

LangChain :-

ACHIEVER

Date: / /

LangChain is an open-source orchestration framework used for application development using large language models (LLMs).

- * LangChain's tools and APIs simplify the process of building LLM-driven applications like chatbots and AI-agents.
- * The LangChain modular environment allows for programs that use one or multiple LLMs:
For ex:, an application that uses one LLM to interpret user queries and another LLM to author the responses.

Families of LLMs

- ① OpenAI → GPT-3.5, GPT-Turbo-40 (comfy)
- ② Anthropic → Claude 1, 2, 3, 3.5
- ③ Google DeepMind → Gemini 1, 1.5, 2.0...
- ④ Meta → Llama 2 (os), Llama3 etc
- ⑤ Qwen → Qwen 1.5, 2.0, 2.5, Qwen 3
- ⑥ Mistral → Mistral 7B, Mistral 8x7B etc.

Integration with LLM's

* LLM's are the brain. Applications are body.
Data + API's are the sense and tools.
orchestrators frameworks are the nervous system that connects it all.

① LLM's are not applications :-

- * They're statistical models trained on massive datasets.
- * They must be paired with an application to serve real-world usecases.

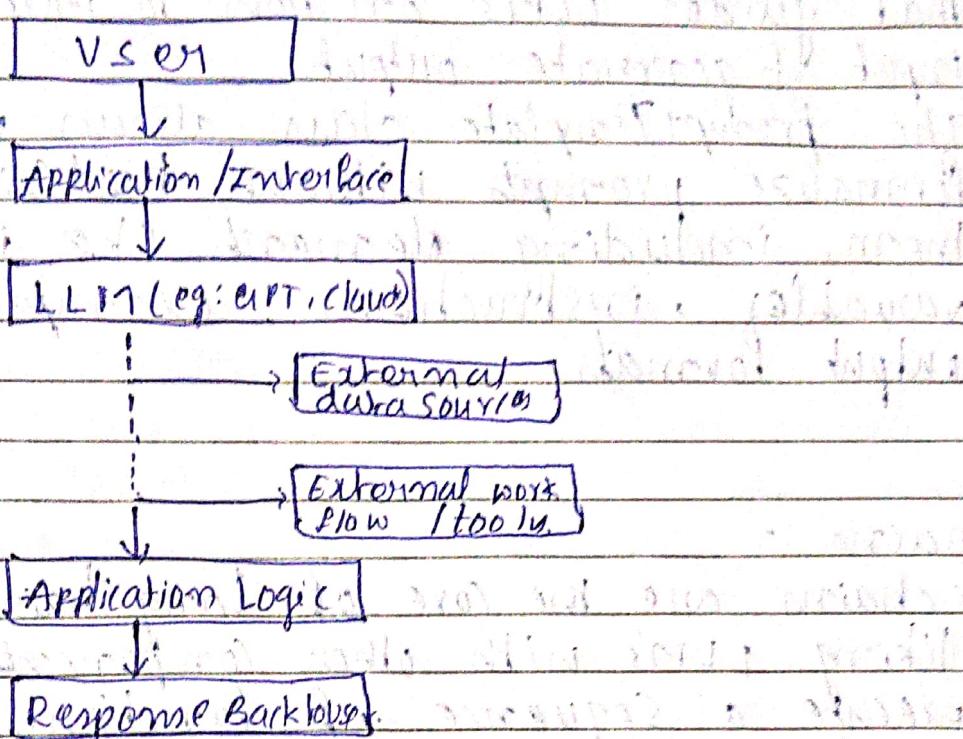
② Need for external data sources :-

- * LLM's are not omniscient.
 - For real-time awareness or domain specific tasks, LLM's must connect to:
 - internal docs, databases or API's
 - news/current sources
 - company specific knowledge.

③ Need for external work flow:-

- * LLM's often need to call API's or tools.
 - Ex: Slack Integration → LLM must connect with slack API

Diagram style work flow :-



Explanation:-

- ① User → sends a query (e.g. summarize this PDF)
- ② App/Interface → collects the query (chat GPT UI) etc
- ③ LLM → interprets input, generate potential response
- ④ External Data → if info is outside 'models' memory (like company docs), it retrieves via API (Retrieval-Augmented Generation)
- ⑤ External tools/work flow → if task requires action (searching web, calling slack API) LLM calls tools.
- ⑥ Orchestration layer → connects everything, manages workflow, decides when to call what [LangChain, IBM Watson]
- ⑦ Response → Sent Back to user in natural language format.

Prompt Template :-

- * prompt templates are structured instructions that guides LLM's on how to interpret input & generate output.
- * The PromptTemplate class allows you to formalize prompts without hardcoding them, including elements like input variables, instructions, examples and output formats.

Chain :-

- * chains are the core of long chain workflow linking LLM's with other components to execute a sequence of functions.

Ex:-

```
chain_ex = LLMChain(llm=Plants, prompt=ExamplePrompt)  
chain_ex.run ("input")
```

- * Simple sequential chain: uses the output of one function as input for the next
- * Each step can use different prompts, tools, parameters or models.

Indexes :-

To achieve certain tasks, LLM will need access to specific external data sources not included in its training dataset, such as internal documents, domains or databases. LangChain collectively refers to such external documentation as "indexes".

① Document Loaders:-

- * LangChain offers a wide variety of document loaders for third-party applications.
- * This allows for easy importation of data from sources like file storage services (like Dropbox, Google Drive, and Microsoft OneDrive), web content (like YouTube, PubMed, or specific URLs), databases (like Pandas, MongoDB, and MySQL), among many others.

② Vector Databases:-

- * Represent data as vector embeddings (numerical vectors).
- * Enable low-latency queries even for large datasets.
- * LangChain supports 9+ embedding methods and 50+ vector stores.

③ Text Splitters:-

- * Divide large texts into smaller, semantically meaningful chunks to speed up processing and reduce computational load.

④ Retrieval & RAI (Retrieval-Augmented Generation) :-

- * Models retrieve relevant info from connected data sources using retriever modules.
- * Agentic RAI systems can also perform tasks beyond retrieval (eg: calculations, emails, data analysis).

LangChain use case:-

① AI Applications :- LangChain enables a wide range of AI applications, from simple Q&A and text generation to complex reasoning tasks.

② Chatbots :- provides context-specific responses, integrate into existing workflows and communication channels using APIs.

③ Summarization :- summarizes complex articles, transcripts or incoming emails.

④ Question Answering :- Retrieve information from documents or specialized knowledge bases (eg: Wolfram, arXiv, PubMed).

⑤ Data Augmentation :- generate synthetic data to expand training datasets for ML.

Wrapper classes for different Providers :-

① Google Gemini :

- ChatGoogleGenerativeAI → chat-based models
- Google Generative AI Embeddings → for Embedding

② Open AI :

- ChatOpenAI → chat-based GPT models
- OpenAIEMBEDDINGS → embedding with oA models

③ Hugging Face :

- Hugging Face Hub → run models hosted on HF Hub
- HuggingFacePipeline → run local Hugging Face model with transformers.

④ Ollama (local models) :

- ollama → run models like Llama, Mstral, etc