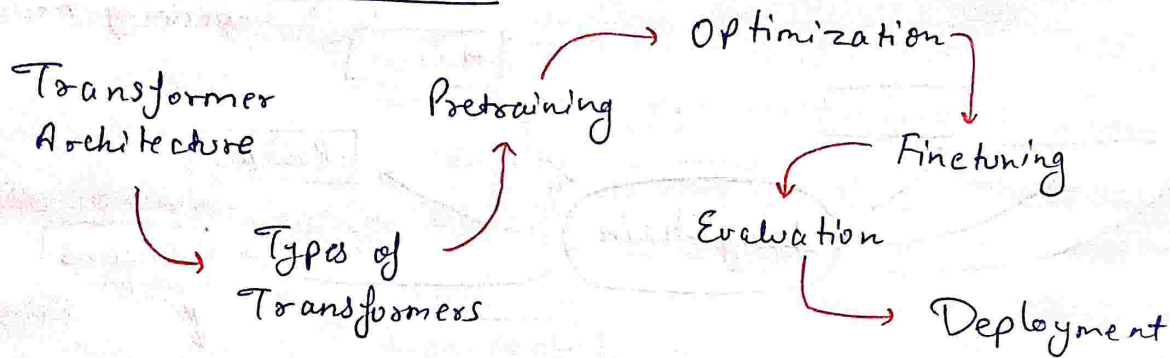


98-06-2026

GEN-AI

## Builder's Perspective



## User's Perspective

- Building Basic LLM Apps
- Prompt Engineering
- Agents
- Fine tuning
- RAG
- LLM Ops

## Lang Chain

- open src framework
- supports all major LLMs  
either open src models or close src
- simplifies developing
- supports many tools & all major GenAI use cases

## Benefits

- concepts of chains
- Model Agnostic Development
- Complete Ecosystem
- Memory & state handling

## Use Cases

- Conversational Chatbots
- AI Agents
- Workflow Automations

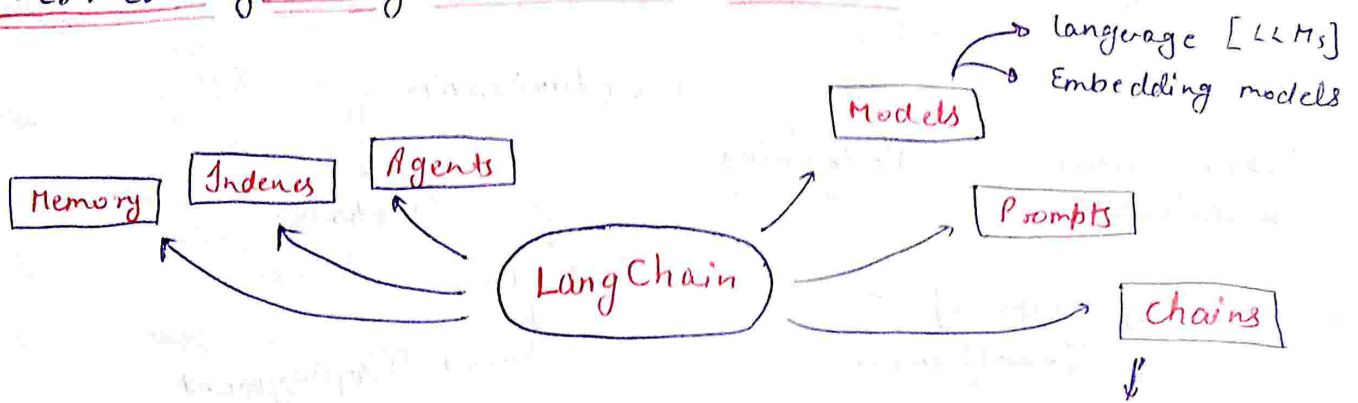
## Alternatives of Lang Chain

- LlamaIndex
- Haystack

RAG: Retrieval Augmented Generation

Ex → Custom Chatbot on a PDF

# Overview of LangChain Components



Output of a component will be input of another

## # Indexes

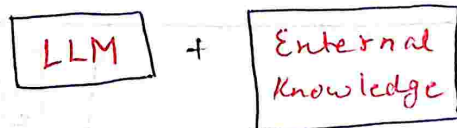
- connect application to external knowledge such as PDF, websites or database

scenario :- If we ask Policy of XYZ company to chat GPT {that chat GPT doesn't train on those policies}

than it fails to answer.

### main things →

- Doc loader
- Text splitter
- Vector store
- Retrievers



## # Memory

LLM API calls are stateless

Q1 → Who is Sunny Deol? → Answer ✓

Q2 → How old is he? → Answer ✗

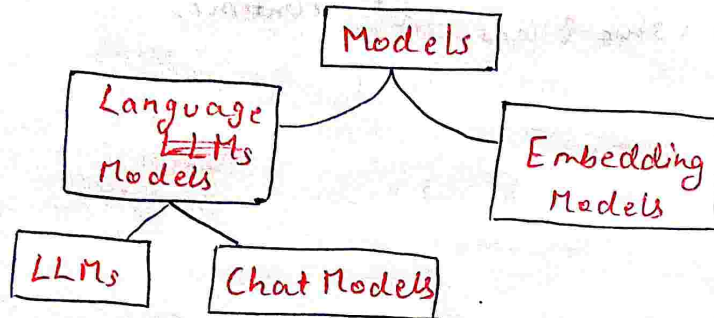
↳ Model/LLM को पता है ना कि हम किसकी बात कर रहे हैं।

- many types, 4 most usable



# # Models

- crucial part of framework
- design to facilitate interactions with various language models & embedding models.



Aspect	LLMs	Chat Models
Purpose :	Free-form text generation	Optimized for multi-turn conversations
Training Data	General text {Books, etc}	Fine tuned on chat datasets (dialogues, user-assistant conversations)
Memory & Content	No built-in memory	Supports structured conversational history
Role Awareness	No understanding of "user" & "assistant" roles	Understands "system", "user", & "assistant" roles
Examples	GPT-3, Llama-2-7B, OPT-1.3B	GPT-4, GPT-3.5-turbo, Llama-2-Chat
Use Cases	Text generation, summarization, translation, creative writing, code generation	Conversational AI, chatbots, virtual assistants, customer support, AI tutors

# Prompts

- Text Based
- Multi Modal {img}
- Default validation optional available

- 2 types

Static



query → Model

Dynamic

query → Template → Model

→ templates → save & load

## Message Placeholder

- used inside ChatPromptTemplate to dynamically insert chat history or a list of messages at runtime.

## Messages

- System Messages
  - Human Messages
  - AI Messages
- Static msg h  
but ~~EH~~ ~~for~~ Dynamic  
msg bna rhta h  
"Chat Prompt Template"

## #Structured Output

LLMs

can

with-structured-output

this fn calls before calling invoke,

methods to define data-formats

can't

output-parsers

- we can also use the output parsers with 'can' models

Typed Dict

Pydantic

json-schema



## Typed Dict

way to define a dictionary in Python where you specify what keys & values should exist. It ensures that your dict follows a specific structure.

ex: `name: str` | अगर age str का फ़ैल तो गलत Typed Dict  
`age: int` | If error नहीं होता

### No-validation

- We can send what about return from LLM by the use of Annotated fn
- We can use pydantic models output as dict & json
- json & typed dict both not have default values feature
- use function calling method in args → method when use openAI LLM else json-mode in args → method

## Output Parsers

raw LLM response → structured format

many parsers but we discuss

(4)  
string opt parser → structured opt par → helps to extract JSON Data from LLM's output  
pydantic opt P → json opt parser  
↓  
json opt parser doesn't enforce a schema

String Output Parser → return simple text

but useful to making chains  
ex → #6 video 20:00

↓  
cons  
↓  
No-validation

## Chains

- + Simple
- + Sequential
- + Parallel
- + Conditional

## Runnables

- unit of work
- common interface
- easy connection
- runnables  $\Rightarrow R_1 + R_2 + R_3 \dots$   $\Rightarrow$  workflow also called Runnable

### TYPES

Task Specific Runnables

- core langchain components
- eg ChatOpenAI, ChatPromptTem...

Runnables Primitive

Runnables Primitive

### # Runnable Sequence

- simple chain

### # Runnable Lambda

- convert any python function to Runnable

### # Runnable Parallel

### # Runnable Passthrough

### # Runnable Branch

- control flow component in langchain that allows you to conditionally route input data to different chains



# RAG

- technique that combines information retrieval with language generation, where a model retrieves relevant docs from a knowledge base & then uses them as content to generate accurate & grounded response.

## Components

- Document Loaders
- Text splitters
- Vector Databases
- Retrievers

- use up-to-date information
- better privacy
- No limit of document size

- TextLoader
- PyPDFLoader
- DirLoader → DirectoryLoader  
If many docs 100s, 1000s, ...  
then use lazy-load()  
↳ use generator
- WebBaseLoader
- CSVLoader

## Document loader ↑

## Text Splitting

Large Text → chunks

## Splitters

- length based → fast but not more useful → chunks overlap
- text structured → paragraph, line, word, character
- Document " " → perform by special words → ex we splitting python code files :-
- semantic meaning

```
class
    def __init__(self):
    def ...
    ...
```

then we specify  
"ndef", "in def",  
"in class", ...



# # Vector Stores

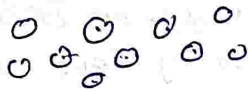
- a system designed to store & retrieve data

- Storage: you can save them

- similarity search

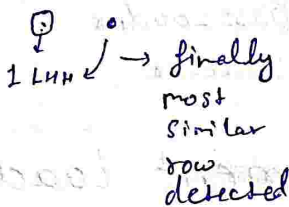
- Indexing: 10 LKH rows

↓  
clustering  
n=10



• → similarity ↑

↓



- CRUD operations

- eg. Chroma, FAISS, etc

## Use - Cases

- Semantic search

- RAG

- Recommender Search

- Image/Multimedia Search

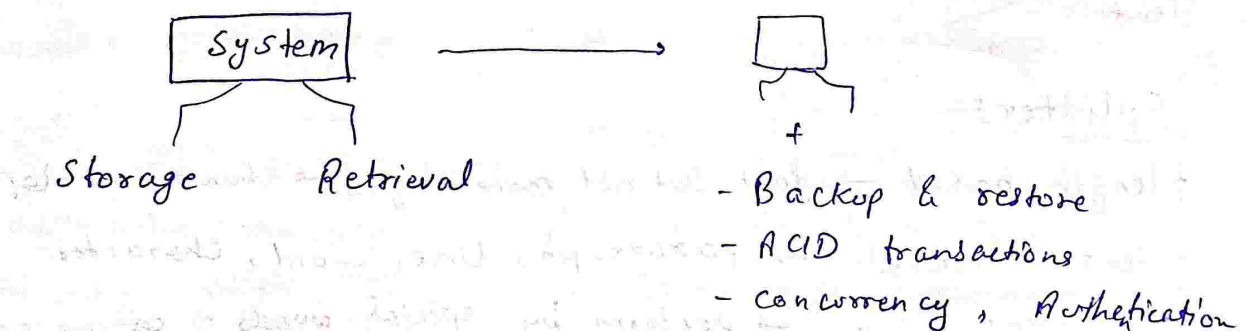
## Chroma Vector Store

- open-src

- small to medium scale production needs

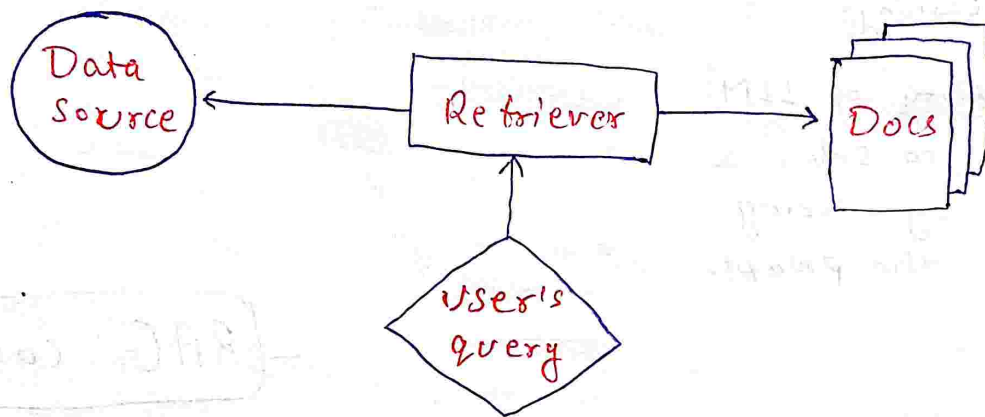
- Lightweight

## Vector Store V/s Vector Database



# Retrievers

- a component in lang chain
- It fetches relevant documents from a data source in response to a user's query.
- many types
- All retrievers in the lang chain are runnables.



## - Types

Wikipedia Ret.	Based on Data
Vector store Ret.	
MMR	Based on Search method
Multi-Query	
etc	

## Contextual Compression Retriever

- compressing docs after retrieval

query → —  
document →   
imp

## MMR

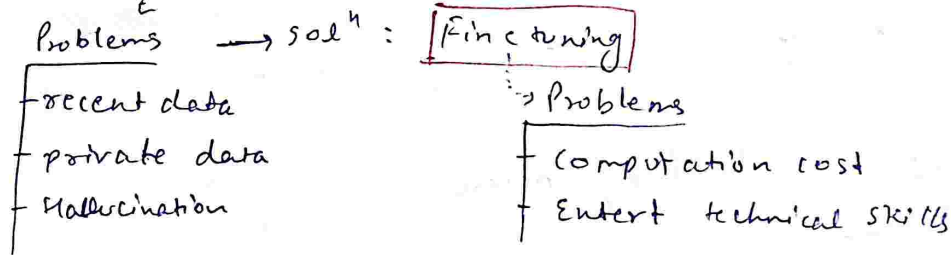
- Maximal Marginal Relevance
- Working  
1<sup>st</sup> Doc → most relevant  
than rest → most relevant & least similar  
& so on.

## MQR

- Multi Query Retriever
- ex: How can I stay healthy?  
could mean
  - + How should I eat?
  - + How often should I exercise?
  - + How can I manage stress?

# RAG Concepts

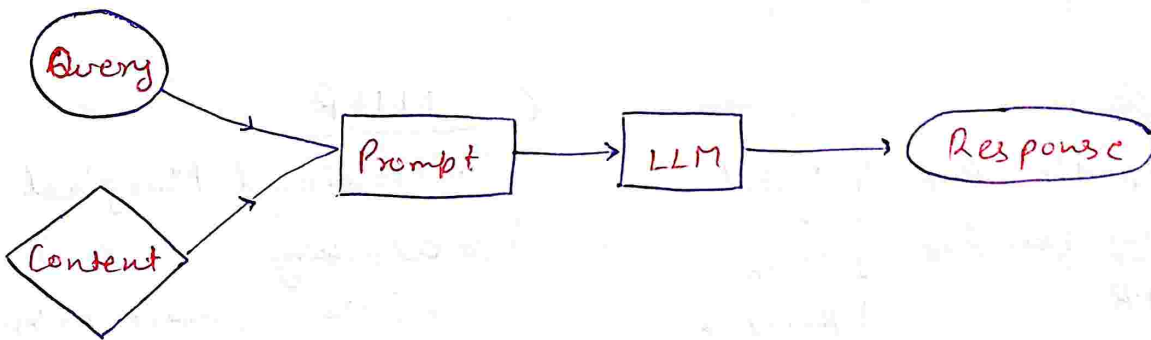
Query  $\rightarrow$  **LLM**  $\rightarrow$  Response



## **In-Content Learning**

- core capability of LLM
- model learns to solve a task purely by seeing examples in the prompt.

RAG: core.c



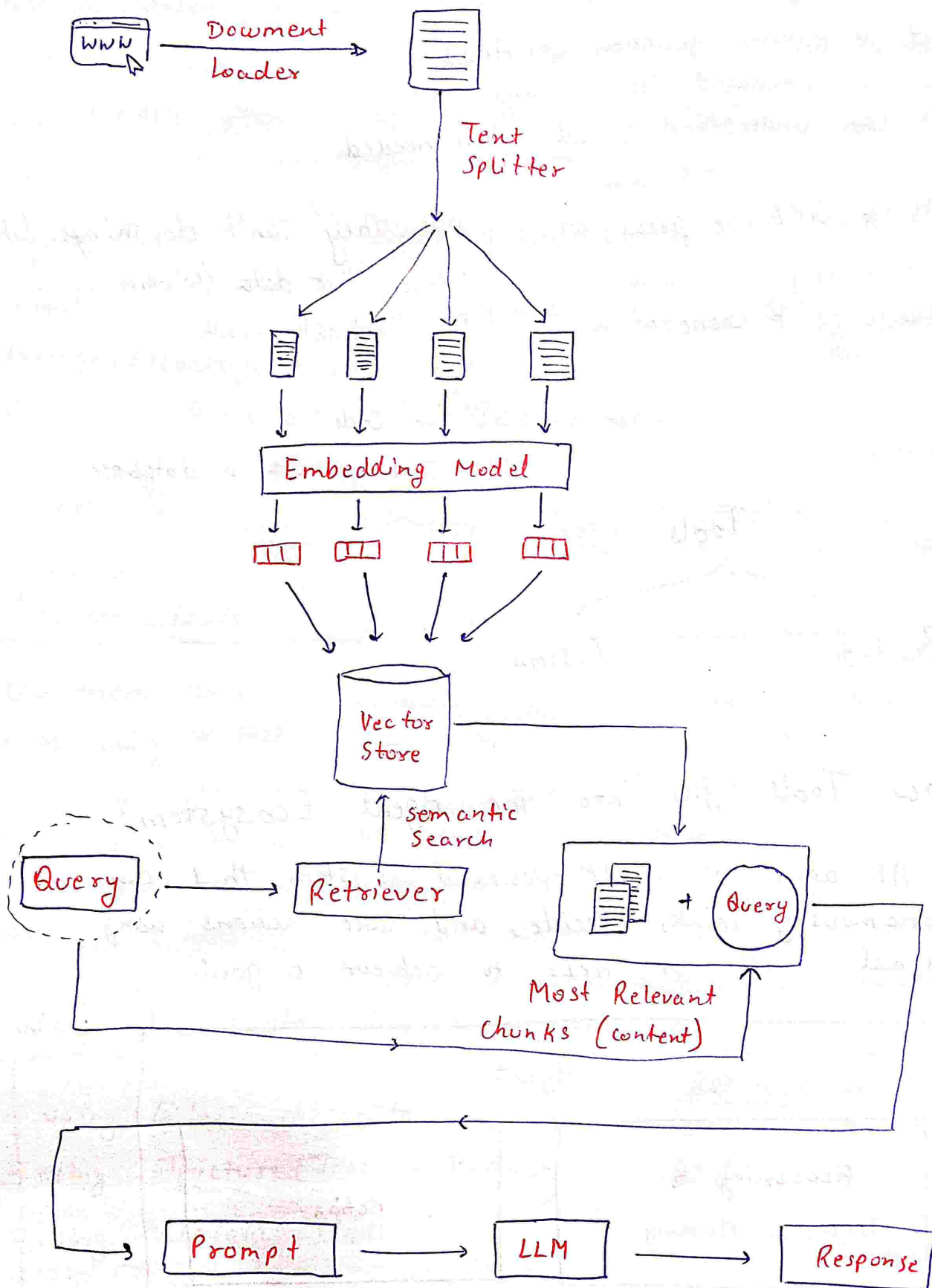
RAG  $\rightarrow$  Information Retrieval + Text-Generation

## Steps

- Indexing
- Retrieval
- Augmentation
- Generation



# Basic Architecture



# Tools

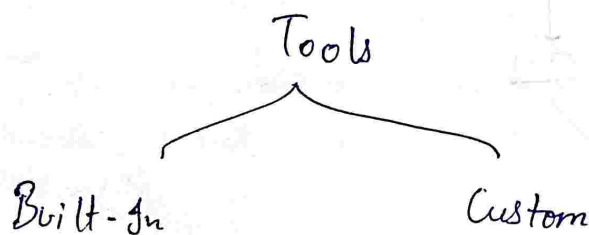
- just a python function (or API) that is packaged in a way that LLM can understand & call when needed

LLMs & GPTs are great at:

- Reasoning
- Language & Generation

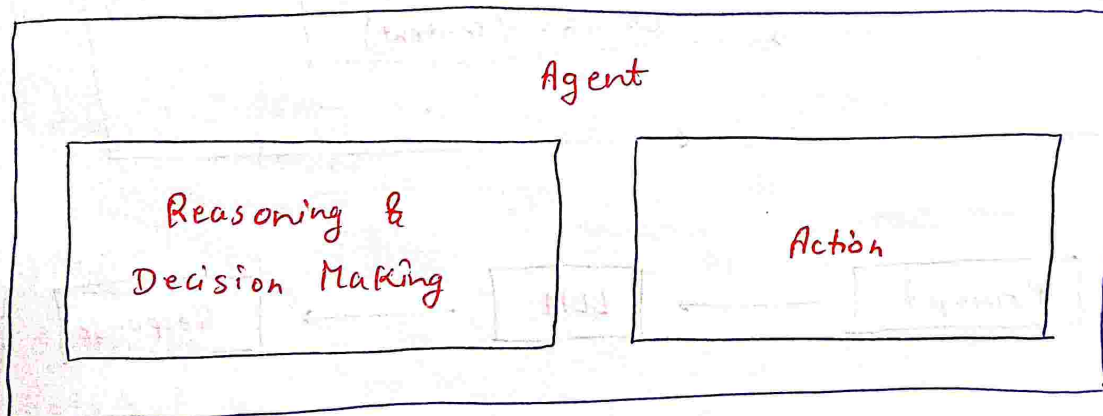
But they can't do things like

- Access live data (weather, news)
- Do reliable math
- Call APIs
- Run code
- Interact with a database



## # How Tools fits into the Agent Ecosystem?

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



## Built In Tools

- DuckDuckGoSearchRun : Web search via duck duck go
- WikipediaQueryRun : Wikipedia summary
- PythonREPLTool : Run Raw Python code
- ShellTool : Run shell commands
- RequestsGetTool : Make HTTP GET Requests
- GmailSendMessageTool : Sends messages via Gmail
- SlackSendMessageTool : Post Message to slack
- SQLDatabaseQueryTool : Run SQL Queries
- !
- more

## Custom Tools

Use them when:

- You want to call your own APIs
- You want to encapsulate business logic
- You want the LLM to interact with your database, product or app.

## Ways to create Custom Tools

- |  |  |
|--|--|
| <ul style="list-style-type: none"><li>+ using @tool decorator</li><li>+ using StructuredTool &amp; Pydantic</li><li>+ using BaseTool class</li></ul> | <p><u>BaseTool</u> is the abstract class for all tools in langchain.</p> |
|--|--|

A Structured Tool in langchain is a special type of Tool where the input to the tool follows a structured schema, typically defined using a Pydantic model.



## Tool Binding

- step where you register tools with a Language Model (LLM) so that:

- The LLM knows what tools are available
- It knows what each tool does (via description)
- It knows what format to use (via schema)

NOTE: All LLMs does not support Tool binding

## Tool Calling

It is the process where the LLM decides, during a conversation or task, that it needs to use a specific tool (function) — and generates a structured output with:

- the name of the tool
- the arguments to call it with

NOTE: LLM does not actually run the tool — it just suggests the tool and the input arguments,

The actual execution is handled by LangChain or you

### Example

query → What's 8 multiplied by 7?

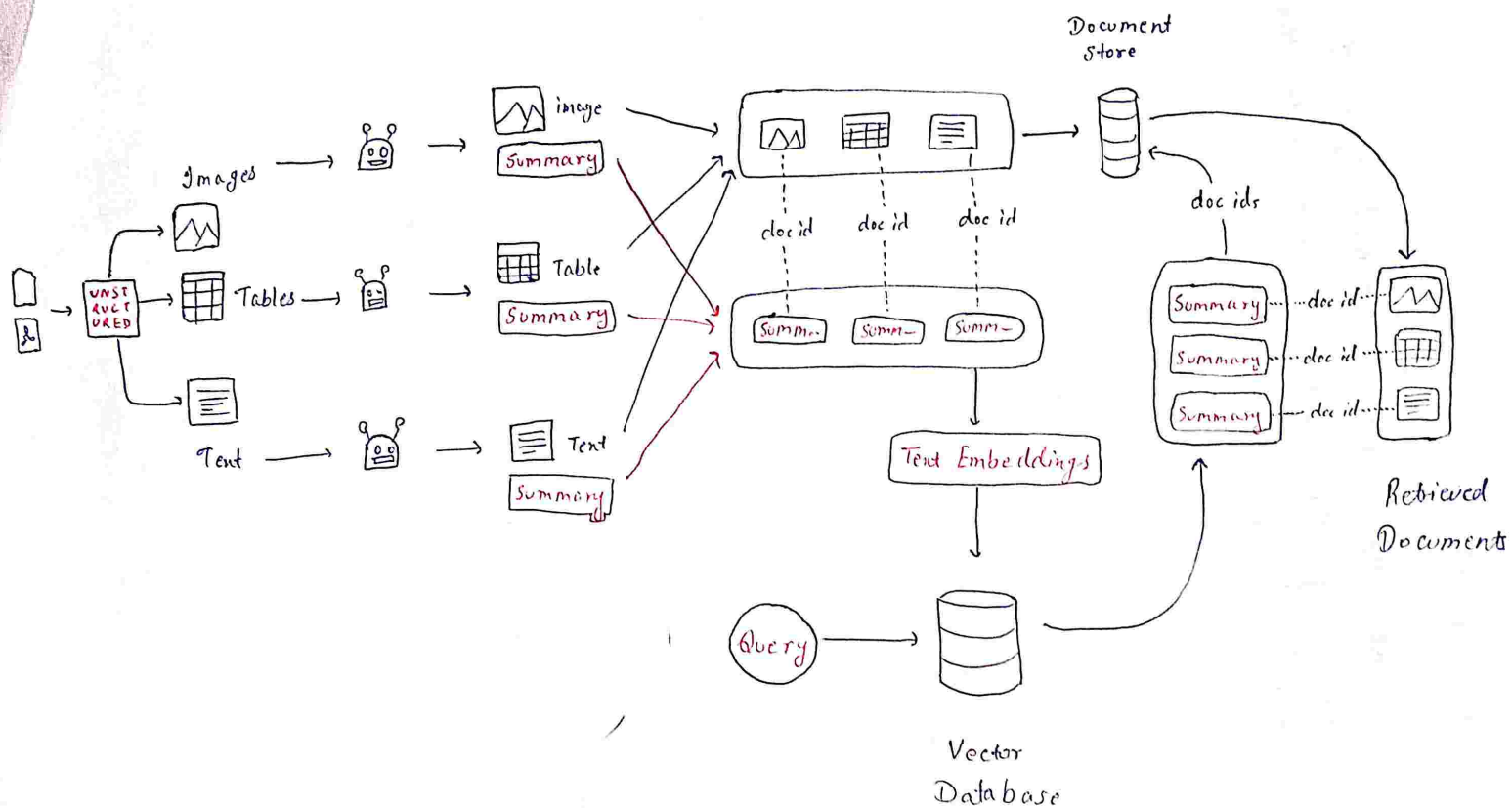
The LLM responds with a tool call:

```
{  
  "tool": "multiply",  
  "args": { "a": 8, "b": 7 }  
}
```

## Tool Execution

It is the step where the actual Python function (tool) is run using the input arguments that the LLM suggested during tool calling.

## Multimodal Retrieval using Unstructured For Extraction





# More Topics

## ★ Evaluation → Metrics

- Rags
- LangSmith

- Faithfulness
- answer relevancy
- content Precision
- content Recall

## ★ Indexing

- Document Ingestion
- Text Splitting
- Vector Store

## ★ Retrieval

### a) Pre-Retrieval

- Query Rewriting using LLM
- Multi-Query generation
- Domain aware routing

### b) During Retrieval

- MMR
- Hybrid Retrieval
- Reranking

### c) Post Retrieval

- Contextual Compression

## ★ Augmentation

- Prompt Templating
- Answer grounding
- Contextual Window Optimization

## ★ Generation

- Answer with Citation
- Guard Railing

## ★ System Design

- Multimodal
- Agentic
- Memory based