

Components of LangChain

1. Introduction to LangChain and its Purpose

LangChain is an open-source framework designed to simplify the development of LLM-powered applications. It addresses the complexity of building such systems from scratch by providing an efficient orchestration mechanism for various components.

- **Necessity of LangChain:** Building LLM-powered applications, such as a system that allows users to "talk to a PDF," involves numerous components and intricate interactions. LangChain streamlines this process, offering "minimal code for maximum output."
- **Orchestration and Pipelines:** A key advantage of LangChain is its ability to "orchestrate very efficiently between all these components" and "build pipelines." This means that the output of one component automatically becomes the input for the next, eliminating the need for manual coding.
- **Model Agnostic Framework:** LangChain is designed to be model-agnostic. This flexibility allows developers to switch between different LLM providers (e.g., OpenAI's GPT models to Google's Gemini models) with "literally one to two lines of code change," minimizing code modifications.
- **Applications Built with LangChain:** LangChain is currently used to create various applications, including "conversational chatbots," "AI knowledge assistants," and "agents."

2. The Six Core Components of LangChain

LangChain comprises six fundamental components that are essential for understanding and utilizing the framework:

1. **Models**
2. **Prompts**
3. **Chains**
4. **Memory**
5. **Indexes**
6. **Agents**

These six components form the backbone of the LangChain framework, and mastering them provides a deep understanding of its capabilities.

3. Deep Dive into LangChain Components

3.1. Models

The "Models" component is considered the most critical in LangChain, serving as the "core interface through which you interact with AI models."

- **Addressing LLM Challenges: Initial NLP Problems (NLU & Text Generation):** Historically, building chatbots faced two major challenges: Natural Language Understanding (NLU) – making the chatbot understand user queries, and context-aware text generation – enabling the chatbot to provide relevant replies. LLMs, trained on vast internet data, simultaneously solved both these problems, developing "understanding of natural language" and "context-aware text generation capability."
- **LLM Size and Accessibility:** LLMs are "huge," with many exceeding 100GB, making them impractical for individuals or small companies to run locally due to infrastructure and cloud costs. This problem was solved by companies like OpenAI providing "APIs," allowing users to interact with LLMs remotely and pay only for usage.

- **API Standardization:** A new problem emerged: different LLM providers (OpenAI, Anthropic, Google) implemented their APIs differently, requiring developers to write "different types of code" to interact with each.
- **LangChain's Solution: Standardized Interface:** LangChain's Model component provides a "standardized interface" that allows communication with "any company's AI model in a standardized fashion." This means developers can switch between providers like OpenAI and Anthropic with "literally two lines of code change." The responses received are also "very similar types," making parsing easier.
- **Types of Models Supported:** LangChain communicates with two types of AI models:
- **Language Models (LLMs):** These models take "text input" (e.g., "How are you today?") and produce "text output" (e.g., "I am good, how about you?"). They are used for applications like chatbots and AI agents.
- **Embedding Models:** These models take "text as input" but output a "vector." Their primary use case is "semantic search."
- **Provider Agnostic:** LangChain supports a wide range of LLM and embedding model providers, including OpenAI, Anthropic, Mistral AI, Azure, Vertex AI, Bedrock (AWS), Hugging Face, IBM, and Llama.

3.2. Prompts

Prompts are the "inputs provided to an LLM," and they are "super important" because an LLM's output is "highly dependent on prompts."

- **Prompt Sensitivity:** Even minor changes in prompts can significantly alter an LLM's output. For example, asking an LLM to "Explain Linear Regression in academic tone" versus "Explain Linear Regression in fun tone" will yield vastly different responses.
- **Prompt Engineering:** The importance of prompts has led to the emergence of "Prompt Engineering" as a field of study and "Prompt Engineers" as a job profile.
- **LangChain's Prompt Component:** LangChain recognizes the importance of prompts and offers a component to "handle prompts very well." It provides "a lot of flexibility" to create "different and powerful types of prompts."
- **Types of Prompts in LangChain: Dynamic and Reusable Prompts:** These prompts use "placeholders" that can be replaced at runtime based on user input, making them adaptable and reusable (e.g., "Summarize this topic in this tone," where "topic" and "tone" are placeholders).
- **Role-Based Prompts:** Developers can define system-level prompts to guide the LLM's persona (e.g., "You are an experienced [profession]," where "profession" is a placeholder), influencing the LLM's response style.
- **Few-Shot Prompts:** This technique involves providing the LLM with "some examples" of input-output pairs before asking a question. This helps the LLM understand the desired output format and behavior for new, similar queries (e.g., classifying customer support tickets based on examples).

3.3. Chains

Chains are a "very important component" in LangChain, so much so that "LangChain is named after it." They enable the creation of "pipelines" in LLM applications.

- **Building Pipelines:** Chains allow developers to represent the flow of an LLM application as a pipeline, where "one stage comes after another."
- **Automatic Input/Output Handling:** The biggest advantage of Chains is that they "automatically make the output of the previous stage the input of the next stage." This eliminates the need for manual coding to pass results between different LLM calls or stages. For example, in a "Hindi summary" application, the English text is translated by one LLM, and that translation is then automatically fed to a second LLM for summarization.

- **Complex Chains:** While simple sequential chains are common, LangChain also allows for the creation of more complex pipeline structures:
- **Parallel Chains:** These involve sending the same input to multiple LLMs simultaneously and then combining their outputs (e.g., generating two reports from an input and then combining them with a third LLM).
- **Conditional Chains:** These execute different processing paths based on a specific condition (e.g., an AI agent handling user feedback, thanking good feedback, but emailing customer support for bad feedback).

3.4. Indexes

Indexes connect LLM applications to "external knowledge sources such as PDFs, websites, or databases." They are crucial for enabling LLMs to answer questions about private or domain-specific data that they were not trained on.

- **Addressing Knowledge Gaps:** Standard LLMs like ChatGPT cannot answer questions about private company policies because they weren't trained on that data. Indexes solve this by providing "external knowledge sources" to the LLM.
1. **Four Main Components of Indexes:**
 - Document Loader:** Loads data from various sources (e.g., PDFs from Google Drive).
 - Text Splitter:** Breaks large documents into "smaller chunks" (e.g., pages, paragraphs) for efficient semantic search.
 - Vector Store:** Converts document chunks into "embeddings" (vectors) using an embedding model and stores them in a "vector database" for future retrieval.
 - Retrievers:** Takes a user query, generates its embedding, performs a "semantic search" in the vector store to find relevant document chunks, and then provides those chunks along with the user query to the LLM for generating a coherent answer.
 - **Simplified Implementation:** LangChain "makes this entire process very easy," significantly reducing the amount of code required to build such systems.

3.5. Memory

The Memory component addresses the "big problem" of "stateless" LLM API calls, where each request is independent and the model has "no memory of the previous request."

- **Stateless Nature of LLMs:** LLM API calls are inherently stateless, meaning they do not retain conversational history. If a user asks "Who is Narendra Modi?" and then "How old is he?", the LLM will not understand "he" refers to Narendra Modi because it has no memory of the previous turn. This leads to a "frustrating" user experience in conversational applications.
- **LangChain's Solution: Adding Memory to Conversations:** LangChain's Memory component enables "adding memory features to your entire conversation."
- **Types of Memories in LangChain:**
 - Conversation Buffer Memory:** Stores the "entire conversation history" and sends it with each subsequent API call. While effective, it can become costly for long conversations due to the increasing amount of text processed.
 - Conversation Buffer Window Memory:** Stores only the "last N interactions," constantly updating to maintain a manageable history size.
 - Summarizer Based Memory:** Generates a "summary" of the entire chat history and sends only the summary with the API call, helping to "save some text" and reduce costs.
 - Custom Memory:** Allows for storing "specialized pieces of information" like user preferences or facts for advanced use cases.

3.6. Agents

Agents are a "next big thing" in AI, allowing developers to build "AI agents" with LangChain. AI agents extend the capabilities of chatbots by enabling them to perform actions.

- **Chatbots vs. AI Agents:** **Chatbots:** Possess NLU and text generation capabilities, making them good for understanding and replying to queries (e.g., "What is the best travel destination in India during summer?").
- **AI Agents:** Are "chatbots with superpowers." They can not only understand and reply but also "do some work for you." For example, an AI agent could find the "cheapest flight" between two cities on a specific date and even "book the flight."
- **Key Capabilities of AI Agents:** AI agents have two crucial features that differentiate them from chatbots:
 1. **Reasoning Capability:** They can "reason what exactly they need to do" by breaking down queries step by step (e.g., using "Chain of Thought prompting").
 2. **Access to Tools:** They can access and utilize "different tools" (e.g., a calculator for mathematical operations, a weather API for fetching weather conditions, or a booking API for booking flights).
- **How AI Agents Work (Example):** An AI agent tasked with multiplying today's temperature of Delhi by three would:
 1. **Reason:** Break down the request into steps: get Delhi's temperature, then multiply it by three.
 2. **Tool Use (Weather API):** Check if it has a tool to get the temperature. It finds a "Weather API," calls it with "Delhi," and gets "25 degrees Celsius."
 3. **Tool Use (Calculator):** Realize it needs a calculator for multiplication. It calls the "Calculator" tool with "25" and "3" and the "multiplication operation."
 4. **Output:** The calculator returns "75," which becomes the final output.

In essence, an AI agent is an "evolved form of chatbot" that can perform actions due to its "reasoning capabilities and access to tools." LangChain "has made it very easy to build AI agents."