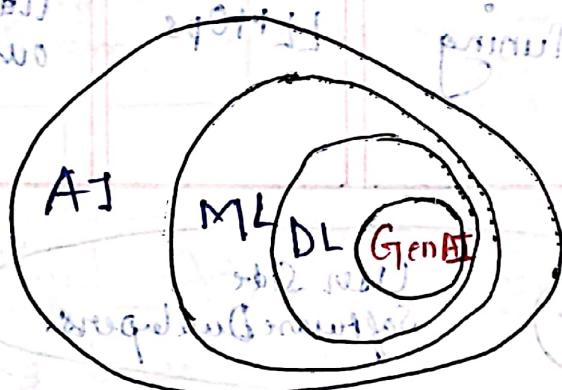


Generative AI

What is ?

Generative AI is a type of AI that creates new content - such as text, images, music or code by learning patterns from existing data, mimicking human creativity!



13th August 2025 *
GenAI ! Start from

1. LangChain ✓

2. Langgraph ✓

3. LangSmith ✓

4. fast-API. ✓

5. Docker vs Kubernetes

6. AWS for Deployment

7. Data Science ✓

8. ML Concept & Algo ✓

9. NLP

10. DL (tensorflow, PyTorch)

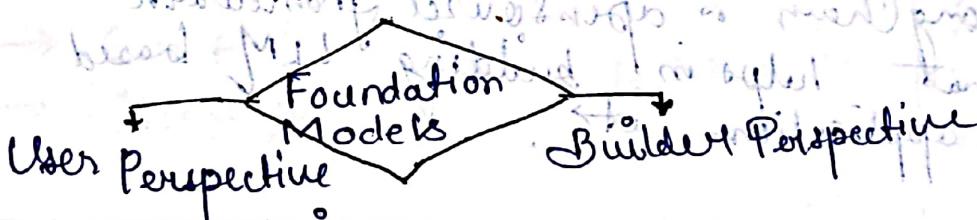
11. Databases ✓

12. Reinforcement Learning

13. Artificial neural n/w

Working area of genAI -

Customer Support, Content Creation, Education, S/W Development



The Builder's Perspective

Transfer Architecture	Pretraining	Fine Tuning
Types of Transformers	Training Objectives Tokenization Strategies Training Strategies Handling Challenges	Task-specific (RLHF) Instruction PEFT Continual Pretraining
Encoder (BERT)	Training Optn. Model compression Optimizing Inference	Evaluation
Decoder (BERT)		Deployment
Batch-based		
LTS (T5)		

Building Basic LLM Apps

OpenSource vs Closed Source LLMs

Using LLM API's

LongChain

Huggingface

Ollama

Prompt

Engineering

for app in LLM

new extract with FA

also no closure, no args

most, rewriting prompt jd

- partition, jobs

! function

fineTuning

Sci Tools

AGENTS

for app in LLM

new extract with FA

also no closure, no args

most, rewriting prompt jd

- partition, jobs

! function

LLMops

partition

Miscellaneous

Builder Side
Data Scientist

AI Engg.

User Side
Software Developers.

LANGCHAIN

→ LongChain is open-source framework that helps in building LLM based application etc.

→ Fundamentals

→ RAG

→ Agents

of NLU - Natural Language Understanding

→ Semantics Search - meaning based searching

By embedding (vector) conversion of the input

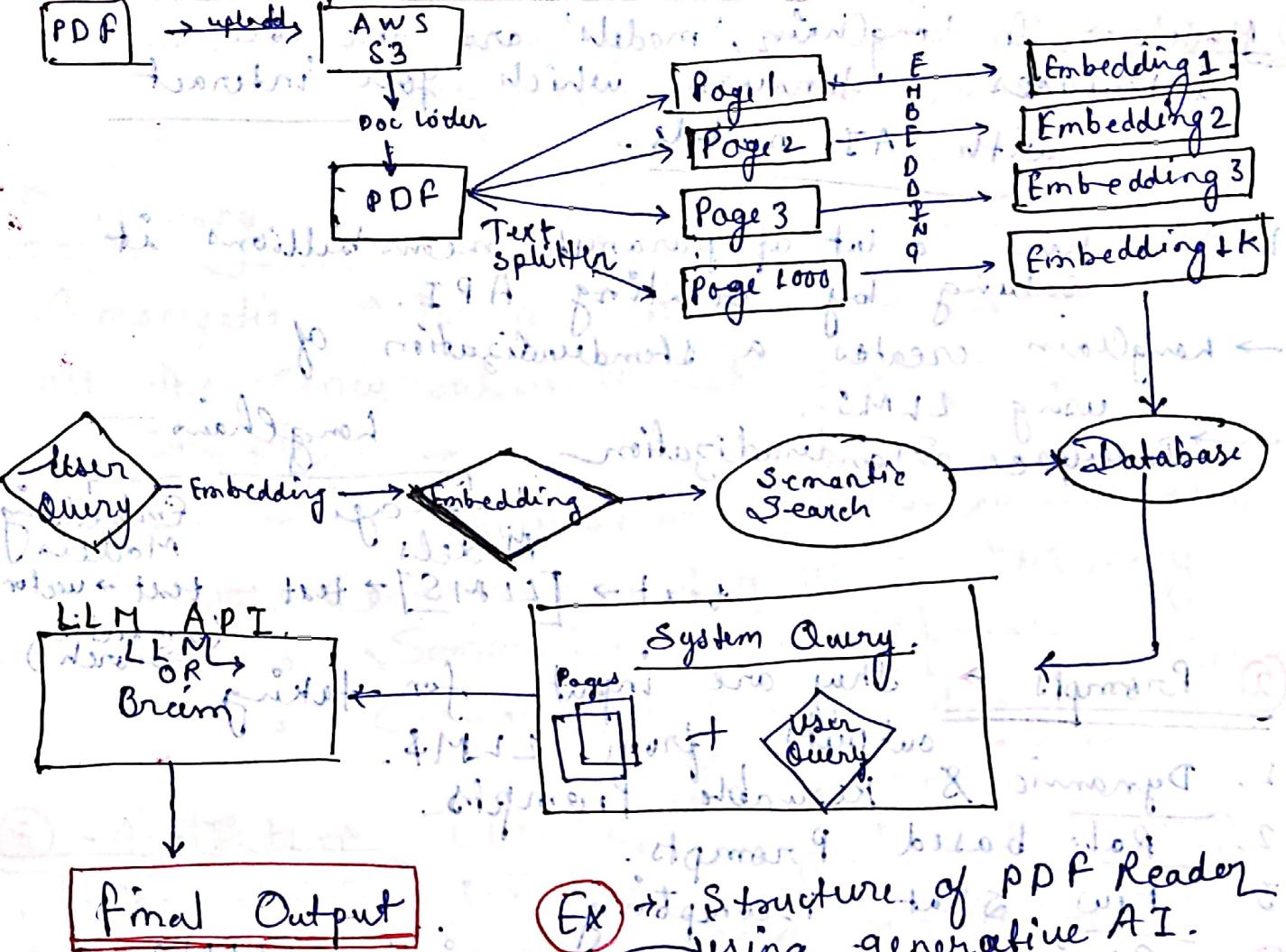
Alternatives → Haystacks, LlamaIndex

→ Context aware text generation

means output/step by step of the

model- or system (Ex- Chatbot)





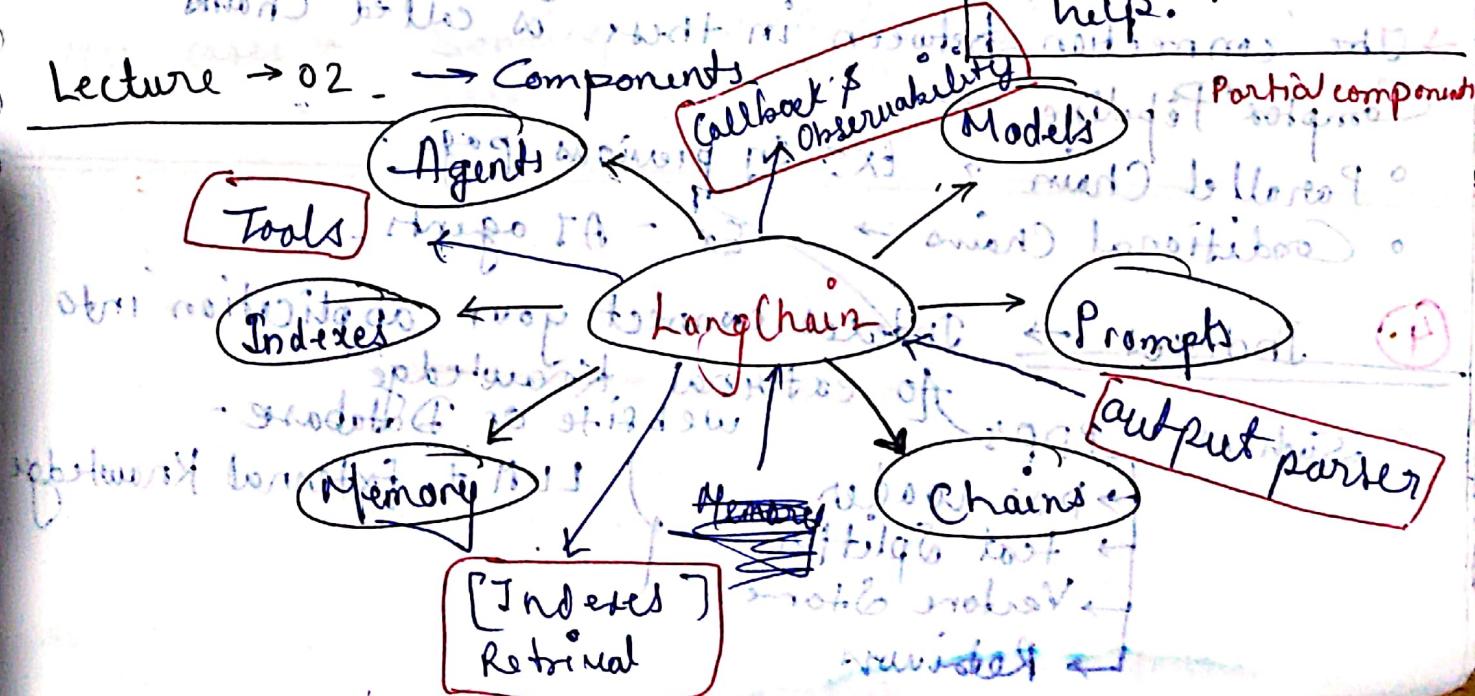
Final Output

Ex

Structure of PDF Reader using generative AI.

- What we can build
- Conversational Chatbots
- AI Knowledge Assistants
- AI Agents
- Workflow Automation
- Summarize / Research help

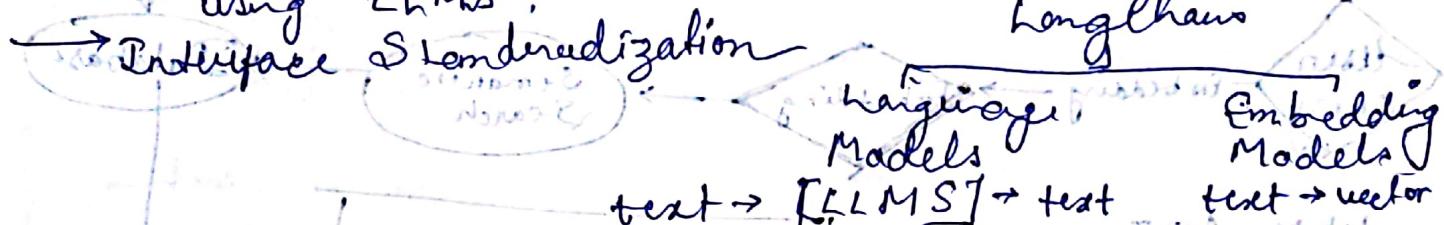
Lecture → 02 → Components



① Models → In LangChain, 'models' are the core interfaces through which you interact with AI models.

LLMS → have a lot of parameters (means billions) it scales by creating API.

→ LangChain creates a standardization of using LLMs.



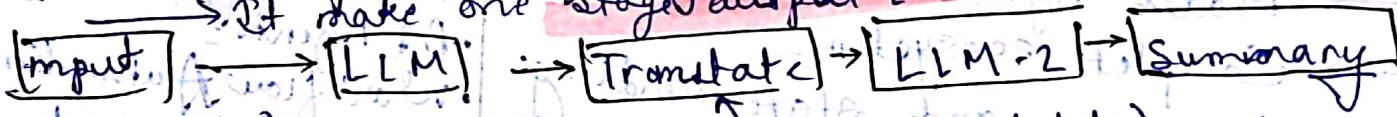
② Prompts → They are inputs for taking outputs from LLMs.

1. Dynamic & Resuable Prompts

2. Rate based Prompts.

3. Few shot Prompting

③ Chains → For pipeline building LLMs. Pipeline build is some parts by using chains concept. It makes one stage output to another as the input.



→ This is pipeline parts or stages (state).

→ The connection between these is called chain.

Complex Pipeline →

- Parallel Chain → Ex: In previous page.
- Conditional Chains → Ex - AI agents

④ Indexes → Indexes connect your application into external knowledge website or Database.

Such as - PPDF's

- Doc Loader
- Text Splitter
- Vector Store
- Retriever

LLM & External Knowledge

Indexes → Is the way of building LLMs which has external knowledge source, using which you can build.

⑤ Memoization → LLM's API calls are stateless.

Memoization → It's a technique to memorize all the conversation or task from start to stopping!

Memories → ° Conversation Buffer Memory → Store history

list n interactions ° Conversation Buffer Window Memory

° Summarizes based - Memory → Store Summary of Chat

° Custom Memory - Specialized or Specific

team, storing things that are crucial!

⑥ Agents → LLMs → NLU → evaluated form text gen. of chatbots.

Chatbots + Done some task
It's a AI that can do some real life tasks for you like your assistant.

→ Reasoning Capabilities
→ They use tools to do task
Tools → It can use API also for task

Why this works on concept → Chain of Thought

use cases → Sending Mail Automatically
also reply them

With no writing code
but & just work to
get

Usage of SA Impact
of writing code
of reading code
of testing code

S-O went without 20+ interplay

new if most required now if there will be
break in middle of work then S-O begin

it will find below blocks so writing to snippets
which will be written in blocks



The model component in longchain is critical part of framework designed to facilitate interactions with various languages models and embedding.

Models

Language Model	Embedding Model
→ return text	→ return vectors [num]
LMs Chatbot	Retrieval augmented gen ^m (RAG)

They are designed to process, generate, and understand natural language text. (NLP)

face from text generation:
Not built for memory
No understanding of rules
Ex - GPT-3, Llama 2-27B.

Optimized for multi-conversation
Supports structured conversation
GPT-4, GPT-3.5 Turbo, Claude

Chatbots

opensource / closedsource

Huggingface is the largest repo of opensource LMs

- opensource.
Using HF Interface API

Running locally.

14-08-2025

Done practical all kind of model type I had to go.

temperature → Is the value from 0-2

Is like that model given output from given input 0 or 1, if 0 means some output is going different or creative as increase value but same or deterministic as value decrease.

Prompts

are the instructions or queries given to a model to guide its output.

↳ Prompts

Text

Multimodal prompts ex - Image, video, audio, file.

but now it is used mostly.

→ Static → all control in your hands

→ Dynamic → Prepare a template so prompt accordingly work

Prompt template → In Langchain is a structured way to create prompts dynamically by inserting values into a predefined template.

Instead of hardcoding prompts.

Why → 1. Default validation } Why its use

2. Reusable; } instead of writing lots of code

3. Langchain Ecosystem

messages → Through an Example.

System → You are a experienced coder along with 25 years.

Human → Build UI and logic for langchain chatbot.

AI → Yes, ofcourse there is professional langchain chatbot by experienced coder with superb UI and a good logic.

→ Chat Prompt templates.

For maintain history

invoke

single message

static threads

Dynamic

list of msg. histories

Static → system

Dynamic → Human

→ AI

Message Placeholder →

It is in langchain so a special type of placeholder used inside a Chat prompt template insert chat history or a list of messages at runtime.

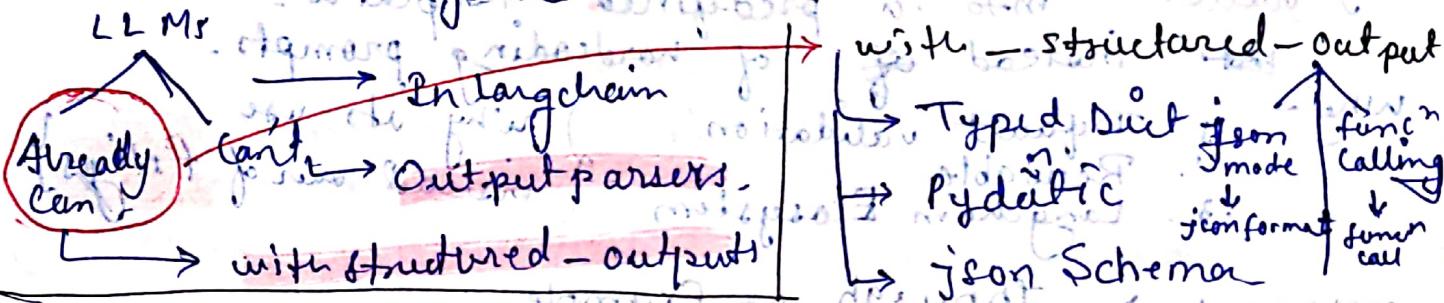
Lecture - 05

Structure Output in LangChain

It refers to having the practice of language models return responses in a well-defined data format. (Ex- JSON) rather free-form text.

This makes the model's O/P easier to work with accordingly.

Why? → Data extraction
API buildings.
Agents → LLMs & Chatbots directly to API & Database



① **Typed Dict** → It's easy to define dict in python while you specify what keys and values should exist, helps to ensure that your dict follows a structure.

- [Dictionary in Structure called Typed Dict].

→ They define value that validate it.

- Annotated → for explaining the LLM about the value returned to adapt better output.
literal → you use literal for give an option specific allowed values. If LLM is that we want within optional for choosing O/P or taken for granted.

② **Pydantic** → It is a data validation

and data parsing library for python.

It ensure that the data you work with

is correct, structured and type-safe

- > It requires valid data type for we define in class or anywhere
- > It can set default value
- > It also have features optional fields
- > It also try to understand data by their own and done implicit data conversion as per own (weise)
- > Built-in Validations → Ex - Email Str.
- > Field functions → type, orange (constraint), description this reduce this hallucinations regex expression because the LLM sees the schema clearly.

o json-schema → Universal data format

When project have multiple language and variety of data

typed but Pydantic

Only type hints.	Need Data Validation	Don't need to import extra python libraries
Don't need validation	Default value of LLM misses field, needed	validation but don't need python object
You trust the LLM to return correct data.	Automatic type conversion between models	Defines structure in standard JSON Data format.

Lecture 06

Structured Output in Langchain 2

15-08

[Builtin]

[Built-in]

✓ { with-structured } outputs { can't structure output-parsers by own }

It also work can't and can structure output
Output Parser → It helps to convert LLM response into structure format like JSON, CSV, etc., Pydantic models and more.

They ensure consistency, validation and ease of use in application.

Output Parsers

- String Output Parser
- JSON O/P Parsers
- Structured O/P Parsers
- Pydantic O/P Parsers

→ String O/P Parser → The String Output Parser is the simplest O/P parsers in LangChain. It is used to parse the O/P of LLM and return it as a plain string.

→ JSON O/P Parser → It is a quickest way to get json format data from LLM & Chatbots.

→ It doesn't enforce schema. We can't define format in that.

→ Structured Output Parser → It solves json problem that helps structured JSON data LLM response based on predefined field schema.

→ It works by defining a list of fields then model should return a structure output by ensuring a follow structure.

→ There is no data validation just defining the data structure.

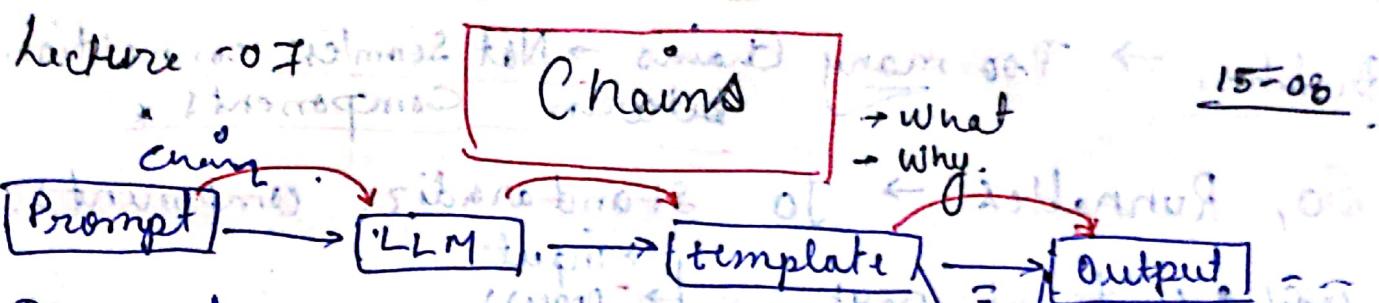
→ Pydantic O/P Parser → It's a structured O/P parser in LangChain that uses pydantic models.

→ Why? → Strict Schema Enforcement

→ Type Safety.

→ Easy Validation.

→ Seam less Integration



It creates a pipelines for the process to work one by one as we defined in a chain.

chain = template + model / part

thus we directly invoke chain

result = chain.invoke(x,y)

This type chain or pipeline is linear or sequential.

Simple Chain ✓

Sequential Chain. ✓

Parallel Chain ✓

Conditional Chain ✓

Types of Chain

Lecture - 08

Runables In LangChain

In simple

17-08

term runnable

means that functions

they invoke func to

execute them [generally]

- ① LLM Chain
 - ② Sequential Chain
 - ③ Simple Sequential
 - ④ Conversational Retrieval
 - ⑤ Retrieval QA Chain
 - ⑥ Router Chain
 - ⑦ MultiPrompt Chain
 - ⑧ Hydr Chain (Hypothetical Document Embeddings)
- ⑨ Agent Executor Chain
 - ⑩ SQL Database Chain
- transfer wider → revisited framework

Problem → Too many chains → Not seamless connection between components *

So, Runnables → To standardize components;

- Unit of work → input process output
- Common interface → invoke() batch() stream()
- Connection between them

→ a couple of runnable or group of runnable & also a runnable or chain is made.

Lecture - 09

Runnables 2:



→ These are core building blocks for structuring execution logic in AI workflows so they can be used in pipelines.

Purpose - Perform task specific execution by defining how different runnables interact (sequential, conditional, parallel)

Ex - ChatOpenAI

Prompt Template.

Retriever - retrieve relevant documents.

Ex - Runnable Sequence

Runnable Parallel

Runnable Map

Runnable Branch

Runnable Lambda

Runnable PassThrough

• read file

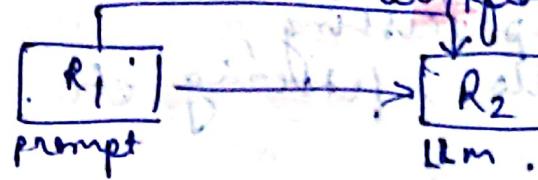
• main function

• strategy pattern

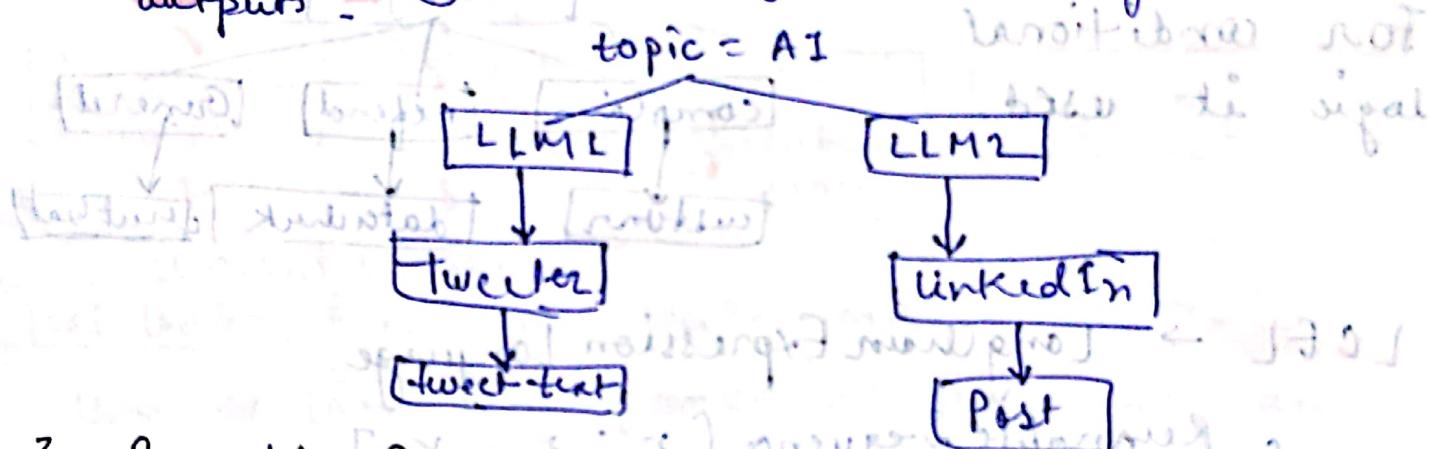
• strategy pattern



1. Runnable Sequence →
→ It is a chain of runnable in LongChain
that execute each after another, passing
the output of one step as the input to
next → Its useful when you need to compose
multiple runnables together in a structured
wayflow



2. Runnable Parallel →
that allow multiple runnables to execute in parallel.
Each R. receive the same input and process it
independently, producing a dictionary of outputs



3. Runnable Passthrough →
It is a special type of —
that simply returns the input as output
without modifying it

Ex - pass = RunnablePassthrough(func)

result = pass.invoke("Hello")

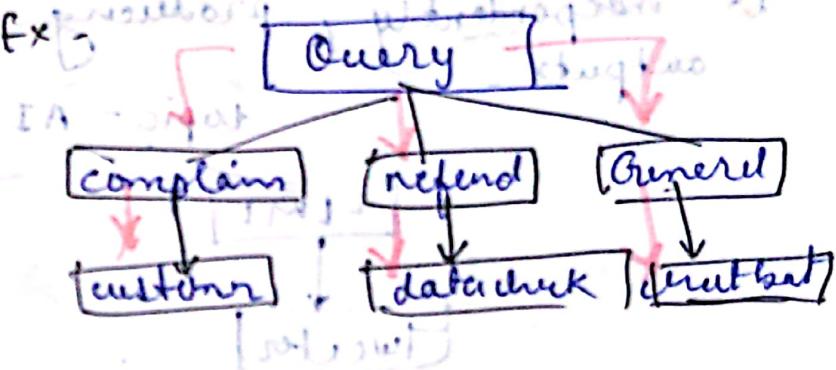
Output = Hello.

4. Runnable Lambda

can apply python custom functions with in an A2 pipeline.
It acts as middleware between different A2 comp. - enabling preprocessing, transformation, API calls, filtering etc.

5. Runnable Branch

It is control flow component that allows you to conditionally iterate input data to different or runnables based on custom logic. For conditional logic it used.



LCEL → Long Chain Expression Language.

Runnable Sequence $[r_1, r_2, r_3]$ is declared directly in LCEL as $[s_1, s_2, s_3, \dots, s_i, s_n]$.
 (s_1, s_2) significant elements = long → 3
 (s_1, s_2) return. step not = short
 $s_1, s_2, s_3, \dots, s_n$ = long

Lecture - 10

2nd Part:

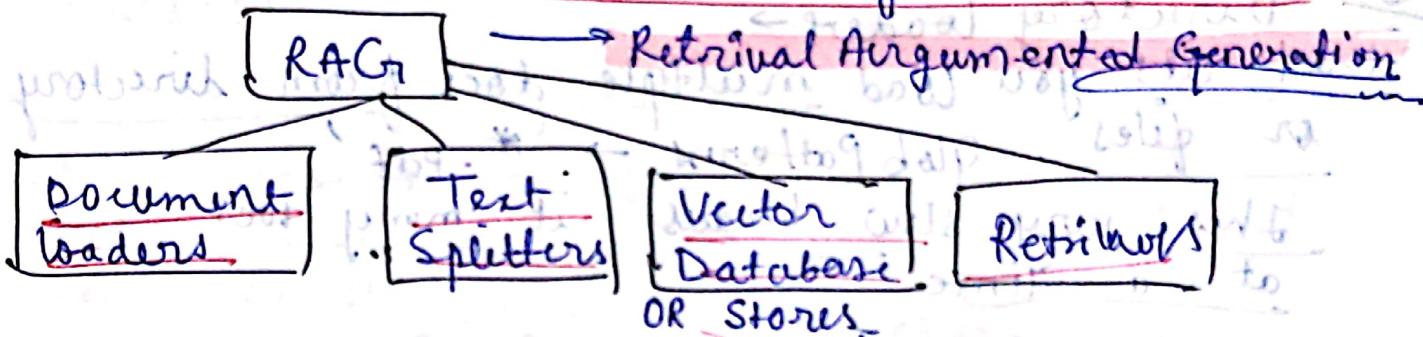
Document Loaders

RAG is a technique that combines info retrieval ~~that~~ with language generations where a model retrieves relevant docs from a knowledge base and then uses them as context to generate accurate and grounded responses.

1. Use of up-to-date info

2. Better privacy

No limit of document size.



Document Loaders →

Text Loader, PyPDF Loader, WebBase Loader, CSV Loader

Used to load data from various sources into a writer standardised format, which can be used for chunking, embedding, retrieval and generation.

Document

page-content = "the actual content of text"
metadata = { 'source': 'filename' }

1. Text Loader → chats, logs, transcripts, code snippets or any plain text data.
→ Works only with .txt files.

1 Text Loader →

18-08-23

2 Py PDF Loader → Simple PDF loader

It read PDF and convert it into document object for each page of PDF

There are issues in multiple scenarios so, there are loaders for multiple

Case → o PDF Plumber loader

Scanned image → o Unstructured PDF loader or Amazon TextractPDF loader

- o PyMuPDF loader ← Need image and layout
- o Unstructured PDF loader ← Want best structure in loader.

3 Directory Loader →

first lets you load multiple docs from directory or files. Glob patterns → *.pdf

This very slow because it many docs at a time.

→ Load vs Lazy load.

Faster loading

loads on demand

returns - a list of docs

A generator of docs obj

load all doc immediately in memory.

Docs. are not loaded at once; they're fetched one at a time as needed.

Best when:

The no of docs is less
you wont everting
upfront

Best when:

You're dealing with large docs and a lot

You want stream processing without lot of memory
(eg. chunks, embeddings)



4 Web based Loader → use to load and extract text and content from web pages (URLs)

It use Beautiful Soup under the hood to parse HTML and extract visible text.

Four blogs, news article
public website where
content is primarily

text

Doesn't JS pages

use Selenium WebDriver
loads only static content

5 CSV Loader → load CSV docs in code.

Lecture - 11.

L8 - 08

Text Splitters

It is process of breaking large chunk of text (files, PDF, articles, HTML pages or book) into smaller, manageable pieces that can handle effectively. {Text-Splitting}



Large-text



chunks

Task	help
Embedding	Short chunks yield more accurate vector
Semantic Search	Search result point to focused information
Summarisation	Prevents hallucination and topic drift

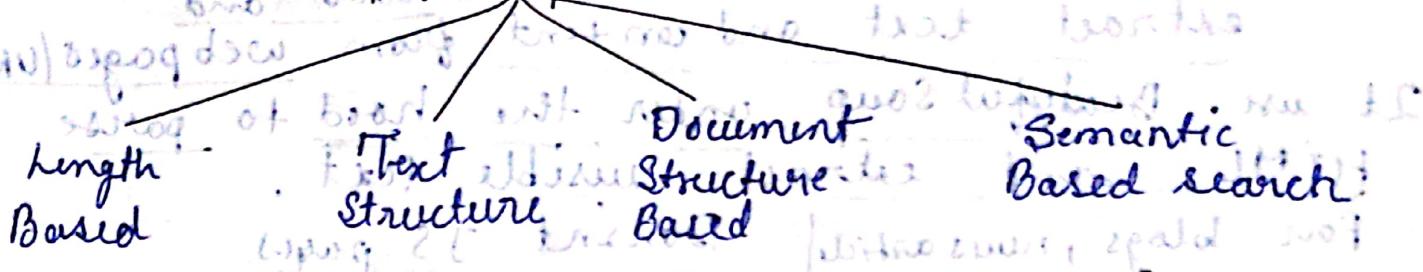
Models have limitations

So it allows us to process docs that would otherwise exceed these limits.

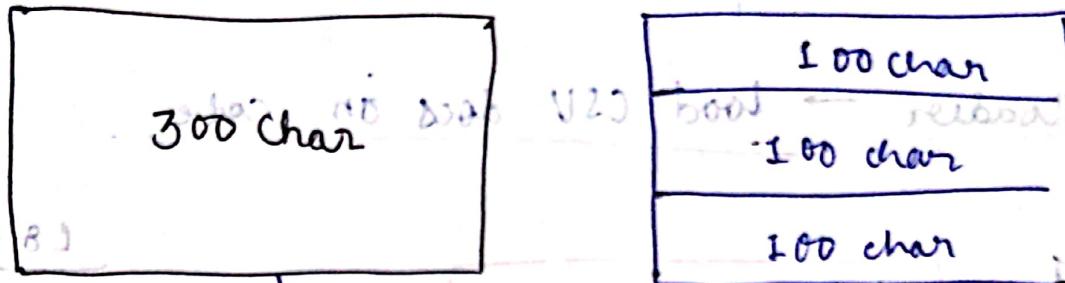
Optimising Computation resources →

Working with smaller chunks of text can be more memory efficient and allow for better parallelization of processing tasks.

Text Splitters



I. Length Based Text Splitting



→ It just split text not modify or delete/edit it.

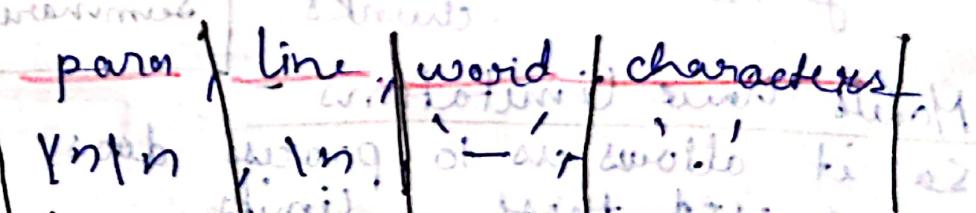
Chunk overlap → that text lie in both chunks so the text is overlap in both them.

Eg. Langchain is a framework this text

to order work to build 11m.
chunk - overlap = 4

2 Text Structure Based → It create chunks

based on - para | line | word | characters



eg -

[My name is. Musa (12)]
[I live in Bhopal (12)] . [n/n - para]
[I am 19 years old (13)]
[And I'm doing coding (15)]

Text processed for multithreading if chunk size = 10

My name is Musa (12) I live in bhopal (12) And I am doing coding (15)
I am 29 years old (13)

In line (n)

My name is Musa (22)

My name is Musa (4)

2 + 4 + 2 + 4 > 10

8 > 10

[My name is] [Musa]

6 14

10

In line (n)

I am 29 years old (13)

1 2

I live in bhopal (12)

1 2

And I am doing coding (15)

1 2

In words ('_') or spaces

or spaces

words

or spaces

20 bath under 10 So total

chunks have 8 atoms

If chunk size = 25 then chunks are 4 - on lines
for some eg.

3. Document Structured Based -
This for special kind of format not
for normal text like english or another
format.
It is for eg - code snippet, markdown

4. Semantic Meaning Based -

It try to understand whole text
then it chunk as per their
relation this is called meaning

based splitter { By using embedding }
for finding similarity

{ in vectors of the text }.

S₁ S₂ S₃ S₄ ... S_n
low -

Lecture - 12

vector stores

It is a system design to store and retrieve data represented as numerical vector.

Features

1. Storage memory
2. Similarity Search
3. Indexing
4. CRUD operations

use case -

1. Semantic Search
2. RAG
3. Recommender System
4. Image/Multimedia Search

vector stores

- Storage System
 - Semantic Search
 - It is library.
- Ex - FAISS.

Vector Database

- Distributed architecture
- Back up & restore
- ACID Transaction
- Autentication
- Concurrency control

Ex - Milvus, Doran, Weaviate

vector stores + Database feature

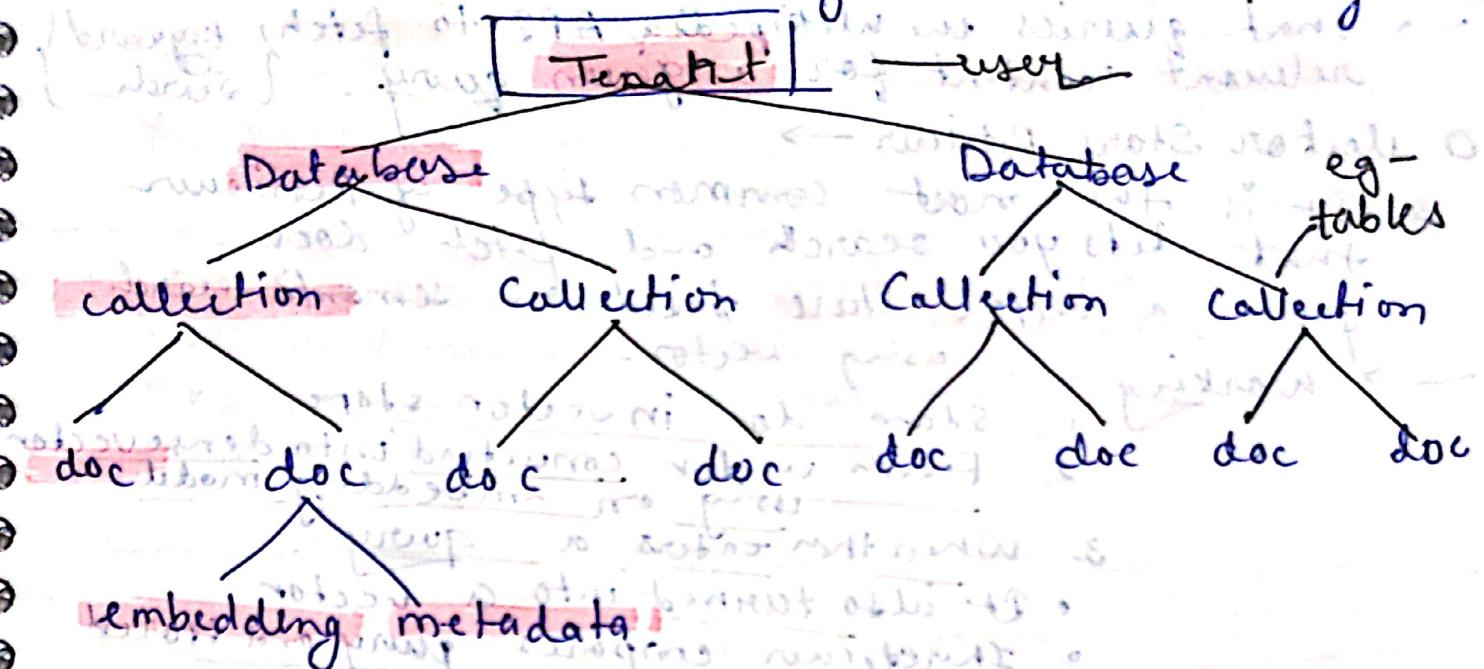
Vector Database

- Its use for clustering, scaling, security metadata, filtering and durability
- Vector Stores In Longchain -

Chroma Vector Store [Database but very fast query to store, i.e. lightweight]

- It is a open source, light weight database that is especially friendly for local development and small to medium scale production needs

→ Chroma Tend only . and DB Hierarchy .



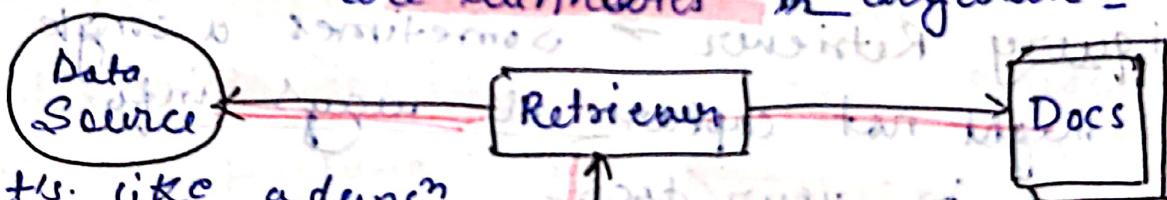
Lecture - 13

Retriever

A retriever is a component in Langchain that fetches relevant docs from a data source in response to user's query.

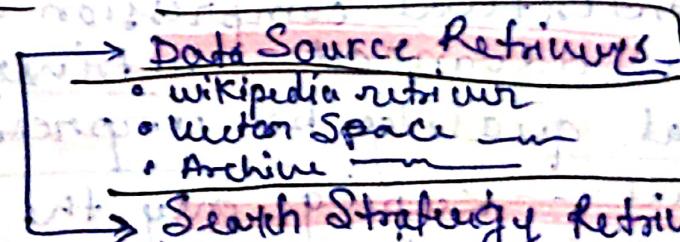
There are multiple type of retrievers. →

All retrievers are cummable in langchain -



It's like a funcn that take user query and give docs object through search in data source for relevant results.

types of →



Search Strategy Retrievers

- MMR (Maximum Margin Relaxation)
- Multiquery
- Contextual Comfusion

o Wikipedia Retriever → that queries the wikipedia API to fetch relevant content for a given query. {Search}

o Vector Store Retriever →

It is the most common type of Retriever that lets you search and fetch docs from a vector store based on semantic search.

Working - using vector -

1. Store doc in vector store.
2. Fetch vector converted into dense vector using an embedding model.
3. When user enters a query:
 - It also turned into a vector
 - Retriever compares query and stored vector
 - It retrieves the top k most similar ones.

o Maximal Marginal Relevance →

It is an info. retrieval algo designed to reduce redundancy in the retrieved results while maintaining high relevance to the query.

Pick relevant query but also they are different each other.

o Multi-query Retriever → Sometimes a single query might not capture all ways info is phrased in your docs.

It is work on query ambiguity to over.

o Contextual Compression Retriever →

The contextual compression retriever in Langchain is an advanced retriever - that improves retrieval qualities by compressing docs after retrieval - keeping only the relevant content based on user's query.

Remove irrelevant content in docs.

- Your does ~~are~~ be long and contain mixed info.
- You want to ~~reduce~~ context length for LLMs
- You need to improve ans accuracy. in RAG pipelines.

B M25 Retriever

Self Query

Arxiv

Multivector

Parent Document

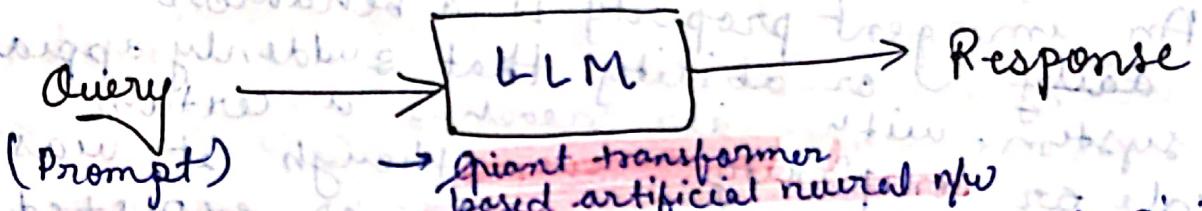
Time weighted vector

Esemble

Lecture - 14

RAG → Theory

18-08



They store data in Parameters. It is called Parametric Knowledge

Hallucination

problems in LLM's.

Private dat

Recent data

Hallucination

Solution → Fine Tuning ⇒ Eg - Student → LLM

Supervise finetuning → Train model by labeled data. engg → Pretrain training → finetuning

Continued/Unsupervised fine tuning → With out labelled data -

RHF → Reinforcement learning + human feedback

1. collect data - desire prompt - eg.

2. Choose a method - full parameter FT, LoRA / OLoRA or parameter-efficient adapter

3. Train for few epochs - exact method

4. Evaluate & Safety test - exact method for safety

fine tuning is not feasible -

→ Expensive to do

→ High technical reqn

→ Not work on high change frequency data.

So, there is alternate →

In context learning → It is a core capability of large language Models (LLMs) like GPT-3/4, Claude and LLaMa where the models learn to solve a task purely by seeing examples in the prompt — without updating its weights.

It is emergent property of LLM's.

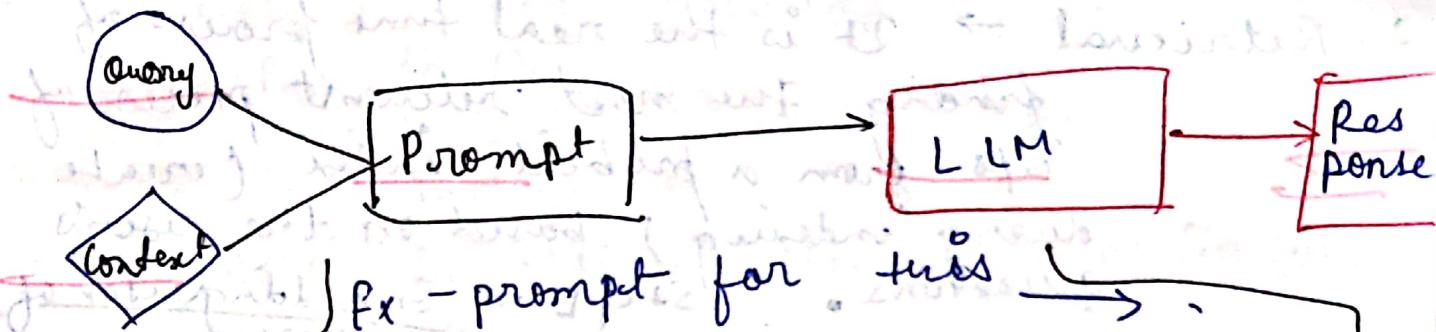
→ An emergent property is a behavior or ~~ability~~ or ability that suddenly appears in system - with it reaches a certain scale or complexity - even though it was not explicitly programmed or expected from the individual components.

→ Instead of just example, retrieve background information, facts, docs, products, manuals etc. inject that into the prompt to argument the model's knowledge

→ This is called RAG

RAG is a way to make a language model [like ChatGPT] smarter by giving it extra information at the time you ask your questions

Basically personalise, protected, lie in given limit and worked only on given environment this language model are known as RAG.



You are helpful { context }
Ans the que only from provided { context }.

If the context is insufficient, just say
don't know.

↓ context of question

Question : { question }

Understanding RAG → Info retrieval → text generation

- Indexing → format data → prepare or create.
- Retrieval → fetch data
- Argumentation → Add extra knowledge on parametric
- Generation → After read prompt → it text genⁿ to give o/p.

→ Indexing → It is the preparing your knowledge base so that it can be efficiently searched at query time.

This step consist 4 parts:

1. Document Ingestion → load your source knowledge

meaningful chunks

2. text chunking → Break large docs into small semantically meaningful chunks.

3. Embedding generation → Convert each chunk into a dense vector that captures its meaning

External Knowledge Base → 4. Storage in vector store → Store the vectors along with the original chunk text + metadata in vector database

2. Retrieval → It is the real time process of finding the most relevant pieces of info from a pre-built index (create during indexing) based on the user's questions.

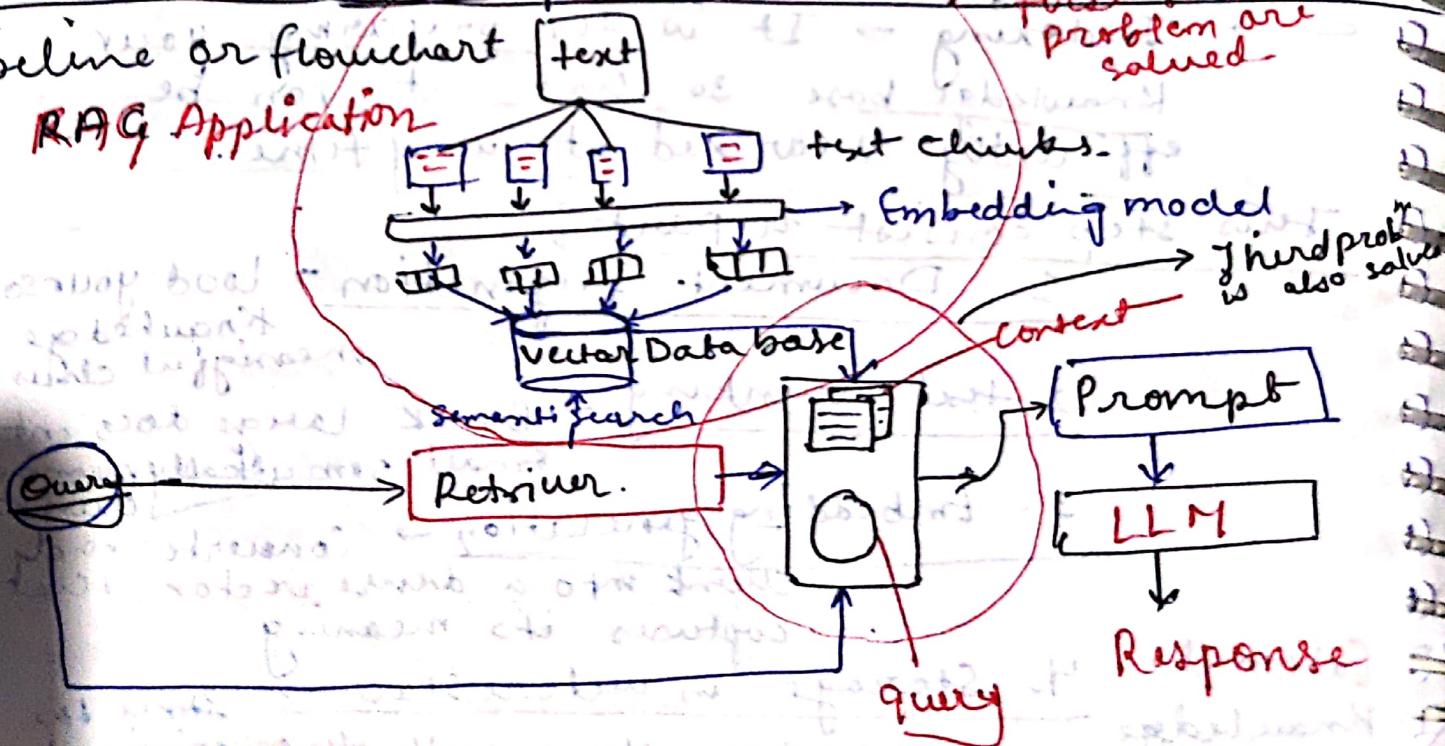
Step 1 → Embedding of query
Step 2 → Searching semantically
Step 3 → Ranking of cosine similarity
Step 4 → Most relevant chunk are fetched

3. Argumentation → refers to the step where the retrieved docs (chunk of relevant context) are combined with user's query to form a new, enriched prompt for LLM.

It creates a [prompt = context + query]

4. Generation → It is the final step where a LLM uses the user's query and retrieved & argumented context to generate a response.

Pipeline or flowchart of RAG Application



→ It clearly solves problems of LLM.

- 1. Private Data
- 2. Recent Data
- 3. Hallucination.

→ So RAG is optimal alternative of finetuning

- 1. Cheaper.
- 2. Easier to implement (not complex).
- 3. Don't require millions of preproc data.
- 4. Don't need train modal.

Lecture - 15

RAG

- Practical

19-08

Done

Improvements *

Metrics →

faithfulness

answer relevance

context - precision

context - recall.

① UI Based Enhancements:

② Evaluation:

- a. Ragas

③ Indexing:

- b. Langsmith

④ Retrieval →

- a. Pre-retrieval,
- b. During retrieval

⑤ Argumentation →

- a. Prompt templating

- b. Answer grounding

⑥ Generation →

- a. Answer with citation

⑦ System Design →

- a. Multimodal

- b. Agentic

- c. Memory Based



Lecture - 16

3rd Part

Tools In longchain

Create
or
fetch
Tools

19-08

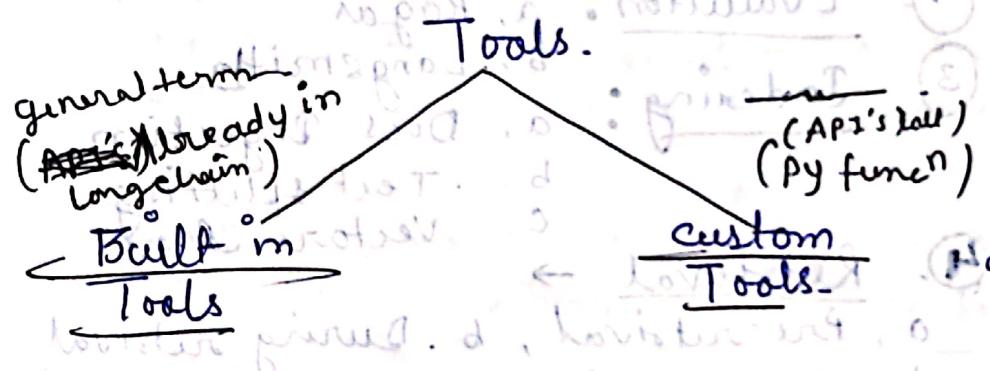


Tool → It is just python function (or API) that is packaged in a way the LLM can understand and call when needed.

LLMs are great at:

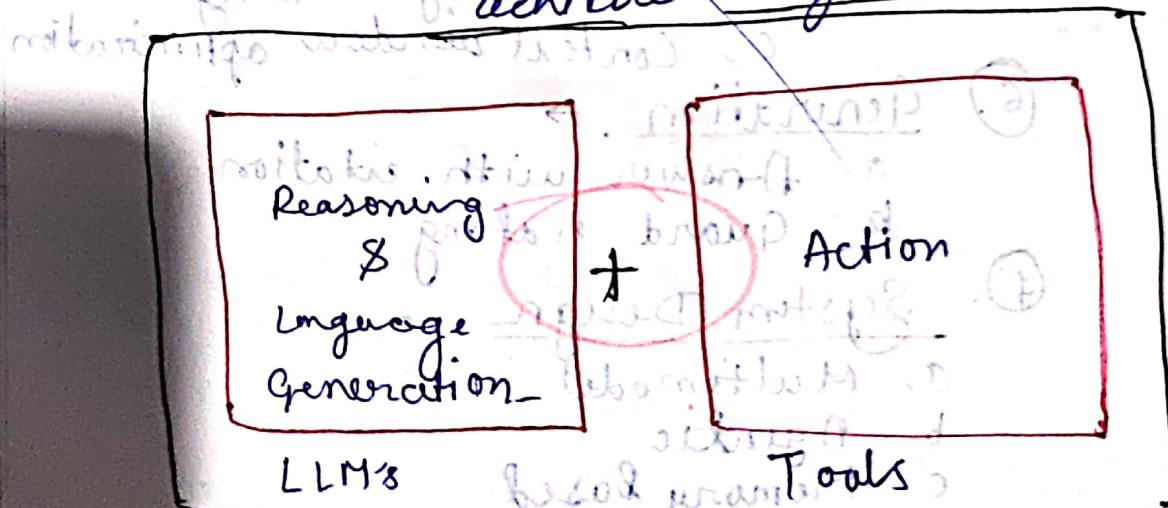
- Reasoning (Think)
- Language (text) generation (Speak)

- But they can't
 - o Access live data
 - o Do reliable math
 - o Call APIs
 - o Run code
 - o Interact with DB



They are also runnable because they have invoke function

An AI agent is an LLM-powered system that can autonomously think, decide and take actions using external tools or APIs to achieve a goal.



AI-AGENTS.

* Built-in Tools →

- A built tool is the already longchain provided you
- It's prebuilt →
 - DuckDuckGoSearchRun → web search
 - WikipediaQueryRun → summary
 - Python RFP Tool → run raw code
 - Shell Tool → run shell commands
 - RequestGetTool → Make HTTP Getreq
 - GmailSendMessageTool.
- send email via gmail →

* Custom Tools → user defined tool.

when :

- when you want to call your own API's
- when you want to encapsulate business logic.

You want the LLM to interact with your database, product or app.

1. Create a function
2. Add type hints
3. add tool decorator.

ways to create tools

using @ tool decorator

using Structured tool & pydantic

using BaseTool class

special type of tool where
the input to tool follow
a structured schema using
a pydantic

It defines the core structure and interface that any tool must follow, whether it's simple one-line or fully customized functions.

All other tool types like @ tool, StructuredTool are built top of: **BaseTool**

Tool kits → A toolkit is just a collection of (bundle of) related tools that serve a common purpose package together to form a convenience and reusability.

In LongChain:
→ A toolkit might be: Google Derive toolkit
→ And it can contain following tools:

→ Google Derive Create File Tool

Search File Tool

Read File Tool

Tools - 2
Calling Tool

Use or call

Bind

Tools

Tool Binding → Tool binding is the step where you register tool with a LLM so that:
1. The LLM knows what tools are available
2. It knows what each tool does (via description)
3. It knows what input format to use (via schema)

Tool Calling → It is the process where the LLM decides, during a conversation or task that it needs to use a specific tool (func) and generates structured with the name of the tool and the arguments to call it with.

→ The LLM doesn't actually run the tool. It suggests the tool and inputs arguments. The actual execution done by or handled by longchain or you.

Toolkits → A toolkit is just a collection of (bundle of) related tools that serve a common purpose package together to form a convenience and reusability.

In Longchain → A toolkit might be: Google Derive toolkit
And it can contain following tools:

→ Google Derive Create File Tool

Search File Tool

Read File Tool

Lectures - 17

Tools - 2
Calling Tool

Use or call

Bind

Tools

Tool Binding → Tool binding is the step where you register tool with a LLM so that 1. The LLM knows what tools are available 2. It knows what each tool does (via description) 3. It knows what input format to use (via schema)

Tool Calling → It is the process where the LLM decides during a conversation or task that it needs to use a specific tool (func) and generates structured with the name of the tool and the arguments to call it with.

→ The LLM doesn't actually run the tool. It suggests the tool and insert arguments. The actual execution done by or handled by longchain or you.

Tool Execution → It is the step where actual python function (tool) is run using the input arguments that the LLM suggest during tool calling.

Tool Message → It is special type of message it get when executing tool with the help of tool call.

Currency converter

Lectures - 18

AI Agents

AI Agent → It is an intelligent system that receive a high level goal from a user and autonomously plans, decides and executes a sequence of action by using external tools API's and Knowledge source.

All while mountain context reasoning over multiple steps adopting to new info and optimizing for the intended outcome.

LLMs + Tools → Agents

Characteristics →

Goal driven

Autonomous planning

Tool using

Context aware

Adaptive

LLM + Tools → Agent

The process is actually React -

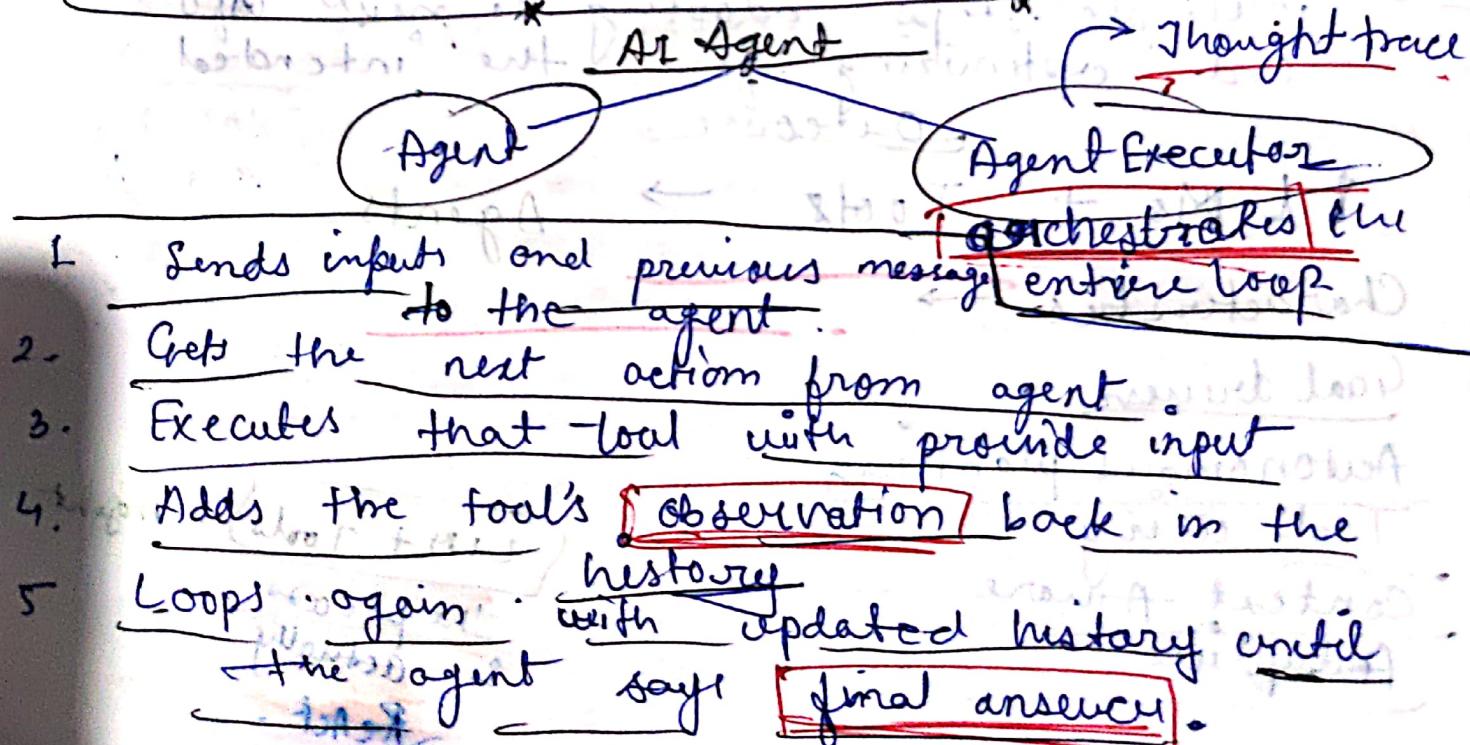
ReAct → It is a design pattern used in AI agents that stands for Reasoning + Acting. It allows a language model (LLM) to interleave internal reasoning (thought) with external actions (like tool use) in a structured multi-step process.

→ Instead of generating ans on go, the model thinks step by step, deciding what it needs to do next and optionally calling tools (APIs, calculator web search etc.) to help it.

useful for →

- Multistep problem
- Tool augmented tasks (web search, DB lookups)
- Making the agent's reasoning transparent and auditable.

Thought
↓
Action
↓
Observation



Creating an Agent

```
agent = create_react_agent(
```

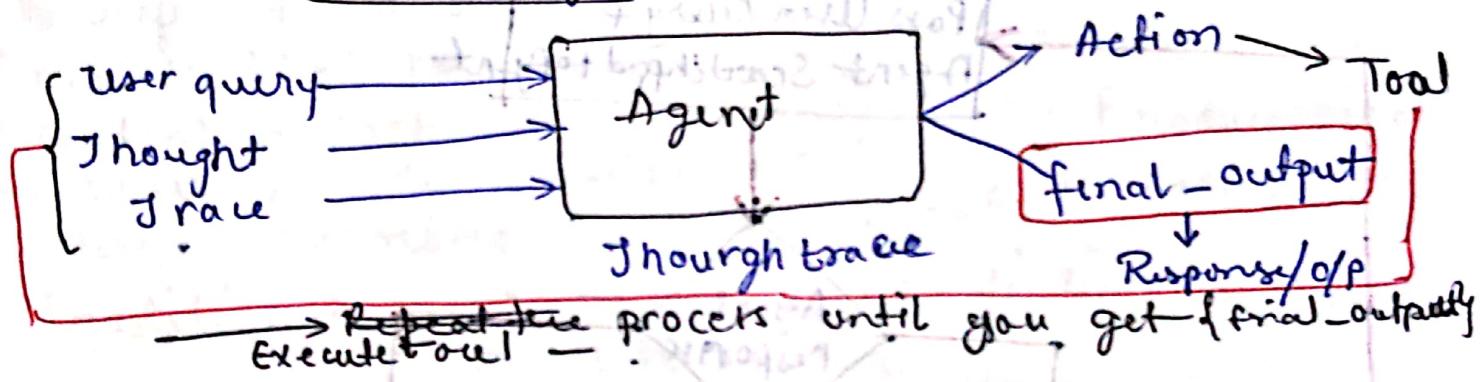
llm = llm, — Your Model - AI

tools = [search_tool], — Tools - in

prompt = prompt — Prefer that from

trigchain-hub use
built in prompt (like [
huchase7/rear])

Also create custom
prompts. but that is
not preferable



Creating an Agent Executor

```
agent_executor = AgentExecutor(
```

.. agent = agent, — Agent calling

tools = [search_tool], — Tool search

herbase = True — Internal process

is display.

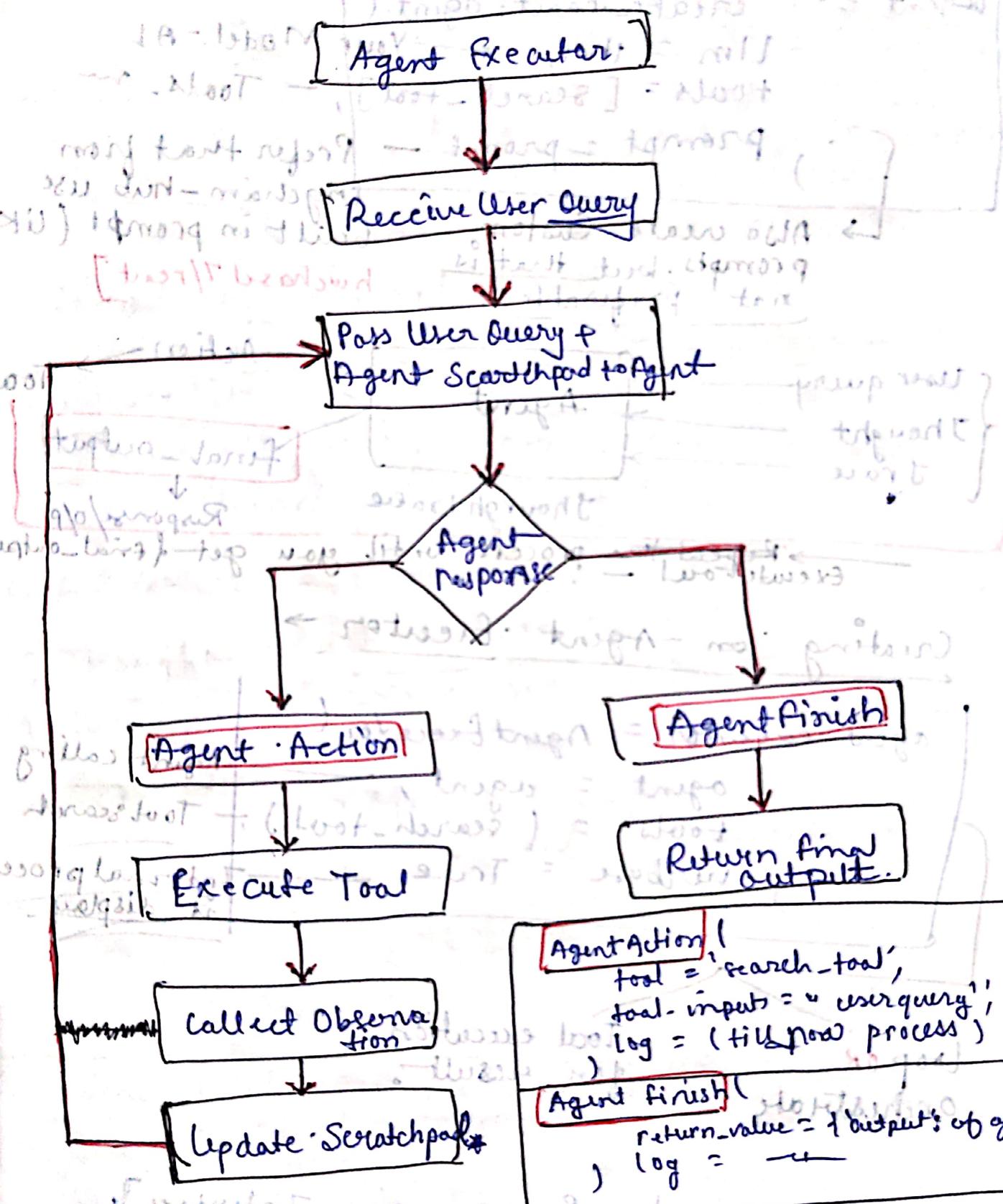
loop OR
Orchestrate

Tool execution
for result.

API { Application Programming Interface }

function ref training to predict sequential acts
as aggregated list see to better storage in
{actions, act, state, reward}. predict next

FLOW CHART



Agent Action
`tool = 'search-tool',
 tool-inputs = "userquery",
 log = (till now process)`

Agent finish
`return-value = {output's of agent},
 log = _____`

WORKING OF AI AGENT.

This longchain library is not efficient for making AI agents instead of use the longGraph or some other library. (longsmitch for monitoring)