

INTRODUCTION TO LANGCHAIN

Yogesh Haribhau Kulkarni

Outline

① INTRODUCTION

② FRAMEWORK

③ WHAT'S NEW

④ CONCLUSIONS

⑤ REFERENCES

Introduction to LangChain

What is LangChain?

LangChain: A comprehensive framework for building LLM-powered applications

- ▶ **Core Purpose:** Simplify development of applications using LLMs
- ▶ **Key Solutions:**
 - ▶ **RAG:** Connect language models to external data sources
 - ▶ **Agentic:** Allow language models to interact with their environment
- ▶ **Main Features:**
 - ▶ Generic interface to various foundation models
 - ▶ Advanced prompt management framework
- ▶ **Availability:** Python and JavaScript libraries
- ▶ **Open Source:** MIT License, created by Harrison Chase
- ▶ **Repository:** <https://github.com/langchain-ai/langchain>

(Ref: Getting Started with LangChain: A Beginner's Guide to Building LLM-Powered Applications)

Why You Need LangChain

- ▶ LLMs alone lack context, memory, and retrieval abilities.
- ▶ Chatbots need to manage chat history and company knowledge bases.
- ▶ Storing, retrieving, and reasoning over data is complex manually.
- ▶ LangChain acts as an abstraction layer for building AI agents.
- ▶ It integrates LLMs, memory, and tools seamlessly.
- ▶ Enables vendor flexibility, easy switch between Open/Close models.
- ▶ Reduces code complexity and accelerates AI app development.
- ▶ Simplifies connecting models, databases, and APIs into a single framework.

LLMs vs Agentic Software

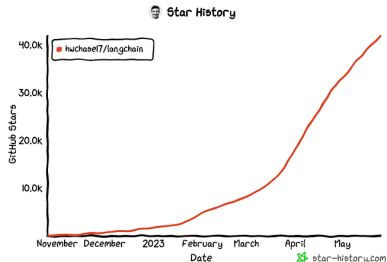
- ▶ LLMs are static, respond from training data without awareness.
- ▶ Agents have autonomy, memory, and tool use for dynamic actions.
- ▶ Example: Refund query, agent retrieves policy, product, and chat history.
- ▶ Agents can access vector databases and retrieve company knowledge.
- ▶ LangChain enables memory persistence across user conversations.
- ▶ Traditional software follows fixed logic; agents make adaptive decisions.
- ▶ Agentic software uses modular components (LLMs, memory, retrievers).
- ▶ LangChain provides prebuilt tools for these agentic capabilities.

LangChain Core Components and Labs

- ▶ Core modules: LLM connectors, memory, embeddings, vector stores.
- ▶ Simplifies LLM API setup, e.g., ChatOpenAI() or ChatGroq()
- ▶ Supports databases like Chroma or Pinecone for knowledge retrieval.
- ▶ Prompt templates manage dynamic, reusable prompts and chat context.
- ▶ LCEL (LangChain Expression Language) composes chains without legacy LLMChain wrappers
- ▶ Enables async, streaming, and batch workflows with type safety.

So, Why LangChain?

- ▶ **RAG Applications:** Build Retrieval Augmented Generation apps with external data
- ▶ **Agent Systems:** Create intelligent agents with tool access
- ▶ **Production Ready:** Logging, callbacks, monitoring built-in
- ▶ **Model Agnostic:** Switch between LLM providers easily
- ▶ **Modular Design:** Compose components as needed



(Ref: LangChain tutorial: Build an LLM-powered app)

LangChain Community & Ecosystem

► **Community Growth:**

- Over 80,000+ GitHub stars (as of 2024)
- 2,000+ contributors
- Millions of monthly downloads
- Active Discord, Twitter/X presence

► **License:** MIT License - freely modifiable and commercial-friendly

► **Ecosystem Components:**

- **LangChain Core:** Foundation library
- **LangGraph:** Stateful, multi-actor applications
- **LangServe:** Deploy as REST APIs
- **LangSmith:** Debugging, testing, monitoring

(Ref: What is Langchain and why should I care as a developer? - Logan Kilpatrick)

Installation & Setup

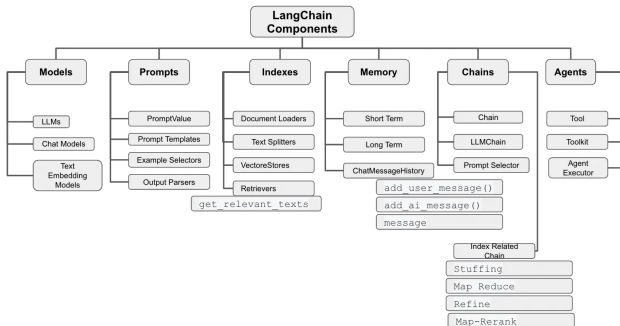
Prerequisites:

- ▶ Python version ≥ 3.9 and < 4.0
- ▶ Modern Installation (2024)
- ▶ (Ref: Official LangChain Documentation 2024)

:

```
1 # Core packages
  pip install langchain langchain-core
3
4 # Provider-specific packages for Groq
5 pip install langchain-groq
6
7 # Alternative providers
8 # pip install langchain-anthropic langchain-google-genai
9
10 # Community integrations for embeddings
11 pip install langchain-community langchain-huggingface sentence-transformers
12
13 # Optional components
14 pip install chromadb # Vector store
15 pip install langchain-text-splitters # Document processing
16
17 //API Keys Setup:
18
19 import os
  os.environ["GROQ_API_KEY"] = "your-groq-api-key-here"
21 # Or use .env file with python-dotenv
```

Core Components Overview



- **Models**: LLMs, Chat Models, Embeddings
- **Prompts**: Template management and optimization
- **Chains/LCEL**: Compose components with pipes
- **Memory**: Conversation state management
- **Retrievers**: Access external data
- **Agents**: Dynamic tool selection and execution

Modern LangChain: LCEL (LangChain Expression Language)

What is LCEL?

- ▶ Modern, declarative way to build chains (introduced 2023)
- ▶ Uses pipe operator `|` to chain components
- ▶ Replaces old `LLMChain` pattern
- ▶ Built-in streaming, async, and batch support

Key Benefits:

- ▶ **Simplicity:** More readable and concise
- ▶ **Streaming:** Built-in streaming by default
- ▶ **Async:** Native async/await support
- ▶ **Observability:** Better tracing and debugging
- ▶ **Fallbacks:** Easy error handling and retries

(Ref: LangChain LCEL Documentation 2024)

LCEL: Basic Example

```
1 // Old Pattern (Deprecated):
2 from langchain.llms import OpenAI
3 from langchain.chains import LLMChain
4
5 llm = OpenAI()
6 chain = LLMChain(llm=llm, prompt=prompt)
7 result = chain.run("input")
8
9 // Modern LCEL Pattern with Groq:
10
11 from langchain_groq import ChatGroq
12 from langchain_core.prompts import ChatPromptTemplate
13 from langchain_core.output_parsers import StrOutputParser
14
15 # Set up Groq model, e.g., Gemma or Llama 3
16 llm = ChatGroq(model_name="gemma-7b-it")
17 prompt = ChatPromptTemplate.from_template("Tell me about {topic}")
18 output_parser = StrOutputParser()
19
20 # Build chain with pipe operator
21 chain = prompt | llm | output_parser
22
23 # Invoke the chain
24 result = chain.invoke({"topic": "LangChain"})
```

LCEL: Advanced Features

```
2 // Streaming Support:
3
4 chain = prompt | llm | output_parser
5
6 # Stream tokens as they're generated
7 for chunk in chain.stream({"topic": "AI"}):
8     print(chunk, end="", flush=True)
9
10 // Async Execution:
11 # Async invocation
12 result = await chain.ainvoke({"topic": "AI"})
13
14 # Async streaming
15 async for chunk in chain.astream({"topic": "AI"}):
16     print(chunk, end="", flush=True)
17
18 // Batch Processing:
19 # Process multiple inputs in parallel
20 results = chain.batch([
21     {"topic": "AI"},
22     {"topic": "ML"},
23     {"topic": "LangChain"}
24 ])
```

LCEL: Complex Chains

```
1 // Multi-step Chain with RunnablePassthrough:
2 from langchain_core.runnables import RunnablePassthrough
3
4 chain = (
5     {"context": retriever, "question": RunnablePassthrough()}
6     | prompt
7     | llm
8     | output_parser
9 )
10
11 result = chain.invoke("What is LangChain?")
12
13 // Parallel Execution with RunnableParallel:
14
15 from langchain_core.runnables import RunnableParallel
16
17 chain = RunnableParallel(
18     summary=prompt1 | llm | output_parser,
19     keywords=prompt2 | llm | output_parser
20 )
21
22 result = chain.invoke({"text": "Long document..."})
23 # Returns: {"summary": "...", "keywords": "..."}

```

Model Integration: Modern Approach

```
1 Updated Import Structure:
3 # OpenAI models
4 from langchain.openai import ChatOpenAI, OpenAIEmbeddings
5
6 # Hugging Face models
7 from langchain.huggingface import HuggingFaceEndpoint
8
9 # Google models (Gemini)
10 from langchain.google.genai import ChatGoogleGenerativeAI
11
12 # Anthropic models
13 from langchain.anthropic import ChatAnthropic
14
15 // Example Usage:
16 from langchain.groq import ChatGroq
17 from langchain_core.messages import HumanMessage, SystemMessage
18
19 # Initialize Groq with Gemma
20 chat = ChatGroq(model="gemma2-9b-it", temperature=0.7)
21
22 messages = [
23     SystemMessage(content="You are a helpful assistant"),
24     HumanMessage(content="Explain LangChain in 2 sentences")
25 ]
26
27 response = chat.invoke(messages)
28 print(response.content)
```


Complete RAG Application: A Quick Start Example

```
1 from langchain_groq import ChatGroq
2 from langchain_community.embeddings import HuggingFaceEmbeddings
3 from langchain_community.document_loaders import WebBaseLoader
4 from langchain_text_splitters import RecursiveCharacterTextSplitter
5 from langchain_community.vectorstores import Chroma
6 from langchain_core.prompts import ChatPromptTemplate
7 from langchain_core.output_parsers import StrOutputParser
8 from langchain_core.runnables import RunnablePassthrough

10 # Load and process documents
11 loader = WebBaseLoader("https://example.com/doc")
12 docs = loader.load()
13 splitter = RecursiveCharacterTextSplitter(chunk_size=1000)
14 splits = splitter.split_documents(docs)

16 embeddings = HuggingFaceEmbeddings()
17 vectorstore = Chroma.from_documents(splits, embeddings)
18 retriever = vectorstore.as_retriever()

20 prompt = ChatPromptTemplate.from_template("""
21 Answer based on context: {context}
22 Question: {question} """)

24 model = ChatGroq(model_name="llama3-8b-8192")

26 # RunnablePassthrough => "context": retriever("What is the main topic?"),
chain = (
28     {"context": retriever, "question": RunnablePassthrough()}
29     | prompt | model | StrOutputParser()
30 )

response = chain.invoke("What is the main topic?")
```

Key Concepts: Chains vs Agents

► Chains (LCEL):

- Predetermined sequence of operations
- Composed with pipe operator: `prompt | llm | parser`
- Fixed execution path
- Best for: Structured, predictable workflows

► Agents:

- Dynamic decision-making with LLM reasoning
- Choose tools based on input
- Adaptive execution path
- Best for: Complex, unpredictable scenarios

When to use what?

- Use **Chains/LCEL** when: Steps are known, workflow is fixed
- Use **Agents** when: Need dynamic tool selection, multi-step reasoning

(Ref: Superpower LLMs with Conversational Agents)

LangChain Use Cases

- ▶ **Retrieval Augmented Generation (RAG):**

- ▶ Document Q&A systems
- ▶ Knowledge base search
- ▶ Semantic search applications

- ▶ **Conversational AI:**

- ▶ Chatbots with memory
- ▶ Customer support agents
- ▶ Personal assistants

- ▶ **Data Analysis:**

- ▶ SQL query generation
- ▶ Tabular data Q&A
- ▶ Report generation

- ▶ **Autonomous Agents:**

- ▶ Web scraping and research
- ▶ Multi-tool workflows
- ▶ API integration

(Ref: LangChain Use Cases Documentation)

LangChain Framework Components

Framework Architecture

Chains: The core of LangChain. Components (and even other chains) can be stringed together to create *chains*.

Prompt templates: Prompt templates are templates for different types of prompts. Like “chatbot” style templates, ELI5 question-answering, etc

LLMs: Large language models like GPT-3, BLOOM, etc

Indexing Utils: Ways to interact with specific data (embeddings, vectorstores, document loaders)

Tools: Ways to interact with the outside world (search, calculators, etc)

Agents: Agents use LLMs to decide what actions should be taken. Tools like web search or calculators can be used, and all are packaged into a logical loop of operations.

Memory: Short-term memory, long-term memory.

Core Building Blocks:

- ▶ **Models:** LLMs, Chat Models, Embeddings
- ▶ **Prompts:** Dynamic template management
- ▶ **Output Parsers:** Structured output extraction
- ▶ **Retrievers:** Document and data access
- ▶ **Memory:** Conversation state persistence
- ▶ **Agents & Tools:** Dynamic reasoning and actions

(Ref: Building the Future with LLMs, LangChain, & Pinecone)

Models

Models in LangChain

Three Types of Models:

- ▶ **LLMs (Large Language Models):**
 - ▶ Input: String (prompt)
 - ▶ Output: String (completion)
 - ▶ Examples: GPT-4, Claude, Gemma, Llama 3, Mixtral.
 - ▶ Use case: Text generation, completion
- ▶ **Chat Models:**
 - ▶ Input: List of messages
 - ▶ Output: Chat message
 - ▶ Examples: ChatGPT, Claude Chat, ChatGroq
 - ▶ Use case: Conversational AI
- ▶ **Embedding Models:**
 - ▶ Input: Text
 - ▶ Output: Vector (list of floats)
 - ▶ Examples: OpenAI Embeddings, HuggingFace, Sentence Transformers
 - ▶ Use case: Semantic search, similarity

Model Integration: Modern Syntax

```
1 from langchain_groq import ChatGroq
2 from langchain_community.embeddings import HuggingFaceEmbeddings
3 from langchain_core.messages import HumanMessage, SystemMessage
4 from langchain_core.output_parsers import StrOutputParser
5
6 # Chat Model with Groq
7 chat = ChatGroq(model_name="llama3-8b-8192", temperature=0.7)
8
9 messages = [
10     SystemMessage(content="You are a helpful assistant"),
11     HumanMessage(content="Explain quantum computing briefly")
12 ]
13 response = chat.invoke(messages)
14
15 # Embedding Model (from Hugging Face)
16 embeddings = HuggingFaceEmbeddings(model_name="all-MiniLM-L6-v2")
17 vector = embeddings.embed_query("Machine learning is...")
18
19 # Using with LCEL
20 from langchain_core.prompts import ChatPromptTemplate
21
22 prompt = ChatPromptTemplate.from_template("Explain {topic}")
23 chain = prompt | chat | StrOutputParser()
24 result = chain.invoke({"topic": "blockchain"})
```


Alternative Model Providers

```
2 // Groq (Recommended for Fast Inference):
  from langchain_groq import ChatGroq

4 # Gemma 2 — Very fast, efficient
  llm = ChatGroq(model_name="gemma2-9b-it", temperature=0.7)
6 response = llm.invoke("What is the Groq LPU?")

8 # Llama 3 — Balanced performance
  llm = ChatGroq(model_name="llama3-70b-8192", temperature=0.7)

10 // Hugging Face (Self-hosted):
12 from langchain_huggingface import HuggingFaceEndpoint
  llm = HuggingFaceEndpoint(
14     repo_id="mistralai/Mistral-7B-Instruct-v0.2",
     temperature=0.7
16 )

18 // Google Gemini (Multimodal):
  from langchain_google_genai import ChatGoogleGenerativeAI
20 llm = ChatGoogleGenerativeAI(model="gemini-pro", temperature=0.7)

22 // Anthropic Claude (Reasoning):
  from langchain_anthropic import ChatAnthropic
24 llm = ChatAnthropic(model="claude-3-5-sonnet-20241022")
```

Prompts

Prompts in LangChain

Prompt Management Strategies:

- ▶ **Prompt Templates:**
 - ▶ Parameterized templates with variables
 - ▶ Dynamic input insertion
 - ▶ Reusable prompt structures
- ▶ **Chat Prompt Templates:**
 - ▶ Multi-message conversations
 - ▶ System, human, AI message roles
 - ▶ Better for chat models
- ▶ **Few-Shot Prompts:**
 - ▶ Include example inputs/outputs
 - ▶ Guide model response style
 - ▶ Improve accuracy

Prompt Templates: Modern Examples

Basic Prompt Template:

```
1 from langchain_core.prompts import PromptTemplate
2
3 template = """You are a {role} assistant.
4 Task: {task}
5 Context: {context}
6 Provide a {format} response."""
7
8 prompt = PromptTemplate(
9     template=template,
10    input_variables=["role", "task", "context", "format"]
11 )
12
13 formatted = prompt.format(
14     role="helpful",
15     task="explain quantum computing",
16     context="for beginners",
17     format="simple"
18 )
19
20 // Chat Prompt Template:
21 from langchain_core.prompts import ChatPromptTemplate
22
23 prompt = ChatPromptTemplate.from_messages([
24     ("system", "You are a {role}"),
25     ("human", "{input}"),
26     ("ai", "I understand. Let me help with that."),
27     ("human", "{follow_up}")
28 ])
```

Few-Shot Prompting

```
from langchain_core.prompts import FewShotPromptTemplate, PromptTemplate

2
# Define examples
4 examples = [
    {"input": "happy", "output": "joyful, cheerful, delighted"},
    {"input": "sad", "output": "unhappy, sorrowful, dejected"},
6 ]
8
# Example template
10 example_prompt = PromptTemplate(
    input_variables=["input", "output"],
    template="Input: {input} \nSynonyms: {output}"
12 )
14
# Few-shot template
16 few_shot_prompt = FewShotPromptTemplate(
    examples=examples,
    example_prompt=example_prompt,
    prefix="Provide synonyms for the following word:",
    suffix="Input: {word} \nSynonyms:",
    input_variables=["word"]
20 )
22
24 # Use in chain with Groq
from langchain_groq import ChatGroq
from langchain_core.output_parsers import StrOutputParser

26
28 chain = few_shot_prompt | ChatGroq(model.name="gemma-7b-it") | StrOutputParser()
result = chain.invoke({"word": "angry"})
```

Memory

Memory in LangChain

Why Memory?

- ▶ LLMs are stateless by default
- ▶ Chatbots need conversation context
- ▶ Memory stores and retrieves conversation history

Memory Types (Consolidated):

- ▶ **ConversationBufferMemory:**
 - ▶ Stores entire conversation history
 - ▶ Simple but can exceed token limits
 - ▶ Best for: Short conversations
- ▶ **ConversationBufferWindowMemory:**
 - ▶ Keeps only last K interactions
 - ▶ Prevents token overflow
 - ▶ Best for: Longer conversations with recent context
- ▶ **ConversationSummaryMemory:**
 - ▶ Summarizes old messages
 - ▶ Reduces token usage
 - ▶ Best for: Very long conversations

Memory: Implementation Examples

Buffer Memory (Stores Everything):

```
1 from langchain.memory import ConversationBufferMemory
2 from langchain.groq import ChatGroq
3 from langchain_core.prompts import ChatPromptTemplate, MessagesPlaceholder
4 from langchain_core.runnables import RunnablePassthrough
5
6 memory = ConversationBufferMemory(
7     return_messages=True,
8     memory_key="history"
9 )
10
11 prompt = ChatPromptTemplate.from_messages([
12     ("system", "You are a helpful assistant"),
13     MessagesPlaceholder(variable_name="history"),
14     ("human", "{input}")
15 ])
16
17 chain = (
18     RunnablePassthrough.assign(
19         history=lambda x: memory.load_memory_variables({})["history"]
20     )
21     | prompt
22     | ChatGroq(model_name="gemma2-9b-it")
23 )
24
25 # Use the chain
26 response = chain.invoke({"input": "Hi, I'm Alice"})
27 memory.save_context({"input": "Hi, I'm Alice"}, {"output": response.content})
```


Memory: Window Memory

Benefits:

- ▶ Fixed memory footprint
- ▶ Maintains recent context
- ▶ Prevents token limit issues

Buffer Window Memory (Keeps Last K):

```
1 from langchain.memory import ConversationBufferWindowMemory
3 memory = ConversationBufferWindowMemory(
4     k=3, # Only keeps last 3 interactions
5     return_messages=True,
6     memory_key="history"
7 )
9 conversations = [ # Add conversations
10     ("Tell me about Python", "Python is a versatile language..."),
11     ("What makes it popular?", "Its simplicity and libraries..."),
12     ("Give an example", "Like NumPy for computing..."),
13     ("What about AI?", "Python excels in AI with TensorFlow...")
14 ]
16 for input_text, output_text in conversations:
17     memory.save_context({"input": input_text}, {"output": output_text})
19 # Retrieves only last 3 interactions
20 history = memory.load_memory_variables({})
21 print(f"Stored messages: {len(history['history'])}") # Will be 6 (3 pairs)
```

Document Loaders & Retrievers

Document Loading

Wide Range of Loaders:

- ▶ **Files:** PDF, Word, PowerPoint, CSV, Markdown
- ▶ **Web:** HTML, URLs, sitemaps, web scraping
- ▶ **Cloud:** S3, GCS, Google Drive, Notion
- ▶ **Databases:** SQL, MongoDB, Elasticsearch
- ▶ **Communication:** Email, Slack, Discord
- ▶ **Code:** GitHub, GitLab, Jupyter notebooks

Modern Document Processing:

```
1 from langchain_community.document_loaders import (  
2     PyPDFLoader,  
3     WebBaseLoader,  
4     TextLoader)  
5 from langchain_text_splitters import RecursiveCharacterTextSplitter  
  
7 loader = PyPDFLoader("document.pdf")  
8 documents = loader.load()  
  
9  
10 splitter = RecursiveCharacterTextSplitter(  
11     chunk_size=1000,  
12     chunk_overlap=200)  
13  
14 chunks = splitter.split_documents(documents)
```

Vector Stores & Retrievers

Popular Vector Stores (2024):

- ▶ **Chroma**: Open-source, local-first, easy setup
- ▶ **Pinecone**: Managed service, production-ready
- ▶ **Weaviate**: GraphQL API, hybrid search
- ▶ **Qdrant**: High performance, filtering support
- ▶ **LanceDB**: Serverless, embedded option

Complete Example with Chroma & HuggingFace Embeddings:

```
2 # from langchain_community.embeddings import HuggingFaceEmbeddings
3 from langchain_huggingface import HuggingFaceEmbeddings
4 from langchain_community.vectorstores import Chroma
5 from langchain.text_splitters import RecursiveCharacterTextSplitter
6 from langchain_community.document_loaders import WebBaseLoader
7
8 loader = WebBaseLoader("https://example.com/article")
9 docs = loader.load()
10
11 splitter = RecursiveCharacterTextSplitter(chunk_size=1000, chunk_overlap=200)
12 splits = splitter.split_documents(docs)
13
14 vectorstore = Chroma.from_documents(
15     documents=splits,
16     embedding=HuggingFaceEmbeddings(model_name="all-MiniLM-L6-v2"))
17
18 retriever = vectorstore.as_retriever(search_kwargs={"k": 3})
```

Retrieval Strategies

Different Retrieval Methods:

```
1 # Returns the most semantically similar documents to the user
2 # query based purely on vector distance (top-k closest matches).
3 retriever = vectorstore.as_retriever(search_type="similarity",
4     search_kwargs={"k": 4})
5
6 # MMR (Maximum Marginal Relevance) Balances relevance and diversity,
7 # ensuring the retrieved documents aren't repetitive by reducing semantic redundancy among the top-k results.
8 retriever = vectorstore.as_retriever(search_type="mmr",
9     search_kwargs={"k": 4, "fetch_k": 20})
10
11 # Filters out documents that fall below a minimum similarity score,
12 # ensuring only high-confidence matches are returned (still limited by k)
13 retriever = vectorstore.as_retriever(search_type="similarity_score_threshold",
14     search_kwargs={"score_threshold": 0.8, "k": 4})
15
16 # Manual search: Directly retrieves the top-k most similar documents
17 # without using a retriever wrapper same as search_type "similarity" but explicitly called.
18 results = vectorstore.similarity_search("What is machine learning?", k=3)
19
20 # With scores: Same as above but also returns the similarity
21 # score for each document, helping you inspect retrieval quality
22 results_with_scores = vectorstore.similarity_search_with_score(
23     "What is machine learning?", k=3)
```

Chains with LCEL

Modern Chains: LCEL Overview

What Changed?

- ▶ **Old:** LLMChain, SimpleSequentialChain (Deprecated)
- ▶ **New:** LCEL with pipe operator |

LCEL Advantages:

- ▶ **Composability:** Chain components naturally
- ▶ **Streaming:** Built-in streaming support
- ▶ **Async:** Native async/await
- ▶ **Batch:** Process multiple inputs
- ▶ **Fallbacks:** Error handling built-in
- ▶ **Parallelization:** Run steps in parallel
- ▶ **Observability:** Better tracing

Core Runnables:

- ▶ RunnablePassthrough: Pass data through
- ▶ RunnableParallel: Execute in parallel
- ▶ RunnableLambda: Custom functions
- ▶ RunnableBranch: Conditional execution

LCEL: Basic Chain Patterns

```
1 // Simple Chain:
2 from langchain_groq import ChatGroq
3 from langchain_core.prompts import ChatPromptTemplate
4 from langchain_core.output_parsers import StrOutputParser
5
6 prompt = ChatPromptTemplate.from_template("Tell me a joke about {topic}")
7 model = ChatGroq(model_name="gemma-7b-it")
8 output_parser = StrOutputParser()
9
10 chain = prompt | model | output_parser
11 result = chain.invoke({"topic": "programming"})
12
13 // Chain with Multiple Steps:
14 # Step 1: Generate topic
15 topic_chain = (
16     ChatPromptTemplate.from_template("Suggest a {genre} topic")
17     | ChatGroq(model_name="gemma-7b-it")
18     | StrOutputParser()
19 )
20
21 # Step 2: Write content
22 content_chain = (
23     ChatPromptTemplate.from_template("Write a story about: {topic}")
24     | ChatGroq(model_name="llama3-8b-8192")
25     | StrOutputParser()
26 )
27
28 # Combine: topic generates input for content
29 full_chain = {"topic": topic_chain} | content_chain
30 result = full_chain.invoke({"genre": "science fiction"})
```


LCEL: RAG Chain Pattern

Retrieval Augmented Generation:

```
1 from langchain_core.runnables import RunnablePassthrough
2 from langchain_groq import ChatGroq
3 # Assume 'retriever' from previous slide is defined
4
5 # Format documents
6 def format_docs(docs):
7     return "\n\n".join(doc.page_content for doc in docs)
8
9 # RAG chain
10 rag_chain = (
11     {
12         "context": retriever | format_docs,
13         "question": RunnablePassthrough()
14     }
15     | ChatPromptTemplate.from_template("""
16         Answer the question based on the context:
17
18         Context: {context}
19
20         Question: {question}
21     """)
22     | ChatGroq(model_name="llama3-8b-8192")
23     | StrOutputParser()
24 )
25
26 # Use it
27 answer = rag_chain.invoke("What is LangChain?")
```

LCEL: Parallel Execution

Benefits:

- ▶ Faster execution
- ▶ Clean code structure
- ▶ Easy to add/remove tasks

Run Multiple Chains in Parallel:

```
1 from langchain_core.runnables import RunnableParallel
2 from langchain_groq import ChatGroq
3 from langchain_core.prompts import ChatPromptTemplate
4 from langchain_core.output_parsers import StrOutputParser
5
6 # Define parallel tasks
7 parallel_chain = RunnableParallel(
8     summary=ChatPromptTemplate.from_template("Summarize: {text}")
9     | ChatGroq(model_name="llama3-8b-8192")
10    | StrOutputParser(),
11
12    keywords=ChatPromptTemplate.from_template("Extract keywords: {text}")
13    | ChatGroq(model_name="gemma-7b-it")
14    | StrOutputParser(),
15
16    sentiment=ChatPromptTemplate.from_template("Analyze sentiment: {text}")
17    | ChatGroq(model_name="gemma-7b-it")
18    | StrOutputParser()
19 )
20
21 # Execute all in parallel
22 results = parallel_chain.invoke({"text": "Long document content..."})
```

LCEL: Error Handling & Fallbacks

```
2 // Fallback to Alternative Model:
3 from langchain_groq import ChatGroq
4 from langchain_anthropic import ChatAnthropic
5
6 primary = ChatGroq(model_name="llama3-70b-8192")
7 fallback = ChatAnthropic(model="claude-3-opus-20240229")
8
9 chain = prompt | primary.with_fallbacks([fallback]) | output_parser
10
11 // Retry Logic: If Groq fails, automatically retries
12 from langchain_core.runnables import RunnableRetry
13
14 chain_with_retry = (
15     prompt
16     | RunnableRetry(max_attempts=3, wait_exponential_jitter=True)
17     | ChatGroq(model_name="gemma-7b-it")
18     | output_parser
19 )
20
21 // Custom Error Handling:
22 from langchain_core.runnables import RunnableLambda
23
24 def handle_errors(x):
25     try:
26         return x
27     except Exception as e:
28         return {"error": str(e)}
29
30 chain = prompt | model | RunnableLambda(handle_errors)
```

LCEL: Streaming

```
1 // Stream Tokens as Generated:
2 from langchain_groq import ChatGroq
3 # Assume prompt and output_parser are defined
4 chain = prompt | ChatGroq(model_name="gemma-7b-it") | StrOutputParser()
5
6 # Stream output
7 for chunk in chain.stream({"topic": "AI"}):
8     print(chunk, end="", flush=True)
9
10 // Async Streaming:
11 import asyncio
12
13 async def stream_response():
14     async for chunk in chain.astream({"topic": "AI"}):
15         print(chunk, end="", flush=True)
16
17 asyncio.run(stream_response())
18
19 // Streaming with Events:
20 # Get detailed streaming events
21 async for event in chain.astream_events({"topic": "AI"}, version="v1"):
22     kind = event["event"]
23     if kind == "on_chat_model_stream":
24         content = event["data"]["chunk"].content
25         print(content, end="", flush=True)
```

Agents & Tools

Modern Agents Overview

What Are Agents?

- ▶ Use LLMs as reasoning engines
- ▶ Dynamically choose which tools to use
- ▶ Iterative: observe, think, act, repeat
- ▶ Best for complex, multi-step tasks

Modern Agent Types (2024):

- ▶ **Tool Calling Agent:** Uses function calling (recommended)
- ▶ **JSON Chat Agent:** Output formatting based on Json
- ▶ **Structured Chat Agent:** For multi-input tools
- ▶ **ReAct Agent:** Reasoning and acting pattern

Deprecated (Don't Use):

- ▶ zero-shot-react-description
- ▶ conversational-react-description
- ▶ self-ask-with-search

(Ref: LangChain Agents Documentation 2024)

Creating Tools: Modern Approach

```
1 // Method 1: Using @tool Decorator:
  from langchain_core.tools import tool
3
4 @tool
5 def search_wikipedia(query: str) -> str:
6     """Search Wikipedia for information about a topic."""
7     from wikipedia import summary
8     try:
9         return summary(query, sentences=3)
10    except:
11        return "Could not find information."
12
13 @tool
14 def calculate(expression: str) -> str:
15     """Evaluate a mathematical expression."""
16     try:
17         return str(eval(expression))
18    except:
19        return "Invalid expression"
20
21 // Method 2: From Function:
22 \begin{lstlisting}[language=python, basicstyle=\tiny]
23 from langchain_core.tools import Tool
24
25 def get_weather(location: str) -> str:
26     """Get weather for a location."""
27     return f"Weather in {location}: Sunny, 72F"
28
29 weather_tool = Tool.from_function(
30     func=get_weather,
31     name="weather",
32     description="Get current weather for a location")
```

Tool Calling Agent Example

```
from langchain.groq import ChatGroq
2 from langchain.agents import create_tool_calling_agent, AgentExecutor
  from langchain.core.prompts import ChatPromptTemplate, MessagesPlaceholder
4 from langchain.core.tools import tool

6 @tool
  def search_wikipedia(query: str) -> str:
8     """Search Wikipedia for a topic."""
    # Dummy implementation
10    return f"Results for {query} from Wikipedia."

12 @tool
  def calculate(expression: str) -> str:
14     """Evaluate a mathematical expression."""
    return str(eval(expression))

16 tools = [search_wikipedia, calculate]

18 prompt = ChatPromptTemplate.from_messages([
20     ("system", "You are a helpful assistant"),
    ("human", "{input}"),
22     MessagesPlaceholder(variable_name="agent_scratchpad")
  ])

24 llm = ChatGroq(model_name="llama3-70b-8192", temperature=0)
26 agent = create_tool_calling_agent(llm, tools, prompt)

28 agent_executor = AgentExecutor(agent=agent, tools=tools, verbose=True)

30 result = agent_executor.invoke({"input": "What is the capital of France?"})
```


Modern Approach with bind_tools

```
2 from langchain.groq import ChatGroq
3 from langchain_core.tools import tool
4
5 @tool
6 def multiply(a: int, b: int) -> int:
7     """Multiply two numbers."""
8     return a * b
9
10 @tool
11 def add(a: int, b: int) -> int:
12     """Add two numbers."""
13     return a + b
14
15 # Bind tools to model
16 llm = ChatGroq(model_name="llama3-8b-8192")
17 llm_with_tools = llm.bind_tools([multiply, add])
18
19 # Invoke
20 response = llm_with_tools.invoke("What is 3 times 4 plus 5?")
21
22 # Check if tool was called
23 if response.tool_calls:
24     tool_call = response.tool_calls[0]
25     print(f"Tool: {tool_call['name']}")
26     print(f"Args: {tool_call['args']}")
```

Output Parsers

Output Parsers Overview

Why Output Parsers?

- ▶ LLMs return unstructured text
- ▶ Applications need structured data
- ▶ Parsers extract and validate output

Common Parser Types:

- ▶ **StrOutputParser**: Basic string extraction
- ▶ **JsonOutputParser**: Parse JSON responses
- ▶ **PydanticOutputParser**: Structured data with validation
- ▶ **StructuredOutputParser**: Multiple fields
- ▶ **CommaSeparatedListOutputParser**: Lists

Modern Best Practice:

- ▶ Use OpenAI's `with_structured_output()` when possible
- ▶ More reliable than prompt-based parsing
- ▶ Leverages function calling internally

Structured Output: Modern Approach

Benefits:

- ▶ Type-safe responses
- ▶ Automatic validation
- ▶ More reliable than prompt-based parsing

```
1 from langchain_groq import ChatGroq
2 from pydantic import BaseModel, Field
3
4 # Define output schema
5 class Person(BaseModel):
6     name: str = Field(description="Person's name")
7     age: int = Field(description="Person's age")
8     occupation: str = Field(description="Person's job")
9
10 # Create model with structured output (relies on tool-calling)
11 llm = ChatGroq(model_name="llama3-70b-8192")
12 structured_llm = llm.with_structured_output(Person)
13
14 # Use in chain
15 response = structured_llm.invoke(
16     "Tell me about a software engineer named Alice who is 28"
17 )
18
19 # Response is a Pydantic object
20 print(response.name)
21 print(response.age)
22 print(response.occupation)
```

Pydantic Output Parser (Alternative)

```
1 from langchain.output_parsers import PydanticOutputParser
2 from langchain.groq import ChatGroq
3 from pydantic import BaseModel, Field, validator
4 from typing import List
5
6 class MovieReview(BaseModel):
7     title: str = Field(description="Movie title")
8     rating: int = Field(description="Rating from 1-10")
9
10 parser = PydanticOutputParser(pydantic_object=MovieReview)
11
12 # Add to prompt
13 prompt = ChatPromptTemplate.from_template("""
14 Review the movie: {movie}
15
16 {format_instructions}""")
17
18 chain = (prompt.partial(format_instructions=parser.get_format_instructions())
19         | ChatGroq(model_name="gemma-7b-it") | parser)
20
21 result = chain.invoke({"movie": "Inception"})
```

Output Parser Error Handling

```
1 from langchain.groq import ChatGroq
  # Assume 'parser' is a defined PydanticOutputParser
3
  // OutputFixingParser (Auto-fix Errors):
5 from langchain.output_parsers import OutputFixingParser
7
  # If parsing fails, use LLM to fix
fixing_parser = OutputFixingParser.from_llm(
9     parser=parser,
    llm=ChatGroq(model_name="gemma-7b-it")
11 )
13
  // Automatically fixes malformed output
  # result = fixing_parser.parse(malformed_output)
15
  // RetryOutputParser (Retry with Context):
17 from langchain.output_parsers import RetryWithErrorOutputParser
19
  retry_parser = RetryWithErrorOutputParser.from_llm(
    parser=parser,
21     llm=ChatGroq(model_name="llama3-8b-8192")
    )
23
  # Retries with both output and original prompt
25 # result = retry_parser.parse_with_prompt(...)
```

LangChain Ecosystem

LangChain Ecosystem Components

- ▶ **LangChain Core:**
 - ▶ Base abstractions and LCEL
 - ▶ Foundation for all other packages
- ▶ **LangChain Community:**
 - ▶ Third-party integrations
 - ▶ Vector stores, document loaders
 - ▶ Community-maintained tools
- ▶ **LangGraph:**
 - ▶ Build stateful, multi-actor applications
 - ▶ Complex agent workflows with cycles
 - ▶ State management for agents
- ▶ **LangServe:**
 - ▶ Deploy chains as REST APIs
 - ▶ Automatic FastAPI generation
 - ▶ Production deployment
- ▶ **LangSmith:**
 - ▶ Debugging and monitoring
 - ▶ Tracing and evaluation
 - ▶ Dataset management

LangGraph: Stateful Agents

- ▶ Build complex, stateful agent workflows
- ▶ Support for cycles and conditional logic
- ▶ Persist state across interactions
- ▶ Multiple agents working together

```
1 from langgraph.graph import StateGraph
2 from typing import TypedDict, Annotated
3 import operator
4
5 class AgentState(TypedDict):
6     messages: Annotated[list, operator.add]
7     next: str
8
9 def call_model(state):
10     response = llm.invoke(state["messages"])
11     return {"messages": [response]}
12
13 def should_continue(state):
14     if len(state["messages"]) > 5:
15         return "end"
16     return "continue"
17
18 workflow = StateGraph(AgentState)
19 workflow.add_node("agent", call_model)
20 workflow.add_conditional_edges("agent", should_continue)
21 workflow.set_entry_point("agent")
22
23 app = workflow.compile()
```

LangServe: Deploy as API

Automatic Features:

- ▶ Interactive playground at /chat/playground
- ▶ OpenAPI docs at /docs
- ▶ Streaming support
- ▶ Batch processing

Deploy Any Chain as REST API:

```
1 from fastapi import FastAPI
2 from langserve import add_routes
3 from langchain-groq import ChatGroq
4 from langchain_core.prompts import ChatPromptTemplate
5
6 # Create chain
7 prompt = ChatPromptTemplate.from_template("Tell me about {topic}")
8 chain = prompt | ChatGroq(model_name="gemma-7b-it")
9
10 # Create FastAPI app
11 app = FastAPI(
12     title="LangChain API", version="1.0",
13     description="API for my LangChain application")
14
15 # Add chain as route
16 add_routes(app, chain, path="/chat")
17
18 # Run: uvicorn app:app --reload
19 # API available at: http://localhost:8000/chat
```

LangSmith: Monitoring & Debugging

What is LangSmith?

- ▶ Platform for debugging LLM applications
- ▶ Trace every step of chain execution
- ▶ Evaluate model performance
- ▶ Manage test datasets
- ▶ Monitor production applications

Setup:

```
1 import os
  from langchain_groq import ChatGroq
3 from langchain_core.prompts import ChatPromptTemplate
  from langchain_core.output_parsers import StrOutputParser
5
  os.environ["LANGCHAIN_TRACING_V2"] = "true"
7 os.environ["LANGCHAIN_API_KEY"] = "your-langsmith-api-key"
  os.environ["LANGCHAIN_PROJECT"] = "my-groq-project"
9
  # Now all chain executions are automatically traced
11 prompt = ChatPromptTemplate.from_template("Tell me about {topic}")
  llm = ChatGroq(model_name="gemma-7b-it")
13 output_parser = StrOutputParser()
  chain = prompt | llm | output_parser
15 result = chain.invoke({"topic": "Large Language Models"})
17
  # View traces at: https://smith.langchain.com
```

Best Practices

Best Practices: Error Handling

```
1 // 1. Use Fallbacks:
2 from langchain_groq import ChatGroq
3 from langchain_anthropic import ChatAnthropic
4
5 primary = ChatGroq(model_name="llama3-70b-8192")
6 fallback = ChatAnthropic(model="claude-3-opus-20240229")
7
8 chain = prompt | primary.with_fallbacks([fallback]) | parser
9
10 // 2. Implement Retries:
11 from langchain_core.runnables import RunnableRetry
12
13 chain_with_retry = (
14     prompt
15     | RunnableRetry(
16         max_attempts=3,
17         wait_exponential_jitter=True
18     )
19     | ChatGroq(model_name="gemma-7b-it")
20     | parser
21 )
22
23 // 3. Graceful Degradation:
24 try:
25     result = chain.invoke(input_data)
26 except Exception as e:
27     # logger.error(f"Chain failed: {e}")
28     # result = fallback_response()
29     pass
```

Best Practices: Token Management with Groq

Groq Models - Fast and Free Tier:

- ▶ gemma2-9b-it: Ultra-fast, 8K context, great for chat
- ▶ llama3-8b-8192: Balanced, 8K context, general purpose
- ▶ llama3-70b-8192: Most capable, 8K context, complex tasks
- ▶ llama-3.1-8b-instant: Fastest, optimized for low latency
- ▶ Groq billed by requests/day on free tier, not tokens

```
1 from langchain.groq import ChatGroq
3 # Choose model based on needs
llm_fast = ChatGroq(
5     model_name="gemma2-9b-it",
6     max_tokens=500,
7     temperature=0.7,
8     max_retries=2
9 )
11 # Monitor context length manually
from transformers import AutoTokenizer
13 tokenizer = AutoTokenizer.from_pretrained("google/gemma-2-9b-it")
token_count = len(tokenizer.encode(text))
15
# Groq's speed allows batch processing without latency concerns
```

Best Practices: Security & Privacy

- ▶ Be aware of the data usage policies of your LLM provider.
- ▶ Groq has a zero-retention policy for API data.
- ▶ Consider self-hosted models for maximum data control.
- ▶ Implement PII detection and redaction before sending data.

```
// 1. API Key Management:
2 import os
  from dotenv import load_dotenv
4
  load_dotenv() # Load from .env file
6 api_key = os.getenv("GROQ_API_KEY")

// 2. Timeouts and Retries:
8 from langchain_groq import ChatGroq
10
  llm = ChatGroq(model_name="llama3-8b-8192",
12     temperature=0,
13     max_retries=2,
14     request_timeout=30) # seconds

// 3. Input Validation:
16 def validate_input(user_input: str) -> str:
18     if len(user_input) > 8000: # Sanitize input
19         raise ValueError("Input too long for the model context.")
20     # Remove potential injection attempts
    return user_input.strip()
```

Best Practices: Choosing the Right Model

Groq Model Selection Guide:

Use Case	Model	Why
Chatbots	<code>gemma2-9b-it</code>	Fast, efficient, good reasoning
Summarization	<code>llama3-8b-8192</code>	Balanced speed/quality
Complex reasoning	<code>llama3-70b-8192</code>	Most capable
Ultra-low latency	<code>llama-3.1-8b-instant</code>	Optimized for speed
Code generation	<code>llama3-70b-8192</code>	Better instruction following

```

1 # Pattern: Start with fast model, fallback to capable model
2 from langchain_groq import ChatGroq
3
4 primary = ChatGroq(model_name="gemma2-9b-it", temperature=0.7)
5 fallback = ChatGroq(model_name="llama3-70b-8192", temperature=0.7)
6
7 chain = prompt | primary.with_fallbacks([fallback]) | parser

```


Best Practices: Async Patterns

Benefits:

- ▶ Faster parallel processing
- ▶ Better resource utilization
- ▶ Improved user experience with streaming

```
1 // Use Async for Better Performance:
import asyncio
3
4 async def process_multiple_queries(queries):
5     # Process queries concurrently
6     tasks = [chain.ainvoke({"input": q}) for q in queries]
7     results = await asyncio.gather(*tasks)
8     return results
9
10 # Run
11 queries = ["Query 1", "Query 2", "Query 3"]
12 results = asyncio.run(process_multiple_queries(queries))
13 // Async Streaming:
14 async def stream_response(input_text):
15     async for chunk in chain.astream({"input": input_text}):
16         print(chunk, end="", flush=True)
17
18 asyncio.run(stream_response("Tell me about AI"))
```

What's New in LangChain Ecosystem

(Oct 2025, Release of 1.0 version)

Why LangGraph 1.0 is Nearly Non-Breaking and Production-Ready

- ▶ LangGraph 1.0 introduces minimal breaking changes from prior releases.
- ▶ Existing LangGraph implementations continue to work as-is.
- ▶ Core runtime stable across Python and TypeScript.
- ▶ Supports durable execution with checkpoints and rollback.
- ▶ State can persist via Postgres or SQLite for reliability.
- ▶ Human-in-the-loop and interrupt handling built-in.
- ▶ Ready for production-grade use; used by Uber, LinkedIn, Klarna, JPMorgan, Cloudflare.

How LangChain 1.0 Simplifies Agents to Just 10 Lines of Code

- ▶ LangChain 1.0 Alpha released; full 1.0 expected by late October.
- ▶ Rewritten as a simplified agent runtime built on LangGraph.
- ▶ Agents can now be created in roughly 10 lines of code.
- ▶ Legacy LangChain moved to “LangChain Classic.”
- ▶ Functionally similar to OpenAI or Pedantic SDK agents.
- ▶ Unified “create_agent” abstraction across languages.
- ▶ Integrates easily with external SDKs like Google’s 80k.

Durable Execution, Streaming, Human-in-the-Loop, and Time Travel

- ▶ Agents persist state with rollback and checkpointing.
- ▶ Supports human approval and manual intervention mid-run.
- ▶ Time travel enables returning to earlier workflow states.
- ▶ Four streaming modes: messages, updates, values, custom.
- ▶ Update streaming ideal for dashboards and UX refresh.
- ▶ State persistence possible locally or via Postgres.
- ▶ Enables retry and branch execution for debugging and auditing.

The New Standardized Content Blocks for Easier Model Switching

- ▶ LangChain Core now uses standardized content/message blocks.
- ▶ Abstracts input/output across providers like OpenAI and Anthropic.
- ▶ Simplifies switching models without rewriting logic.
- ▶ Unifies multimodal inputs, reasoning traces, and tool calls.
- ▶ Middleware layer ensures consistent formatting and metadata.
- ▶ Supports normalized logging across different model providers.
- ▶ Key enabler for multi-model experimentation and portability.

Real-World Use Cases

- ▶ Meeting notes enrichment and summarization pipelines.
- ▶ Financial EDI approval flows with human-in-loop steps.
- ▶ Real-time dashboards using streaming “update” mode.
- ▶ Persistent agent state across cloud or local environments.
- ▶ Dynamic model selection for efficiency (e.g., switching from Claude to Llama3).
- ▶ Internal RAG/Graph-RAG systems leveraging PGVector.
- ▶ Used at scale in production by enterprises and startups alike.

Example: Defining a Simple Agent in LangChain 1.0

```
1 from langchain.agents import create_agent
2
3 agent = create_agent(
4     tools=[search, summarize],
5     model="gpt-4-turbo"
6 )
7
8 result = agent("Summarize today's meeting notes")
9 print(result)
```


Anti-Patterns to Avoid

- ▶ Token-only streaming without structured updates.
- ▶ Relying on ephemeral in-memory state, persist externally.
- ▶ Single-function “flat” agents instead of structured graphs.
- ▶ Provider lock-in design for model and SDK portability.
- ▶ Neglecting checkpointing and recovery strategies.
- ▶ Ignoring observability, use LangSmith, LangFuse, or OpenTelemetry.
- ▶ Skipping standardized content blocks creates fragility.

Conclusions

What is LangChain?

LangChain: A Comprehensive Framework for Building LLM-Powered Applications

- ▶ **Core Purpose:** Simplify development of complex applications using Large Language Models
- ▶ **Key Components:**
 - ▶ Unified interface for multiple language models
 - ▶ Advanced prompt management
 - ▶ Flexible data integration
 - ▶ Intelligent agent and tool systems
- ▶ **Founding:** Created by Harrison Chase as an open-source project
- ▶ **Availability:** Python and JavaScript libraries

Feature	Description
Data Awareness	Connect LLMs to external data sources
Agentic Capability	Enable dynamic interaction with environment
Model Flexibility	Support for multiple LLM providers

(Ref: LangChain Framework Overview)

Why LangChain?

Addressing Challenges in LLM Application Development

► Limitations in Current LLM Tooling:

- Fragmented model interfaces
- Complex prompt management
- Limited external data integration
- Lack of flexible reasoning mechanisms

► LangChain Solutions:

- Model-agnostic framework
- Standardized prompt engineering
- Seamless external data augmentation
- Intelligent agent orchestration

► Key Advantages:

- Rapid prototyping of AI applications
- Simplified model and tool integration
- Scalable architecture for complex use cases

Enabling Developers to Build Sophisticated AI Systems with Ease

(Ref: LangChain Development Principles)

One pager

- **Models:** Building blocks supporting different AI model types - LLMs, Chat, Text Embeddings.
- **Prompts:** Inputs constructed from various components. LangChain offers easy interfaces - Prompt Templates, Example Selectors, Output Parsers.
- **Memory:** Stores/retrieves messages, short or long term, in conversations.
- **Indexes:** Assist LLMs with documents - Document Loaders, Text Splitters, Vector Stores, Retrievers.
- **Chains:** Combine components or chains in order to accomplish tasks.
- **Agents:** Empower LLMs to interact with external systems, make decisions, and complete tasks using Tools.



(Ref: Building Generative AI applications made easy with Vertex AI PaLM API and LangChain - Anand Iyer, Rajesh Thallam)

Resources & Next Steps

Official Documentation:

- ▶ LangChain Docs: <https://python.langchain.com>
- ▶ LangChain Blog: <https://blog.langchain.dev>
- ▶ LangChain Academy: <https://academy.langchain.com>
- ▶ API Reference: <https://api.python.langchain.com>

GitHub Repositories:

- ▶ Core: <https://github.com/langchain-ai/langchain>
- ▶ Templates:
<https://github.com/langchain-ai/langchain/tree/master/templates>
- ▶ LangGraph: <https://github.com/langchain-ai/langgraph>

Community:

- ▶ Discord: <https://discord.gg/langchain>
- ▶ Twitter: @LangChainAI
- ▶ YouTube: LangChain official channel

Practice:

- ▶ Start with simple LCEL chains
- ▶ Build a RAG application
- ▶ Create custom tools and agents
- ▶ Deploy with LangServe

Closing Thoughts

- ▶ Langchain's usefulness in solving problems today
- ▶ Possibility of LLM APIs expanding capabilities over time
- ▶ Potential for Langchain to become an interface to LLMs
- ▶ Acknowledgment of Langchain's valuable contributions and community efforts
- ▶ Appreciation for the work done by Harrison and the Langchain team

(Ref: What is Langchain and why should I care as a developer? - Logan Kilpatrick)

References

Many publicly available resources have been refereed for making this presentation. Some of the notable ones are:

- ▶ Intro to LangChain for Beginners - A code-based walkthrough- Menlo Park Lab
- ▶ LangChain Crash Course (10 minutes): Easy-to-Follow Walkthrough of the Most Important Concepts - Menlo Park Lab
- ▶ Official doc: <https://docs.langchain.com/docs/>
- ▶ Git Repo: <https://github.com/hwchase17/langchain>
- ▶ LangChain 101: The Complete Beginner's Guide Edrick
<https://www.youtube.com/watch?v=P3MAbZ2eMUI> A wonderful overview, don't miss
- ▶ Cookbook by Gregory Kamradt(Easy way to get started):
<https://github.com/gkamradt/langchain-tutorials/blob/main/LangChain%20Cookbook.ipynb>
- ▶ Youtube Tutorials: https://www.youtube.com/watch?v=_v_fgW2SkkQ
- ▶ LangChain 101 Course (updated 2024 with LCEL) - Ivan Reznikov
<https://github.com/IvanReznikov/DataVerse/tree/main/Courses/LangChain>
- ▶ A Complete LangChain Guide <https://nanonets.com/blog/langchain/>
- ▶ 7 Ways to Use LangSmith's Superpowers to Turboboost Your LLM Apps - Menlo Park Lab
- ▶ LangChain official blog: <https://blog.langchain.dev>
- ▶ LangChain Academy: <https://academy.langchain.com>
- ▶ LangChain templates repo:
<https://github.com/langchain-ai/langchain/tree/master/templates>
- ▶ Groq Documentation: <https://console.groq.com/docs>
- ▶ Groq Model Playground: <https://groq.com/>

Thanks ...

- ▶ Office Hours: Saturdays, 3 to 5 pm (IST); Free-Open to all; email for appointment to yogeshkulkarni at yahoo dot com
- ▶ Call + 9 1 9 8 9 0 2 5 1 4 0 6



(<https://www.linkedin.com/in/yogeshkulkarni/>)



(<https://medium.com/@yogeshharibhaukulkarni>)



(<https://www.github.com/yogeshhk/>)

Pune AI Community (PAIC)

▶ Two-way communication:

- ▶ Website [puneaicommunity dot org](https://puneaicommunity.org)
- ▶ Email [puneaicommunity at gmail dot com](mailto:puneaicommunity@gmail.com)
- ▶ Call + 9 1 9 8 9 0 2 5 1 4 0 6
- ▶ LinkedIn:
<https://linkedin.com/company/pune-ai-community>

▶ One-way Announcements:

- ▶ Twitter (X) @puneaicommunity
- ▶ Instagram @puneaicommunity
- ▶ WhatsApp Community: Invitation Link
<https://chat.whatsapp.com/LluOrhyEzuQLDr25ixZ>
- ▶ Luma Event Calendar: [puneaicommunity](https://puneaicommunity.org)

▶ Contribution Channels:

- ▶ GitHub: [Pune-AI-Community and puneaicommunity](https://github.com/pune-ai-community)
- ▶ Medium: [pune-ai-community](https://medium.com/pune-ai-community)
- ▶ YouTube: @puneaicommunity



Website

Pune AI Community (PAIC) QR codes



Website



Medium Blogs



Twitter-X



LinkedIn Page



Github Repository



WhatsApp Invite



Luma Events

YouTube Videos

Instagram