**Langchain - Detailed Explanation**

**what is LangChain?**

- **Langchain** is a **framework** used to build applications powered by **Large Language Models**

- LangChain is a powerful framework designed to build applications using Large Language Models (LLMs).

- It helps developers connect language models with external data sources, tools, memory, and logic in a structured and reusable way.

- Traditional LLM usage is limited to prompt–response interactions, but LangChain extends this by enabling document understanding, retrieval-based question answering, conversational memory, tool usage, and agent-based reasoning.

- In real-world applications, data is often stored in PDFs, databases, websites, or APIs.

- LangChain acts as a bridge between raw data and intelligent AI responses. It is widely used for building chatbots, document Q&A systems, search engines, summarizers, and AI assistants.

- like GPT,Gemini,Claude,etc

- it helps in:

  - Handling Long documents
  - Connecting models+data+pr`mompts+tools
  - Building chatbots,RAG systems,agents,assistants

**Load the Documents-Documents Loaders**

**why needed?**

- LLMs cannot directly read files like PDF,Word,CSV,websites,etc.

- so Langchain provides Document Loaders

- Large Language Models like GPT, Gemini, or LLaMA are powerful but stateless and context-limited.

- They cannot:

- Remember long conversations

- Access private documents directly

- Perform multi-step reasoning reliably

- Interact with tools or APIs by default

**what they do:**

- Load raw data From different sources

- Convert it into a standard document format

**Document Loaders (In-Depth)**

- Document Loaders are responsible for loading data from different sources and converting it into text format that LLMs can understand.

- Supported sources include:

  - PDF files

  - Word documents

  - Text files

  - CSV and Excel files

  - Websites and APIs

  - lDatabases

**Examples**

- PDF
- TXT
- DOCX
- CSV
- Web URLs
- Databases

**Common Loaders:**

- PyPDFLoader
- TextLoader
- CSVLoader
- WebBaseLoader

**LangChain solves these problems by:**

- Adding memory to conversations

- Connecting LLMs to vector databases

- Structuring prompts systematically

- Enabling reasoning through chains and agents

**This makes LangChain suitable for enterprise-grade AI applications.**

**Langchain Load Different Models(LLMs)**

**Large Language Models (LLMs)**

- LLMs generate the final response using the prompt and retrieved context.
- LLMs do not search data themselves; they rely entirely on provided context, making retrieval accuracy critical.

**why need**

- different applications need different models:
  - Chatbots
  - Summarization
  - Q&A
  - Code generation

**Models LangChain Supports:**

- OpenAI(GPT-4, GPT-3.5)
- Google Gemini
- Anthropic Claude
- HuggingFace models
- Local LLMs (LLaMA, Mistral)

**Role:**

- Generate text
- Answer questions
- Reason over data

**Langchain provides a common interface, so model switching is easy**

**Provide Structure – Prompt Templates**

**Problem without prompt templates:**

- Hard-Coded prompts are:
  - Messy
  - Not reusable
  - Error-prone

**What Prompt Templates do:**

- Create structured prompts
- Insert variables dynamically
- Make prompts reusable

**Benefits**

- Clean code
- Dynamic inputs

- Easy debugging

**Divide Data into Multiple Parts – Chunking**

**Chunking (Text Splitting)**

- Chunking is the process of dividing long text into smaller, overlapping chunks.

- Key concepts:

  - Chunk size: Number of characters or tokens
  - Chunk overlap: Shared content between chunks to maintain context
- Common text splitters:

  - RecursiveCharacterTextSplitter
  - CharacterTextSplitter
  - TokenTextSplitter

**Why chunking is important?**

- LLMs have token limits
- Long documents cannot be passed directly

**Chunking does:**

- Splits large documents into small pieces
- Keeps overlap for context continuity
- Proper chunking improves retrieval accuracy and ensures important context is not lost.

**Exanple**

- Chunk size: 500 tokens
- Overlap: 50 tokens

**Result**

- Document → chunks → embeddings **this is mandatory for RAG systems**

**Store Embeddings – Vector Database**

**What are embeddings?**

- Embeddings = numerical representation of text meaning
- Embeddings convert text into numerical vectors that represent semantic meaning.
- Similar texts produce similar vectors.

**Why embeddings are important:**

- Enable similarity search

- Allow semantic understanding instead of keyword matching

**Popular embedding models:**

- OpenAI Embeddings
- HuggingFace Sentence Transformers
- Google Gemini Embeddings
- Each chunk is converted into an embedding vector and stored for future retrieval.
  **Why store embeddings?**
- Fast similarity search
- Semantic retrieval
- Context-aware answers

**Vector Databases:**

- Vector databases store embeddings efficiently and allow fast similarity search.
- Common vector databases: - FAISS - Chroma - Pinecone - Weaviate - Milvus

**What they store:**

- Embedding vectors

- Original text chunks

- Metadata **Role in RAG:**

- User question → embedding → similarity search → relevant chunks

**Save Memory – Memory in LangChain**

**Problem:**

- LLMs are stateless
- They forget previous conversation

**Memory solves:**

- Conversation history
- Context continuity
- Personalized responses

**Types of Memory:**

- ConversationBufferMemory
- ConversationSummaryMemory
- VectorStoreMemory

**Use cases:**

- Chatbots
- Assistants
- Multi-turn conversations

**Retrieval Process – Generate the Answer**

**This is the heart of RAG**

**Flow:**

- User asks a question
- Question converted to embedding
- Vector DB retrieves most relevant chunks
- Chunks + question passed to LLM
- LLM generates grounded answer

**Benefit**

- Answers are based on your data
- Reduces hallucinations
- More accurate responses

**Agents Generation**

**What are Agents?**

- Agents are decision-making LLMs.
- They can:
  - Decide which tool to use
  - Decide what step to take next
  - Perform multi-step reasoning

**Example**

- User asks:

  -"Analyze this PDF and summarize key points"
- Agent decides:

  - Load document
  - Chunk data
  - Retrieve relevant parts
  - Summarize

**Agents = LLM + tools + reasoning loop**

**Chains – Connecting Everything**

**What are Chains?**

- Chains connect: - Prompt → Model → Output **Advanced Chains:**
- RetrievalQAChain
- ConversationalRetrievalChain
- SequentialChain

**Why chains are powerful:**

- Automate workflows
- Connect prompts to Python functions
- Build end-to-end pipelines

**Final Output and Real-World Applications**

- The final output is an accurate, context-aware, and reliable answer generated using retrieved knowledge.
- Real-world applications:
    - Document Q&A systems

    - Chatbots

    - Knowledge assistants

    - Legal and medical research tools

    - Customer support automation

**Conclusion**

- LangChain transforms simple LLMs into powerful AI systems by integrating documents, embeddings, retrieval, memory, tools, and reasoning. -It is a foundational framework for building next-generation AI applications.

**End-to-End LangChain Flow (Summary)**

Document Loader ↓ Chunking ↓ Embeddings ↓ Vector Database ↓ Retriever ↓ Prompt Template ↓ LLM ↓ Memory ↓ Final Answer

```
In [ ]:
```