains are components using which we can create a pipeline.

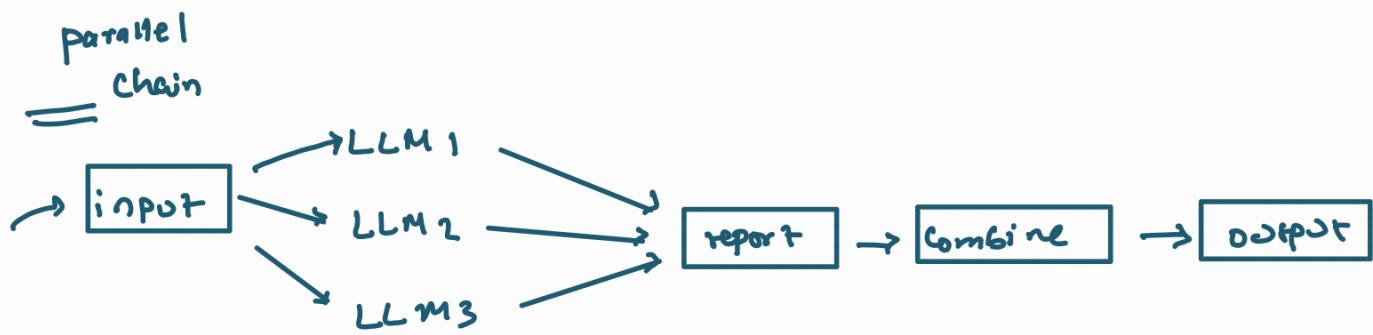
eg



sequential chain

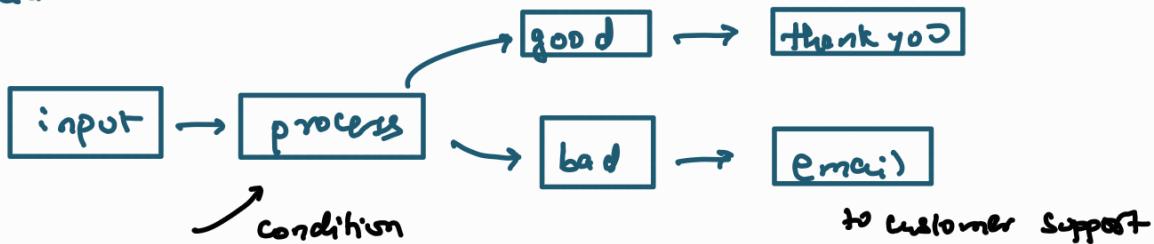
chains → output of one can be given as input to another.

Complex pipelines



conditional chaining

AI agent feedback



4) Indexes

Indexes connect your application to external knowledge such as pdf, websites or databases

Components of Indexes (4)

- i) document loader
- ii) Text splitter
- iii) Vector store
- iv) Retrievers

flow?

Load document → Embed document → store in vector store



chain for further ← search & retrieve ← create a retriever
processing

why indexes?

it cannot answer to questions that require access to private data

(we will need to provide it, in the form of pdf, doc (rag) kind of)

①

chatgpt



leave policy in xyz → cannot answer

②

chatgpt

← pdf, doc (external knowledge)



leave policy in xyz



Output

③ Memory : for storage

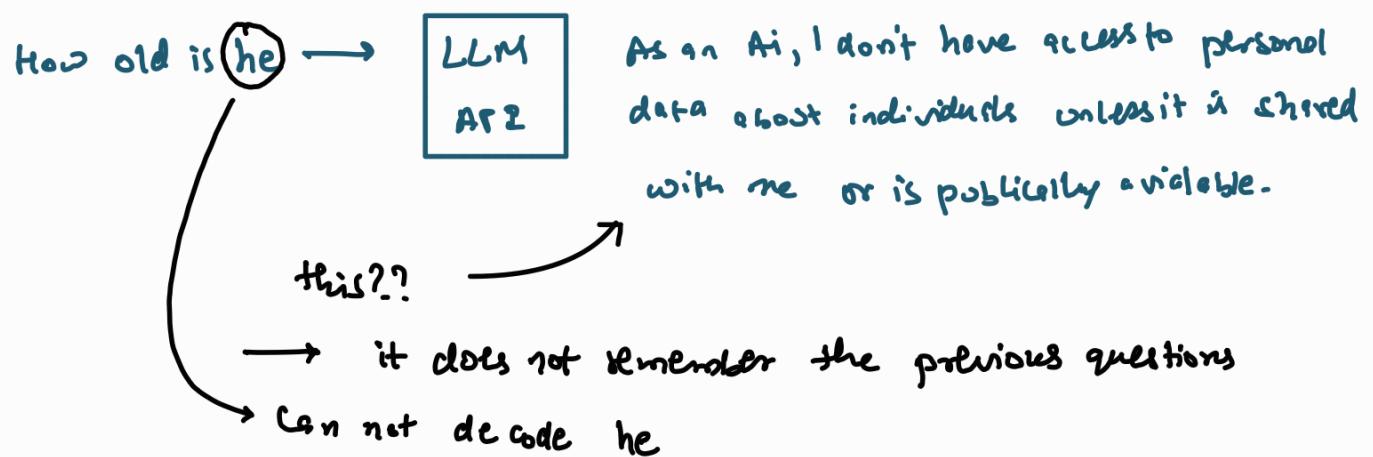
④ LLM API calls are stateless

who is
Narendra Modi?

→ **LLM API**

→ Narendra Modi is an Indian Politician
serving as the 14th and current PM of India

Follow Up



Issues

i) frustrating chatbot experience (since we need to remind it all the time)

ii) to fix this memory issue, we use this component.

Types of memory:

i) Conversation buffer memory:

→ Stores a transcript of recent messages, great for short

chats but can grow large quickly

chat history is stored and it sends the history of the chat with next API call

ii) Conversation Buffer Window Memory:

→ Stores last (n) interactions

→ Sends only those (n) interactions instead of entire history

iii) Summarizer-Based memory:

→ Periodically summarizes older chat segments to keep a condensed memory footprint

iv) Custom memory:

→ for advanced use cases, we can store specialized state (e.g. the user's performance or

key facts about them) in a custom memory class.

6) AI Agents:

Components designed to perform task autonomously by interacting with external tools or APIs.

→ LLMs (NLU + context-aware text generation)

↳ chatbot

eg:

↳ travel website

↳ AI agent (chatbot + agent)

user:

find cheapest tickets to Shimla on 25 Jan
(API call)

user:

book the tickets
(API call)

AI agent

↳ Reasoning Capabilities

↳ Access to tools.

tools: calculator, weather API

user: can you multiply today's temperature with 3?

agent: chain of thoughts

i) to do multiply today's temperature by 3

ii) what is today's temperature → weather API → 25°C

iii) multiply it by 3 → calculator $25 \times 3 = 75$