# LangChain Introduction

## 🔷 Summary (At a Glance)

- **LangChain** is a **Python framework** that helps developers build applications using **Large Language Models (LLMs)** like GPT (ChatGPT), Claude, etc.

- It **chains** together different components (like prompts, memory, tools, documents, agents) to create more **intelligent, interactive applications**.

- Think of it like **LEGO blocks for AI** apps – each block does a specific job, and you combine them to build smart applications like chatbots, personal assistants, search engines, and more.


✅ **LangChain** (building)

✅ **LangSmith** (debugging, evaluating, monitoring)

✅ **LangServe** (deployment)

## 🔷 Why LangChain?

Most LLMs (like GPT) work in a very **basic, one-shot way**:

You give a prompt → it gives a response. That's it.

But real apps need more:

- Memory of past conversations

- Access to tools like Google or a calculator

- Access to your files or documents

- Multi-step reasoning

- Using APIs or databases

**LangChain solves this** by connecting (chaining) these parts together.

## Step-by-Step Flow:

1. User: "Who is Elon Musk?"
2. LangChain checks memory → nothing yet
3. It sends the question to Wikipedia tool
4. Gets data → sends to LLM (GPT)
5. GPT makes a nice answer
6. LangChain saves this conversation in memory
7. User: "What companies does he own?"
8. LangChain checks memory → knows you meant Elon Musk
9. Sends to LLM again, with context
10. Returns smart answer

# 🔧 What LangChain Does?

LangChain simplifies the process of building LLM-powered applications by focusing on:

1. **Prompt Management**

   - Create, test, and manage dynamic prompts.

   - Easily swap between different models (OpenAI, Anthropic, etc.).

2. **Chains**

   - Combine multiple components (e.g., LLM calls, APIs, tools) into workflows.

   - For example: Take user input → Search web → Generate a summary.

3. **Agents**

   - Allow models to decide which tools to use and in what sequence.

   - Useful for dynamic tasks like data retrieval or multi-step reasoning.

4. **Memory**

   - Maintain state across interactions.

   - Useful for chatbots or applications needing context over time.

5. **Integrations**

   - Connect to external data sources (databases, APIs, vector stores).

- Popular integrations: Pinecone, Weaviate, OpenAI, Hugging Face, etc.

## 🔍 Use Cases

- Conversational AI (chatbots with memory)

- Retrieval-Augmented Generation (RAG)

- Question Answering over documents

- Workflow automation (e.g., using agents + APIs)

- Summarization, translation, code generation

## 🔷 Core Concepts of LangChain (Explained Simply)

| Concept | What it is (Simple Explanation) | Example |
|---|---|---|
| **LLM** | Large Language Model (e.g., GPT-4) | ChatGPT→ `OpenAI()` , `ChatAnthropic()` |
| **Prompt** | The input we give the LLM | "Tell me a joke" → `PromptTemplate("Tell me about {topic}")` |
| **Chain** | A sequence of steps, like a pipeline | Prompt → LLM → Response ( `LLMChain` , `SequentialChain` ) |
| **Agent** | A smart decision-maker that uses tools | "Use a calculator to solve 2+2" ( `initialize_agent(tools=[...])` ) |
| **Tool** | An external utility like a web search or Python function | Google Search |
| **Memory** | Saves conversation history | Like ChatGPT remembers the chat → `ConversationBufferMemory()` |
| **Document Loader** | Loads data from files (PDF, CSV, etc.) | Reading your resume |
| **Retriever** | Finds relevant documents when asked a question | Like a librarian for documents |
| **Indexes** | Connects LLMs to external data | `VectorStoreRetriever()` |

# LangChain + LangSmith: The Full Development Cycle

> LangSmith is LangChain's **observability + debugging platform** for LLM apps.

LangChain isn't just about chaining LLM calls—it offers **end-to-end tools** for building, testing, and deploying LLM apps.

## 📦 Features LangSmith adds:

| Feature | What It Does |
|---------|--------------|
| 🔍 **Debugging** | See how each chain, agent, tool, LLM step behaves internally |
| 📊 **Evaluation** | Run A/B tests, grade answers, track metrics over time |
| 📈 **Monitoring** | Log real usage from your deployed app (traces, latency, errors, token usage) |
| 🛠️ **Traces** | View full history of every call in your chain or agent |
| 🧪 **Custom evaluators** | Plug in your own logic to grade LLM answers (e.g., accuracy, helpfulness) |

You can think of LangSmith as:

> "The microscope and the dashboard" for your LLM apps.
>
> It shows **why something is working or breaking**, and helps you **fix and improve** it.

## 🧠 Why It's a Big Deal?

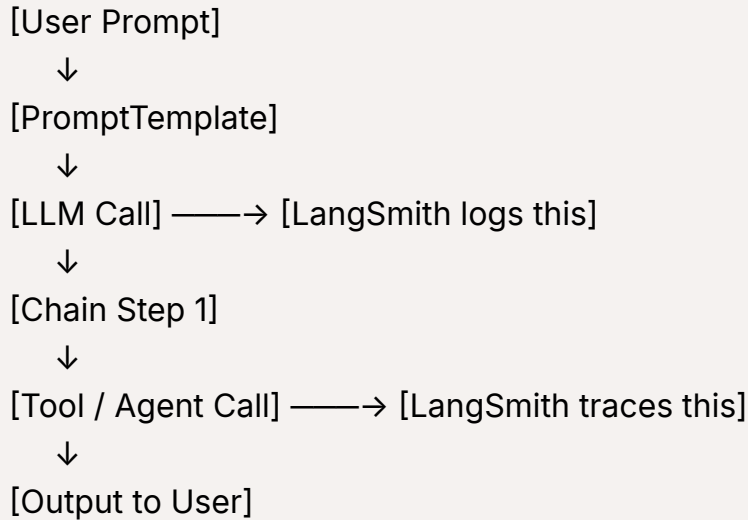Most LLM apps are **unpredictable**:

- Responses vary by temperature, input wording, token limits
- Debugging a multi-step agent is **almost impossible** with just print statements

LangSmith solves this:

- Logs **every chain, prompt, variable, output**
- Lets you **replay, analyze, and test** flows visually

- Helps you **compare different models** and see **what actually works**

## 🧪 Visual Diagram:

```
[User Prompt]
    ↓
[PromptTemplate]
    ↓
[LLM Call] ——→ [LangSmith logs this]
    ↓
[Chain Step 1]
    ↓
[Tool / Agent Call] ——→ [LangSmith traces this]
    ↓
[Output to User]
```

## 🔧 LangChain + LangSmith Workflow in Real Life:

| Stage | What You Use |
|---|---|
| Build logic | LangChain ( chains , agents , etc.) |
| Run it live | In your app / notebook / API |
| Log + Trace | LangSmith dashboard |
| Debug errors | See where chains/agents failed |
| Evaluate quality | LangSmith or your own graders |
| Improve | Tune prompts, models, parameters |

# LangServe

## What is LangServe?

- **LangServe** lets you **turn any LangChain logic into a REST API**, instantly.

It's like:

> "Take your LangChain chain/agent, and deploy it as a ready-to-use backend server."

## 🛠️ Why Use LangServe?

| Need | LangServe Solves It |
|------|---------------------|
| You built a chatbot in LangChain | But want to call it from a front-end (React/Flutter) |
| You made a document Q&A agent | And want to expose it as a web service |
| You have custom logic (tools, agents, memory) | And want to deploy without building a server manually |
| You want it to scale or connect via HTTPS | LangServe can be hosted (e.g., on Hugging Face Spaces, Render) |

# Templates

## What Are LangChain Templates?

- Templates in LangChain are **pre-made app blueprints** — best practices wrapped into a starter project.
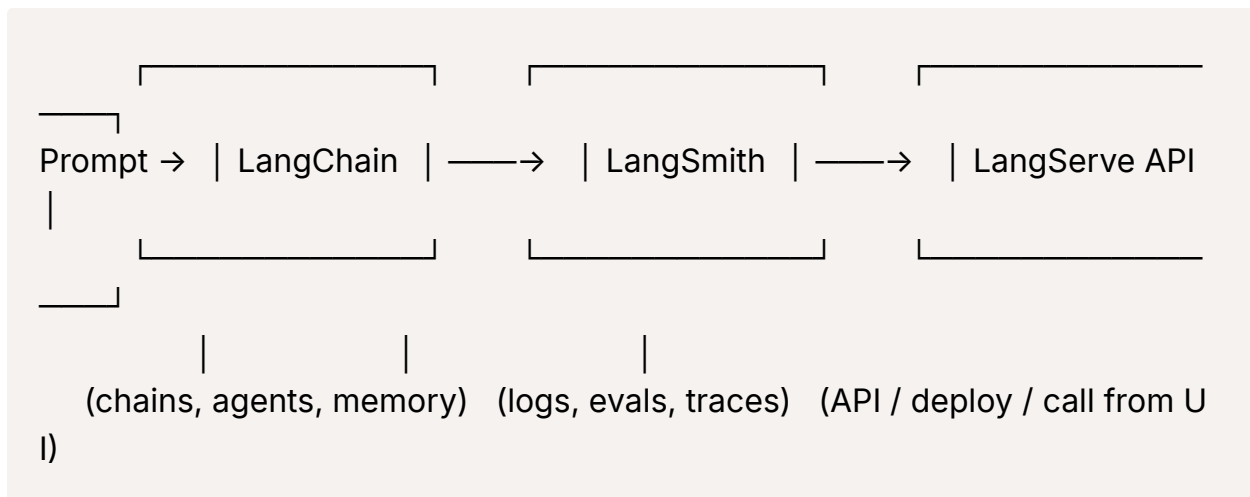
- They are ideal for:

  > "Just give me a working chatbot / RAG app / tool-using agent I can start modifying."

### 🚀 Example Templates (official and community)

| Template Name | What It Does |
|---------------|--------------|
| rag-langchain | RAG (Retrieval-Augmented Generation) setup |
| chat-with-tools | Chatbot that can use search, calculator, Python, etc. |

| Template Name | What It Does |
|---|---|
| multi-agent-debate | 2+ LLMs debating with memory and rules |
| document-QA | Upload + ask questions about PDF/docx/markdown |

## ✅ Visual Overview

```
        ┌──────────────┐      ┌──────────────┐      ┌──────────────┐
    ┌───┘
Prompt → │ LangChain │ ────→ │ LangSmith │ ────→ │ LangServe API
    │
        └──────────────┘      └──────────────┘      └──────────────┘
    ┌───┘

              │              │              │
      (chains, agents, memory)   (logs, evals, traces)   (API / deploy / call from U
I)
```

## 🔷Trivia

- LangChain was **created in 2022** by Harrison Chase.

- It became very popular because it was **the first framework** focused on building **complex LLM apps**.

- You can combine it with **tools like Streamlit, Gradio, FastAPI, or Flask** to make real web apps.

# Python Code:

### Free Inference API (Hugging Face Hub)

```
!pip install langchain huggingface_hub
```

```
from langchain_community.llms import HuggingFaceHub
```

```
llm = HuggingFaceHub(
    repo_id="deepseek-ai/deepseek-llm-7b",
    huggingfacehub_api_token="YOUR_HF_TOKEN",  # Get at: https://huggingfa
ce.co/settings/tokens
    model_kwargs={"temperature": 0.7}
)

response = llm("Explain quantum computing briefly")
print(response)
```

❗ **Limits**: Uses your Hugging Face free tier (**1k requests/day**).