

# LangChain Expert Learning Guide

| For students aspiring to master LangChain development

---

## 1. Summary: What is LangChain?

**LangChain** is an open-source framework that simplifies building applications with Large Language Models (LLMs) by providing a standardized interface for connecting LLMs with external data sources, tools, and memory systems. It abstracts away the complexity of managing prompts, model interactions, and multi-step workflows, allowing developers to focus on application logic rather than LLM plumbing. Whether you're building chatbots, RAG systems, or autonomous agents, LangChain provides battle-tested components that work seamlessly together. The framework uses a composable architecture where components can be chained together using **LangChain Expression Language (LCEL)**, making it the industry standard for LLM application development alongside tools like LlamaIndex and Anthropic's SDK.

### Who should use it?

- Developers building production LLM applications
  - Teams needing rapid prototyping of AI features
  - Projects requiring RAG (Retrieval-Augmented Generation)
  - Anyone building agents or complex multi-step LLM workflows
- 

## 2. Core Components Deep Dive

### Component 1: Models (LLMs, Chat Models, Embeddings)

Aspect	Details
Name	Models / Language Models
Description	LangChain's model interface provides a unified abstraction for different LLM providers (OpenAI, Anthropic, Ollama, local models, etc.). It handles API calls, token management, and response parsing. Embeddings are a special model type that convert text into dense numerical vectors for semantic search and similarity.
Use case	When you need to call an LLM, generate embeddings for vector search, or support multiple LLM providers interchangeably.

---

Aspect	Details
Key classes	<code>ChatOpenAI</code> , <code>Ollama</code> , <code>LLMChain</code> (legacy), <code>BaseLanguageModel</code> , <code>Embeddings</code> , <code>OpenAIEmbeddings</code>

Code example:

```
python

from langchain_openai import ChatOpenAI
from langchain_ollama import OllamaEmbeddings

# Chat model for text generation
llm = ChatOpenAI(model="gpt-4", temperature=0.7)
response = llm.invoke("What is LangChain?")

# Embeddings for semantic search
embeddings = OllamaEmbeddings(model="nomic-embed-text")
vector = embeddings.embed_query("semantic search example")
```

Component 2: Prompts (Templates, Few-shot, Selectors)

Aspect	Details
Name	Prompts / Prompt Templates
Description	Prompt templates manage how you format instructions and context for LLMs. They use placeholders for dynamic values, making prompts reusable and maintainable. Few-shot templates include examples to improve model outputs through in-context learning. Selectors dynamically choose relevant examples based on query similarity.
Use case	Building consistent prompt formatting, implementing few-shot learning, managing complex multi-variable prompts, and maintaining prompt consistency across applications.
Key classes	<code>PromptTemplate</code> , <code>ChatPromptTemplate</code> , <code>FewShotPromptTemplate</code> , <code>FewShotChatMessagePromptTemplate</code> , <code>PromptSelector</code>

Code example:

```
python
```

```

from langchain_core.prompts import ChatPromptTemplate, FewShotChatMessagePromptTemplate

# Simple prompt template
template = ChatPromptTemplate.from_messages([
    ("system", "You are a helpful assistant."),
    ("user", "{user_input}")
])

# Few-shot prompt with examples
examples = [
    {"input": "Happy", "output": "😊"},
    {"input": "Sad", "output": "😞"}
]
few_shot = FewShotChatMessagePromptTemplate(
    examples=examples,
    example_prompt=ChatPromptTemplate.from_messages([
        ("user", "{input}"),
        ("assistant", "{output}")
    ])
)

```

### Component 3: Chains (LCEL, Legacy Chains)

Aspect	Details
<b>Name</b>	Chains / LCEL (LangChain Expression Language)
<b>Description</b>	Chains compose multiple components (LLMs, prompts, parsers, tools) into sequential workflows. <b>LCEL</b> is the modern, recommended way to build chains using the `
<b>Use case</b>	Orchestrating multi-step workflows, chaining prompts → LLM → output parser, building RAG pipelines, creating complex reasoning chains.
<b>Key classes</b>	<code>LCEL</code> (implicit via `

### Code example:

```
python
```

```
from langchain_core.prompts import ChatPromptTemplate
from langchain_openai import ChatOpenAI
from langchain_core.output_parsers import StrOutputParser

# Modern LCEL approach (recommended)
prompt = ChatPromptTemplate.from_template("Explain {topic} in 2 sentences.")
llm = ChatOpenAI(model="gpt-4")
parser = StrOutputParser()

chain = prompt | llm | parser
result = chain.invoke({"topic": "Machine Learning"})
```

Component 4: Memory (Conversation History)

Aspect	Details
Name	Memory
Description	Memory manages conversation history and context. LangChain provides multiple memory types: <code>ConversationBufferMemory</code> (stores all messages), <code>ConversationSummaryMemory</code> (summarizes old messages), <code>ConversationTokenBufferMemory</code> (keeps within token limit), and specialized variants. Memory is essential for stateful conversations and multi-turn interactions.
Use case	Building chatbots with conversation context, maintaining user state across interactions, implementing conversation summarization for long chats.
Key classes	<code>ConversationBufferMemory</code> , <code>ConversationSummaryMemory</code> , <code>ConversationTokenBufferMemory</code> , <code>ConversationKGMemory</code> , <code>PostgresChatMessageHistory</code>

Code example:

```
python
```

```
from langchain.memory import ConversationBufferMemory
from langchain_openai import ChatOpenAI
from langchain_core.prompts import ChatPromptTemplate

memory = ConversationBufferMemory(return_messages=True)
llm = ChatOpenAI(model="gpt-4")

# Add messages to memory
memory.chat_memory.add_user_message("Hi, I'm learning LangChain")
memory.chat_memory.add_ai_message("Great! I can help you learn.")

# Retrieve conversation history
history = memory.load_memory_variables({})
print(history["history"])
```

Component 5: Retrieval (Document Loaders, Vector Stores, Retrievers)

Aspect	Details
Name	Retrieval / RAG (Retrieval-Augmented Generation)
Description	Retrieval components handle data ingestion and semantic search. Document loaders ingest various file types (PDFs, CSVs, web pages). Vector stores (Pinecone, Weaviate, FAISS, Chroma) store embeddings for fast similarity search. Retrievers wrap vector stores with advanced search logic. Together, they enable RAG—augmenting LLM responses with relevant external documents.
Use case	Building RAG chatbots over custom documents, implementing semantic search, creating knowledge bases, integrating external data into LLM responses.
Key classes	TextLoader, PyPDFLoader, WebBaseLoader, FAISS, Chroma, PineconeVectorStore, VectorStoreRetriever, BM25Retriever, EnsembleRetriever

Code example:

```
python
```

```

from langchain_community.document_loaders import TextLoader
from langchain_text_splitters import RecursiveCharacterTextSplitter
from langchain_community.vectorstores import FAISS
from langchain_openai import OpenAIEmbeddings

# Load and chunk documents
loader = TextLoader("document.txt")
docs = loader.load()
splitter = RecursiveCharacterTextSplitter(chunk_size=500)
chunks = splitter.split_documents(docs)

# Create vector store
embeddings = OpenAIEmbeddings()
vector_store = FAISS.from_documents(chunks, embeddings)

# Retrieve similar documents
retriever = vector_store.as_retriever()
relevant_docs = retriever.invoke("How does LangChain work?")

```

## Component 6: Agents (Tools, Agent Types)

Aspect	Details
Name	Agents / Tool Use
Description	Agents are LLM-powered decision makers that use tools to accomplish tasks. Unlike chains (pre-defined workflows), agents decide dynamically which tools to call and in what order using reasoning loops. Tools are functions the agent can invoke (search, calculator, database queries). Modern agents use <b>tool_choice</b> to force tool use and structured outputs.
Use case	Building autonomous AI assistants, implementing ReAct (Reasoning + Acting) patterns, creating systems that decide which tools to use, building code execution engines.
Key classes	<code>Tool</code> , <code>@tool</code> decorator, <code>AgentExecutor</code> , <code>create_react_agent</code> , <code>ToolMessage</code> , <code>BaseTool</code>

## Code example:

```
python
```

```

from langchain.tools import tool
from langchain_openai import ChatOpenAI
from langchain.agents import create_react_agent, AgentExecutor
from langchain_core.prompts import ChatPromptTemplate

@tool
def multiply(a: int, b: int) -> int:
    """Multiply two numbers."""
    return a * b

@tool
def add(a: int, b: int) -> int:
    """Add two numbers."""
    return a + b

tools = [multiply, add]
llm = ChatOpenAI(model="gpt-4")

# Create and execute agent
agent = create_react_agent(llm, tools, ChatPromptTemplate.from_messages([
    ("system", "You are a helpful math assistant."),
    ("user", "{input}")
]))
executor = AgentExecutor(agent=agent, tools=tools)
result = executor.invoke({"input": "What is 5 times 3, plus 2?"})

```

## Component 7: Callbacks (Logging, Streaming)

Aspect	Details
Name	Callbacks / Event Handlers
Description	Callbacks are hooks into the execution lifecycle that enable logging, monitoring, streaming, and debugging. <code>BaseCallbackHandler</code> is the base class for implementing custom handlers. Common use cases include streaming token outputs in real-time, logging LLM calls for monitoring, and implementing cost tracking.
Use case	Streaming LLM responses to users, logging API calls for debugging, tracking token usage and costs, implementing real-time progress bars, custom monitoring integrations.

Aspect	Details
Key classes	<code>BaseCallbackHandler</code> , <code>StreamingStdOutCallbackHandler</code> , <code>LangChainTracer</code> , <code>CustomCallbackHandler</code>

Code example:

```
python

from langchain_openai import ChatOpenAI
from langchain_core.callbacks import StreamingStdOutCallbackHandler
from langchain_core.prompts import ChatPromptTemplate

# Streaming callback - prints tokens as they arrive
streaming_handler = StreamingStdOutCallbackHandler()

prompt = ChatPromptTemplate.from_template("Explain {topic}")
llm = ChatOpenAI(model="gpt-4", streaming=True)

chain = prompt | llm
result = chain.invoke(
    {"topic": "quantum computing"},
    config={"callbacks": [streaming_handler]}
)
```

Component 8: Output Parsers

Aspect	Details
Name	Output Parsers
Description	Output parsers transform raw LLM text responses into structured Python objects. Types include: <code>StrOutputParser</code> (string), <code>JsonOutputParser</code> (JSON), <code>PydanticOutputParser</code> (validated Pydantic models), <code>XMLOutputParser</code> (XML structures). With newer LLM models supporting structured outputs natively, parsers are increasingly handling formatting logic.
Use case	Extracting structured data from LLM outputs, validating model responses, parsing JSON/XML/CSV from LLM generations, implementing data validation.
Key classes	<code>StrOutputParser</code> , <code>JsonOutputParser</code> , <code>PydanticOutputParser</code> , <code>XMLOutputParser</code> , <code>CommaSeparatedListOutputParser</code> , <code>OutputFixingParser</code>



Code example:

```
python

from langchain_openai import ChatOpenAI
from langchain_core.output_parsers import JsonOutputParser
from pydantic import BaseModel, Field

class Person(BaseModel):
    name: str = Field(description="Person's name")
    age: int = Field(description="Person's age")

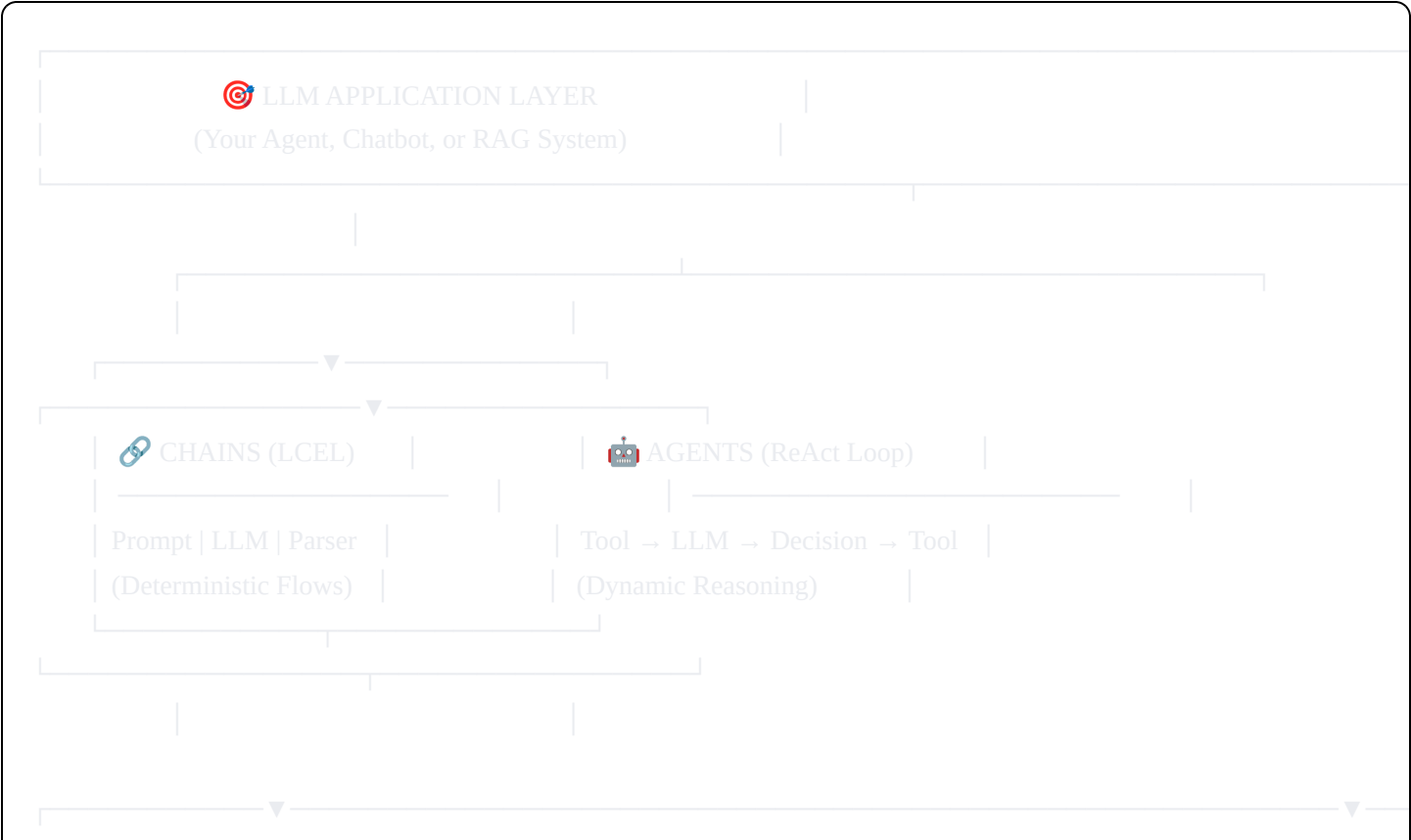
parser = JsonOutputParser(pydantic_object=Person)
llm = ChatOpenAI(model="gpt-4")

prompt = "Extract person info: John is 30 years old"
chain = llm | parser

result = chain.invoke(prompt) # Returns Person(name="John", age=30)
```

3. 🏗️ Architecture Overview

Component Relationship Diagram





## CORE COMPONENTS



MODELS



PROMPTS



MEMORY



TOOLS

• LLMs (OpenAI)

• Templates

• Conv Buf

• Custom

• Chat Models

• Few-shot

• Summary

• Search

• Embeddings

• Selectors

• KG Memory

• Calculator



RETRIEVAL LAYER (RAG)

Document

Vector

Retrievers

Loaders

Stores

(BM25, Ens)

• PDF, Web

(FAISS, etc)

• CSV, JSON



OUTPUT & CALLBACKS LAYER

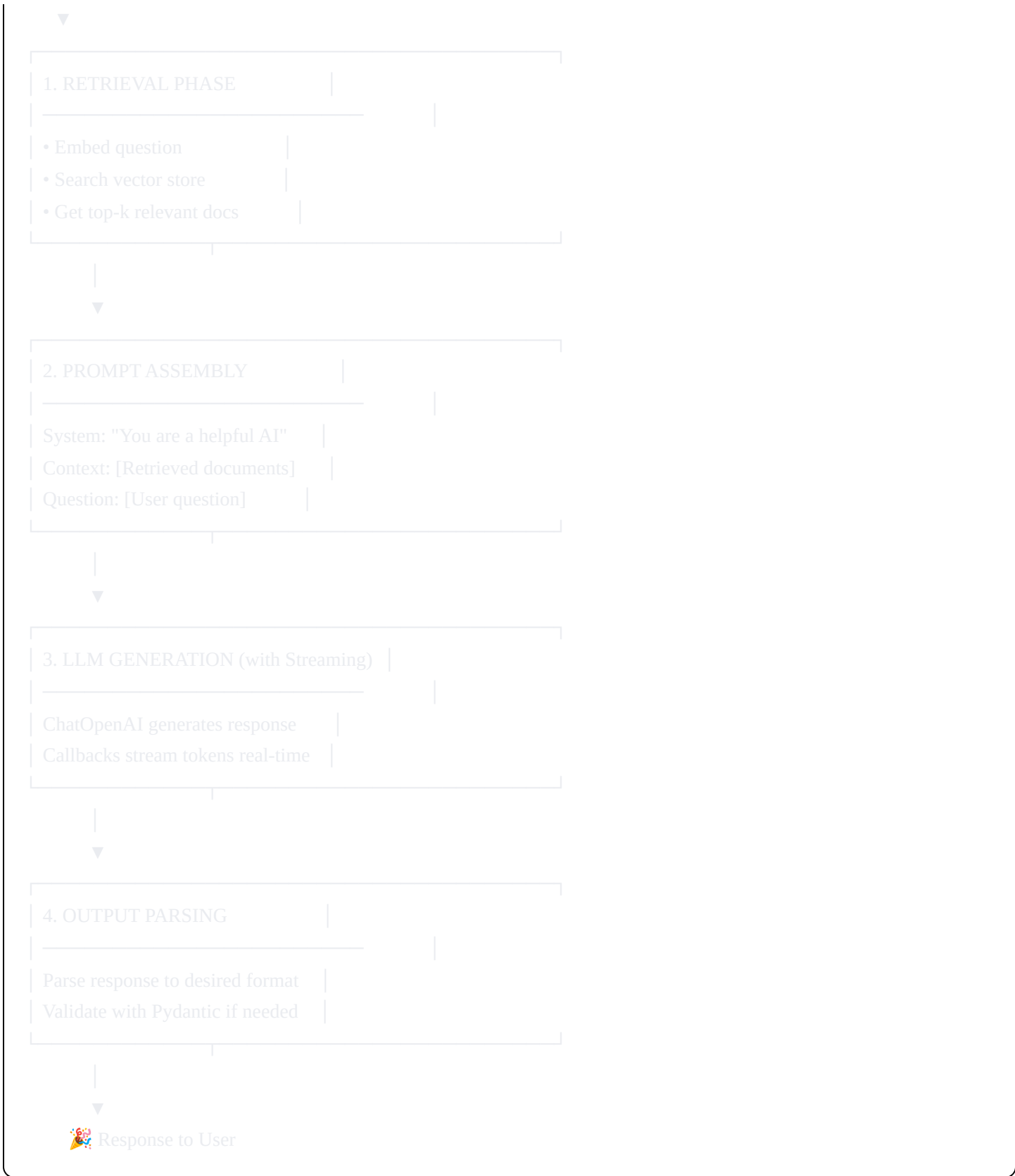
• Output Parsers (JSON, Pydantic, XML)

• Streaming Callbacks (Real-time output)

• Logging & Monitoring (Cost, tokens, latency)

## Data Flow Example: RAG Chatbot

User Question



#### 4. 🚀 Quick Start Example: Building Your First Chain

##### Setup

```
bash
```

pip [install](#) langchain langchain-openai python-dotenv

## Complete Example: Q&A Over Documents

```
python
```

```
import os
from dotenv import load_dotenv
from langchain_community.document_loaders import TextLoader
from langchain_text_splitters import RecursiveCharacterTextSplitter
from langchain_community.vectorstores import FAISS
from langchain_openai import OpenAIEmbeddings, ChatOpenAI
from langchain_core.prompts import ChatPromptTemplate
from langchain_core.output_parsers import StrOutputParser

# 1. Load and chunk documents
loader = TextLoader("my_document.txt")
docs = loader.load()
splitter = RecursiveCharacterTextSplitter(chunk_size=1000, chunk_overlap=200)
chunks = splitter.split_documents(docs)

# 2. Create vector store
embeddings = OpenAIEmbeddings()
vector_store = FAISS.from_documents(chunks, embeddings)
retriever = vector_store.as_retriever(search_kwargs={"k": 3})

# 3. Build RAG chain with LCEL
prompt = ChatPromptTemplate.from_template("""
You are a helpful assistant. Answer the question based on the context provided.

Context:
{context}

Question: {question}

Answer:
""")

llm = ChatOpenAI(model="gpt-4", temperature=0.7)
parser = StrOutputParser()

chain = (
    {
        "context": retriever | (lambda docs: "\n".join([d.page_content for d in docs])),
        "question": lambda x: x["question"]
    }
    | prompt
    | llm
    | parser
```

)

# 4. Run the chain

```
result = chain.invoke({"question": "What is LangChain?"})
```

```
print(result)
```

---

## 5. Learning Path Recommendations

### Phase 1: Fundamentals (Week 1-2)

- ☐ Understand LLM basics and tokens
- ☐ Explore basic prompts and templates
- ☐ Build simple LLM chains with LCEL
- ☐ Experiment with different LLMs (OpenAI, Ollama)

### Phase 2: Intermediate (Week 3-4)

- ☐ Implement memory for conversations
- ☐ Learn vector stores and embeddings
- ☐ Build your first RAG system
- ☐ Work with output parsers

### Phase 3: Advanced (Week 5-6)

- ☐ Create agents with tools
- ☐ Implement custom callbacks
- ☐ Build multi-step workflows
- ☐ Deploy to production

### Phase 4: Mastery (Week 7+)

- ☐ Implement advanced retrieval strategies
  - ☐ Create custom agent types
  - ☐ Optimize for cost and latency
  - ☐ Build production monitoring
-

## 6. Common Patterns & Best Practices

### Pattern 1: Simple Q&A Chain

```
python  
  
prompt | llm | output_parser
```

### Pattern 2: RAG Chain

```
python  
  
retriever | format_context | prompt | llm | parser
```

### Pattern 3: Agent Loop

```
python  
  
prompt | llm | parse_tool_calls | execute_tool | (loop back)
```

### Pattern 4: Streaming Response

```
python  
  
chain.stream(input, config={"callbacks": [streaming_handler]})
```

### Best Practices

1. **Use LCEL over legacy chains** - Cleaner, more composable, better type support
  2. **Always add error handling** - LLM calls can fail; implement retries and fallbacks
  3. **Monitor costs** - Track token usage, especially with expensive models
  4. **Version your prompts** - Treat prompts as code; use version control
  5. **Test with local models first** - Use Ollama to reduce API costs during development
  6. **Implement proper memory management** - Long conversations need summarization
  7. **Cache embeddings** - Reuse embeddings to reduce API calls
  8. **Use structured outputs** - Force JSON/Pydantic validation when possible
- 

## 7. Resources for Deep Learning

- **Official Docs:** <https://python.langchain.com/docs/>

- **LCEL Guide:** [https://python.langchain.com/docs/expression\\_language/](https://python.langchain.com/docs/expression_language/)
  - **LangSmith** (Debugging): <https://smith.langchain.com/>
  - **GitHub:** <https://github.com/langchain-ai/langchain>
  - **Community Discord:** Join for real-time help
  - **YouTube:** Harrison Chase (creator) explains architecture
- 

## 8. ? Common Mistakes to Avoid

Mistake	Why It's Bad	Solution
Not using LCEL	Harder to read, maintain, and type-safe	Always use <code> &gt;</code> operator for chaining
Hardcoding API keys	Security risk, environment pollution	Use <code>.env</code> files and environment variables
No output parsing	Strings from LLMs are unpredictable	Always parse to structured objects
Ignoring memory limits	Tokens = costs; long convos explode tokens	Implement token buffer or summarization
Not testing locally	Expensive API calls during development	Use Ollama or mock responses
Missing error handling	Production systems crash silently	Wrap calls in try-except, implement retries
Prompt versioning chaos	Hard to track which prompt works best	Version prompts like code in git


---

## 9. 🎯 Your Next Steps

1. **Clone this repo:** `git clone https://github.com/langchain-ai/langchain`
  2. **Install locally:** `pip install -e .`
  3. **Run examples:** Explore `examples/` directory
  4. **Build a mini-project:** RAG over your favorite PDF
  5. **Join the community:** Contribute to LangChain or build with it
  6. **Share your learnings:** Write a blog post about what you built
- 

**Happy learning! You're now equipped with the knowledge to build production-grade LLM applications. Start small, iterate fast, and scale gradually.**



 *Built with clarity for learning developers, by someone who learned the hard way.*