

AI VIETNAM
All-in-One Course
(TA Session)

Building LLMs application with LangChain

Extra Class: LLMs



AI VIET NAM
[@aivietnam.edu.vn](https://aivietnam.edu.vn)

Dinh-Thang Duong – TA
Nguyen-Thuan Duong – TA

Objectives



LangChain

```
Curl
curl -X 'POST' \
  'http://0.0.0.0:5000/generative_ai' \
  -H 'accept: application/json' \
  -H 'Content-Type: application/json' \
  -d '{
    "question": "what is BERT?"
}'

Request URL
http://0.0.0.0:5000/generative_ai

Server response
Code Details

200 Response body
{
  "answer": "BERT is a bidirectional transformer model for natural language processing pre-trained on a large corpus of text. It stands for Bidirectional Encoder Representations from Transformers. BERT was designed to be similar to OpenAI GPT, and it outperforms other language representation models in various tasks. The table shows the GLUE test results for BERT and OpenAI GPT, and the table shows the SQuAD 1.1 results for different models. BERT comes in two sizes: BERT BASE and BERT LARGE. The larger model has more parameters and generally performs better. Fine-tuning BERT on different tasks is illustrated in Figure 4. Additional details and ablation studies are presented in Appendices A, B, and C."
}
```

In this lecture, we will discuss about:

1. What is LLMs in Production?
2. What is Langchain?
3. Basic components of LangChain.
4. How to use LangChain to deploy an API?
5. How to use LangChain to deploy a RAG application?

Outline

- Introduction
- LangChain
- RAG with LangChain
- Question

Introduction

Introduction

❖ Getting Started



ChatGPT



Examples

"Explain quantum computing in simple terms" →

"Got any creative ideas for a 10 year old's birthday?" →

"How do I make an HTTP request in Javascript?" →



Capabilities

Remembers what user said earlier in the conversation

Allows user to provide follow-up corrections

Trained to decline inappropriate requests



Limitations

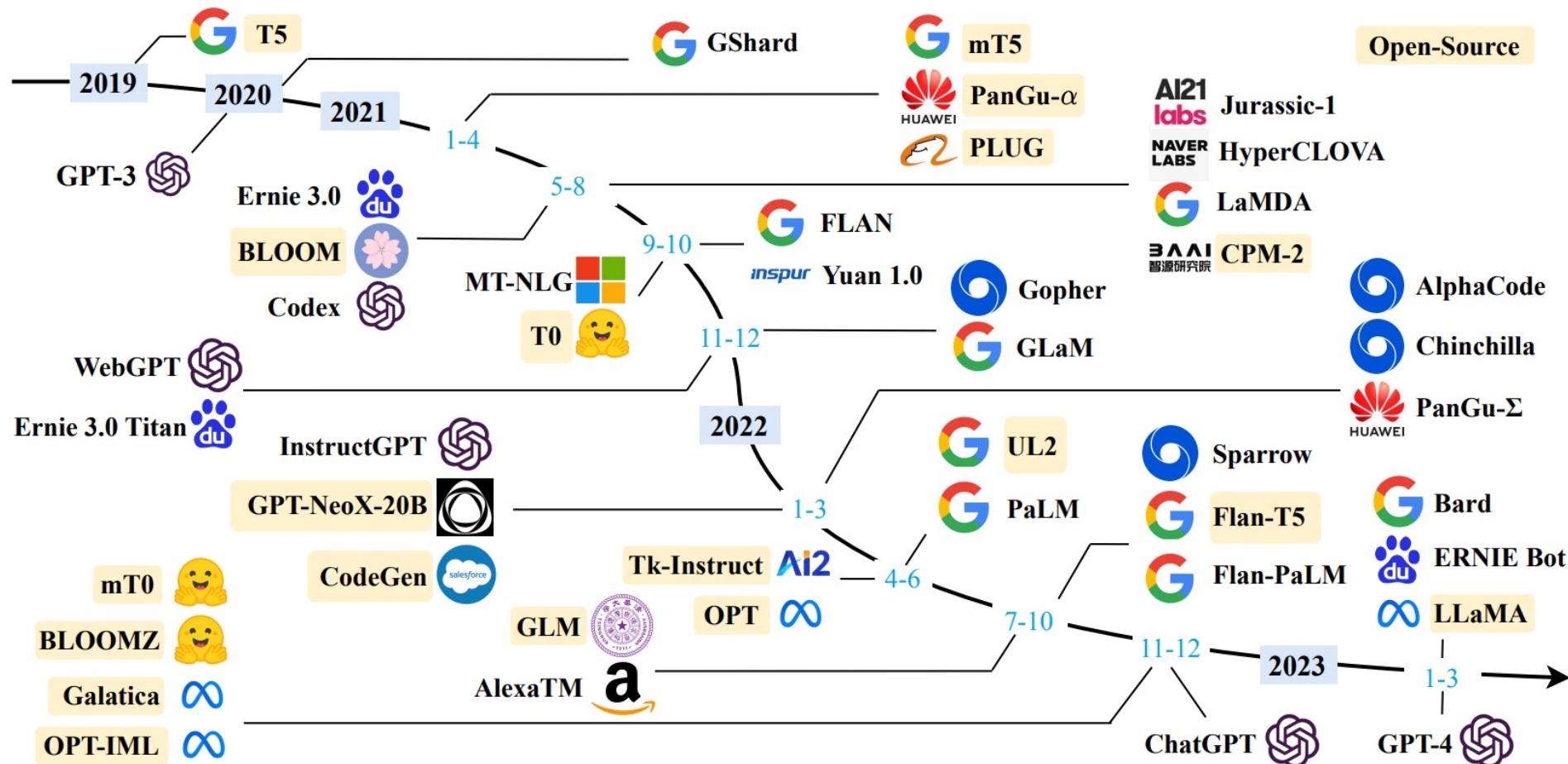
May occasionally generate incorrect information

May occasionally produce harmful instructions or biased content

Limited knowledge of world and events after 2021

Introduction

❖ LLMs size over time



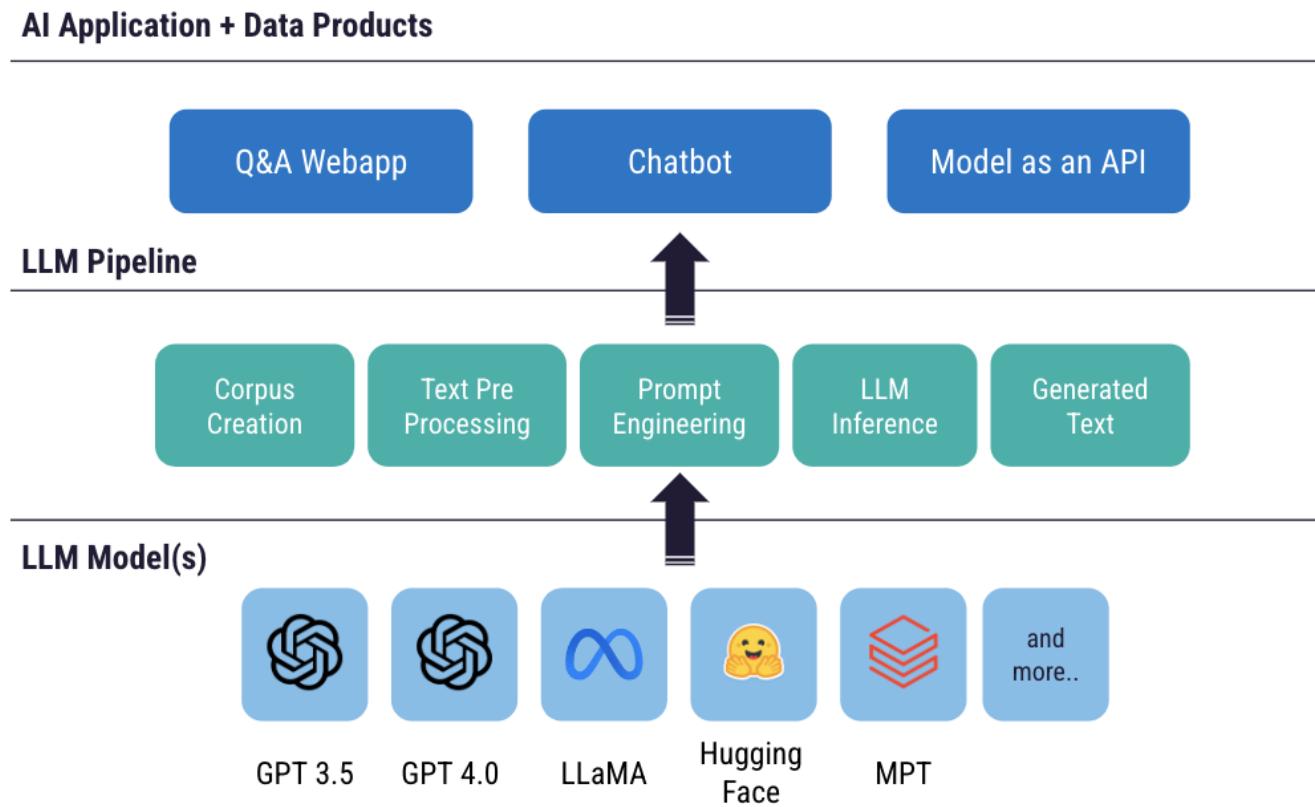
Introduction

❖ LLMs Applications



Introduction

❖ Getting Started



Overview of LLMs in Production

LangChain

LangChain

❖ Introduction

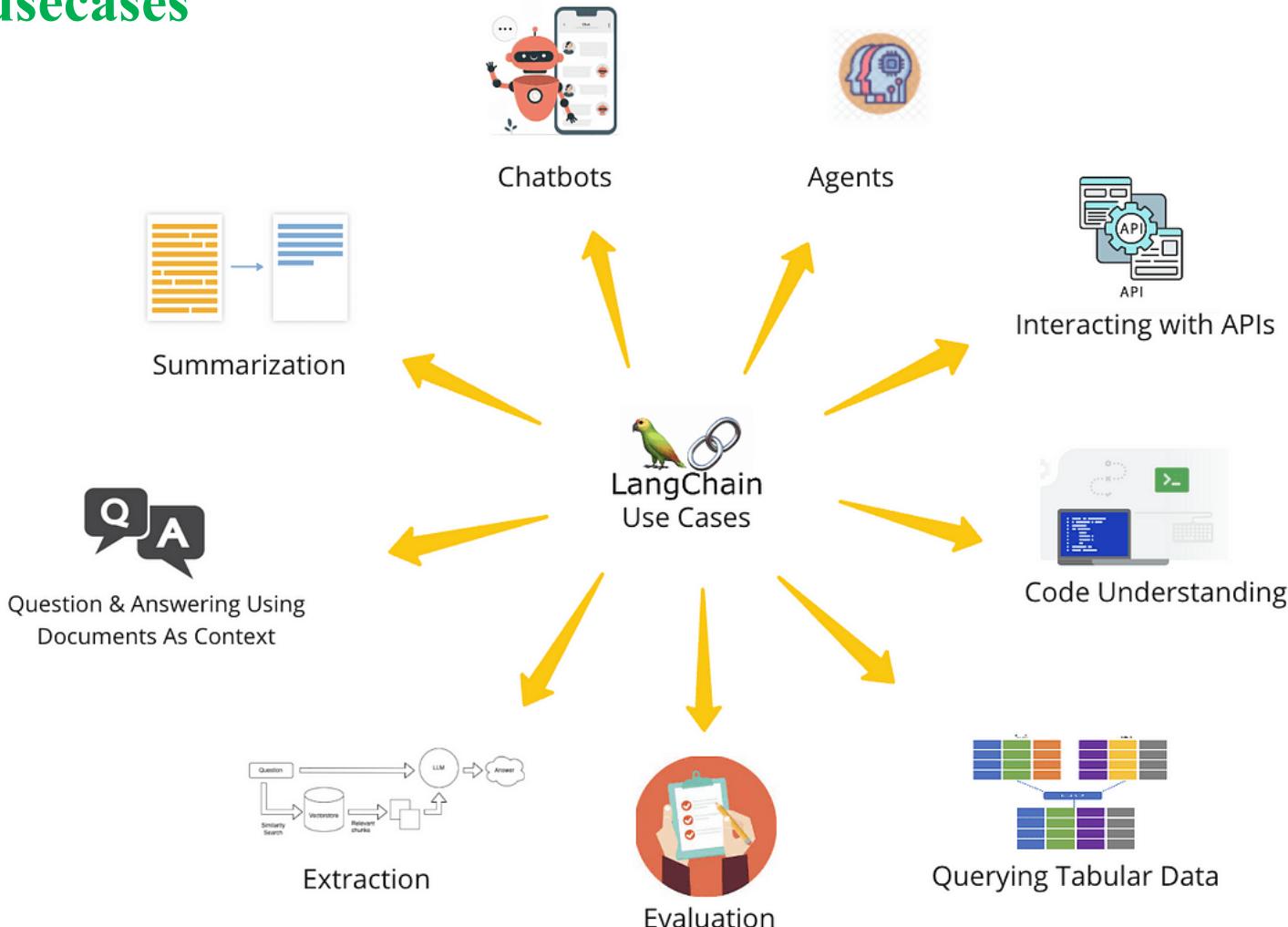


LangChain

LangChain: A framework for developing applications powered by large language models (LLMs). LangChain simplifies every stage of the LLM application lifecycle: Development, Productionization, Deployment.

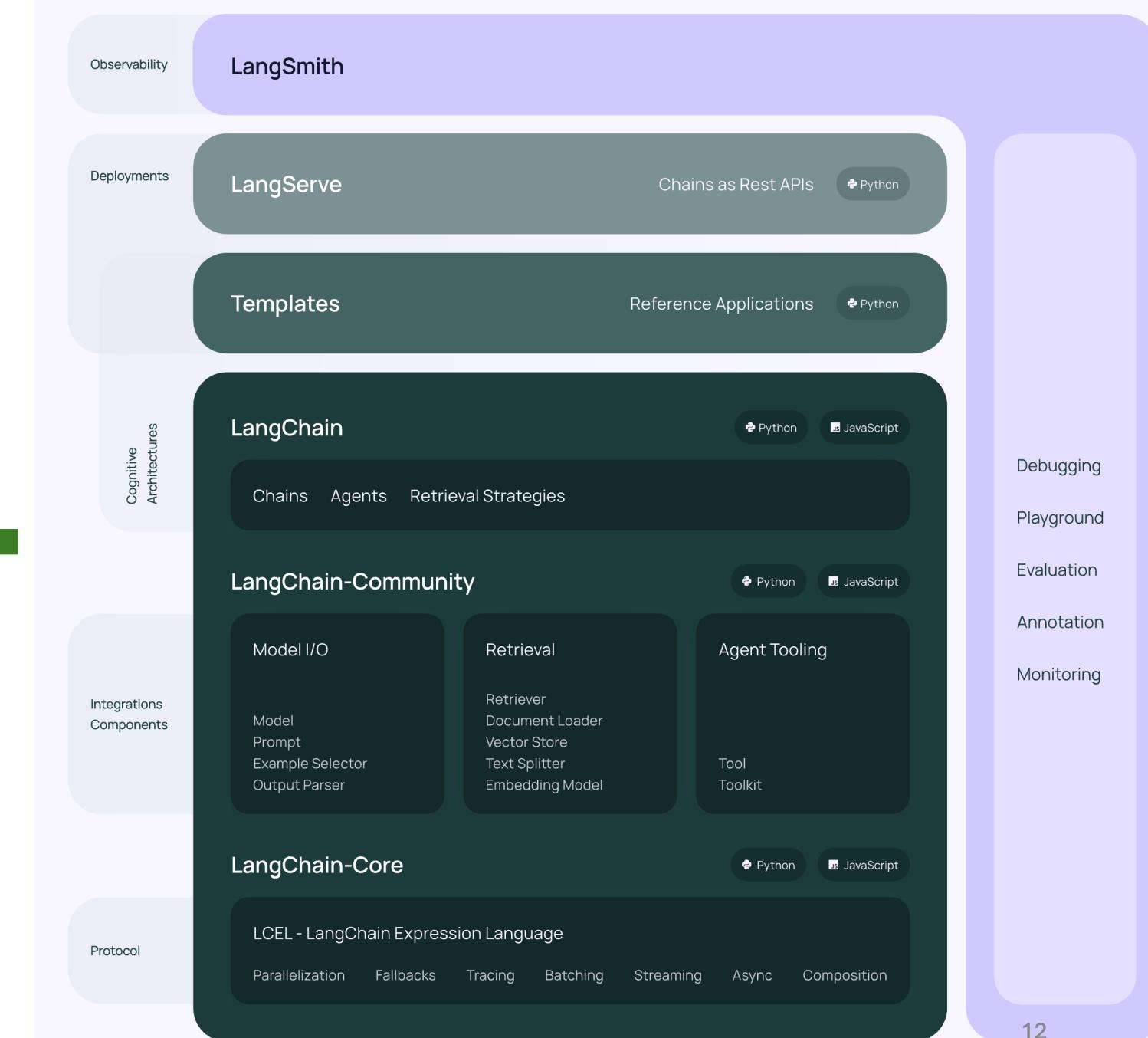
LangChain

❖ LangChain usecases



LangChain

❖ Introduction



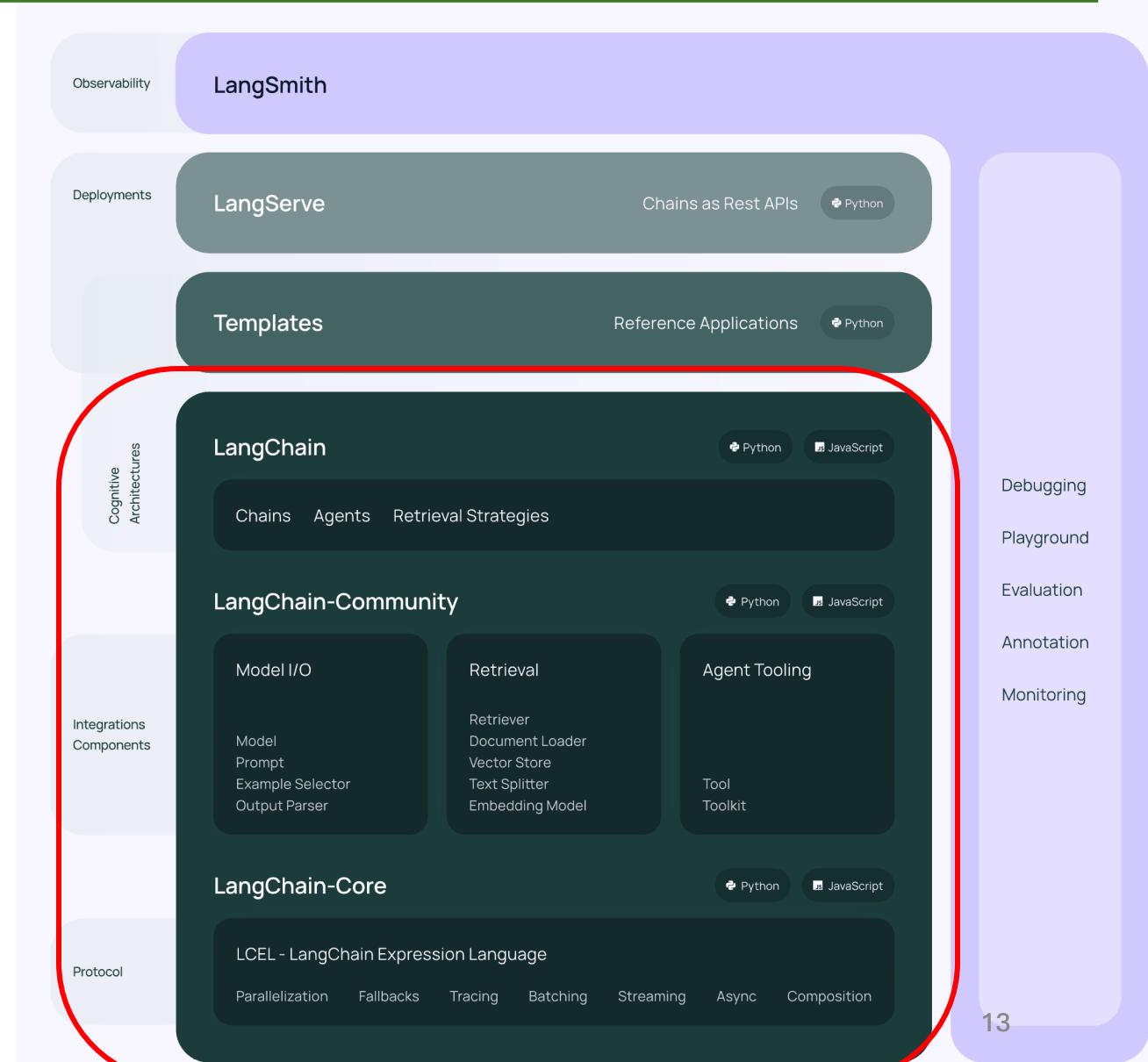
LangChain

❖ Introduction

Development: Build your applications using LangChain's open-source **building blocks** and **components**. Hit the ground running using **third-party integrations** and **Templates**.



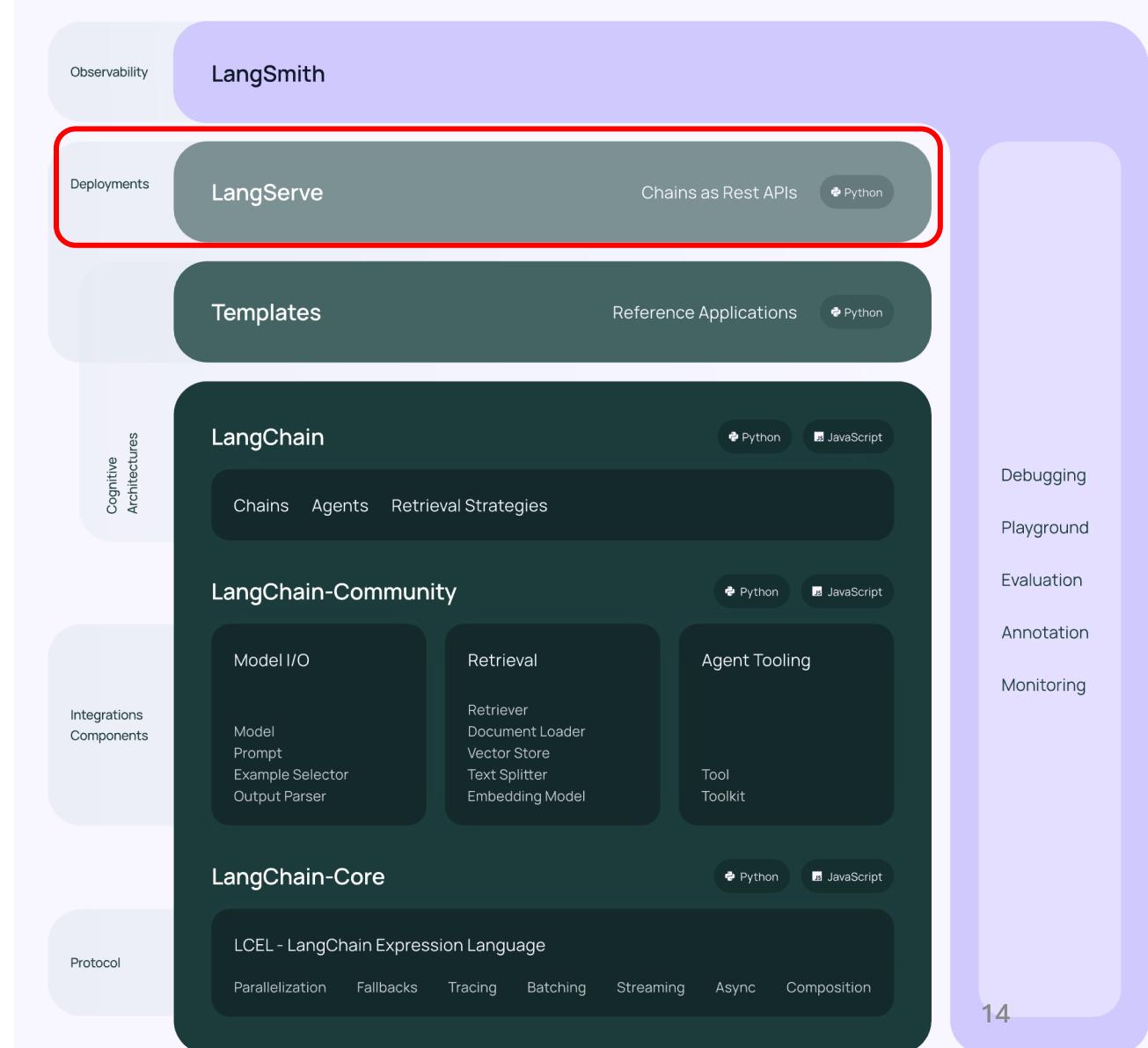
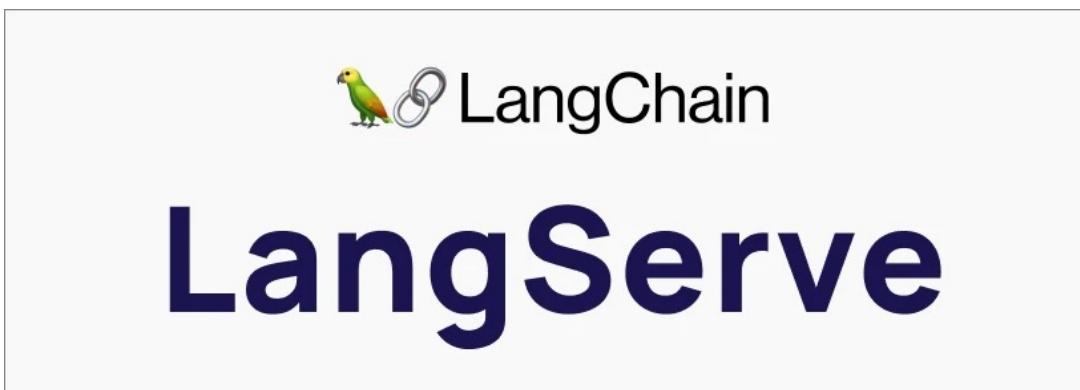
LangChain



LangChain

❖ Introduction

LangServe helps developers deploy LangChain runnables and chains as a REST API. This library is integrated with FastAPI and uses pydantic for data validation.



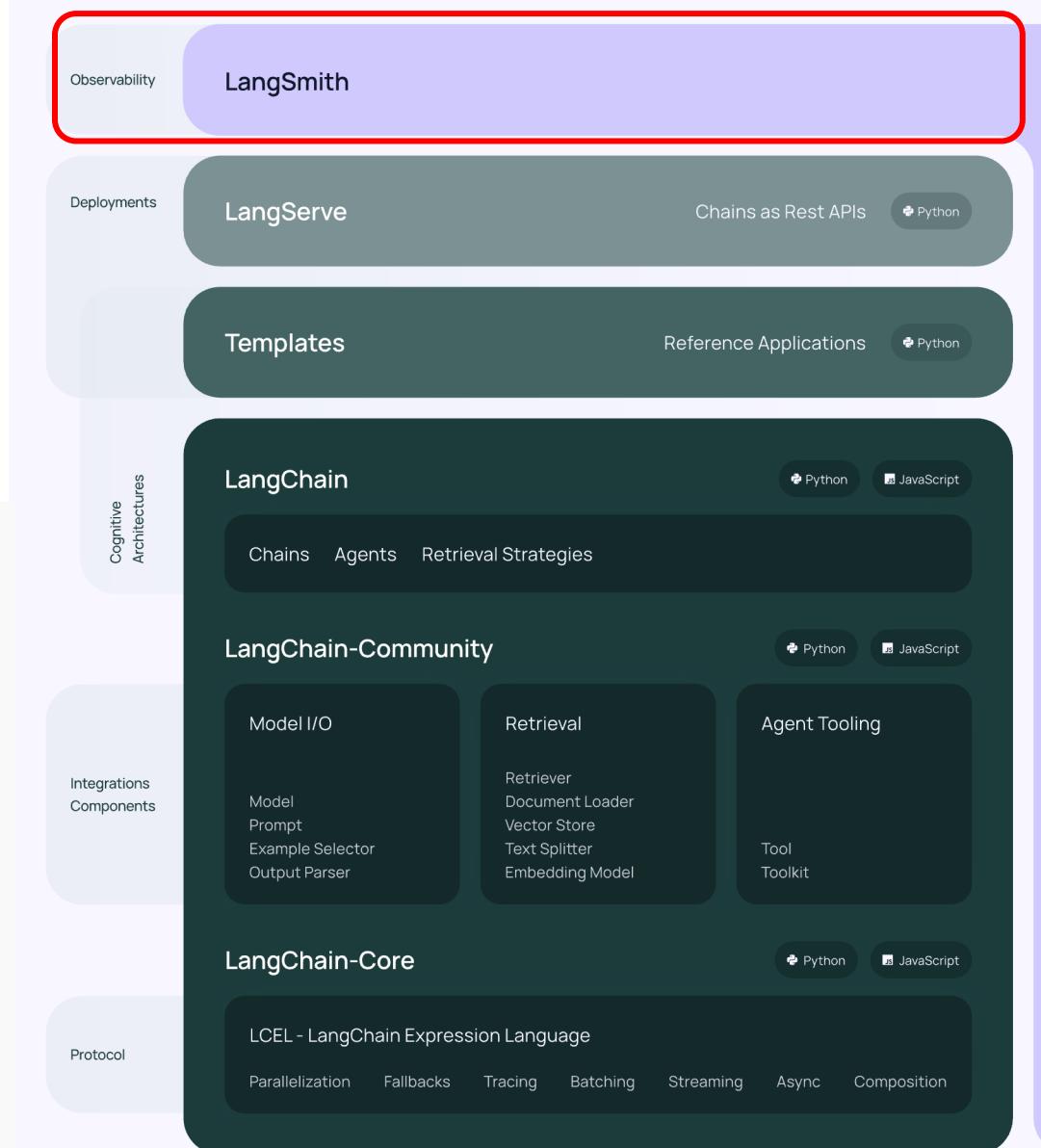
LangChain

❖ Introduction

LangSmith: used to inspect, monitor and evaluate your chains, so that you can continuously optimize and deploy with confidence (productionization).



LangSmith



LangChain

❖ LangSmith

The screenshot shows the LangSmith web application's project management interface. At the top, there is a navigation bar with icons for user profile, search, and developer mode. Below the navigation is a sidebar with various project-related icons. The main area is titled "Projects" and contains a table with one row. The columns are: Name (sorted by Most Recent Run), Feedback (7D), Run Count (7D), Error Rate (7D), % Streaming (7D), Total Tokens (7D), Total Cost (7D), P50 Latency (7D), P99 Latency (7D), and Most Recent Run (7D). The single project listed is "default". The table includes sorting arrows for each column header. A note at the bottom of the table states: "Current sorting (Most Recent Run) excludes projects with no runs. Sort by Name to see all." There are also navigation arrows and a "Show 10" button.

Name ↑	Feedback (7D)	Run Count (7D)	Error Rate (7D) ↑	% Streaming (7D)	Total Tokens (7D)	Total Cost (7D)	P50 Latency (7D) ↑	P99 Latency (7D) ↑	Most Recent Run (7D) ↓
default	1	0%	0%		540	7.31s	7.31s		5/5/2024, 5:29:20 AM

Current sorting (Most Recent Run) excludes projects with no runs. [Sort by Name](#) to see all.

< > Show 10

LangChain

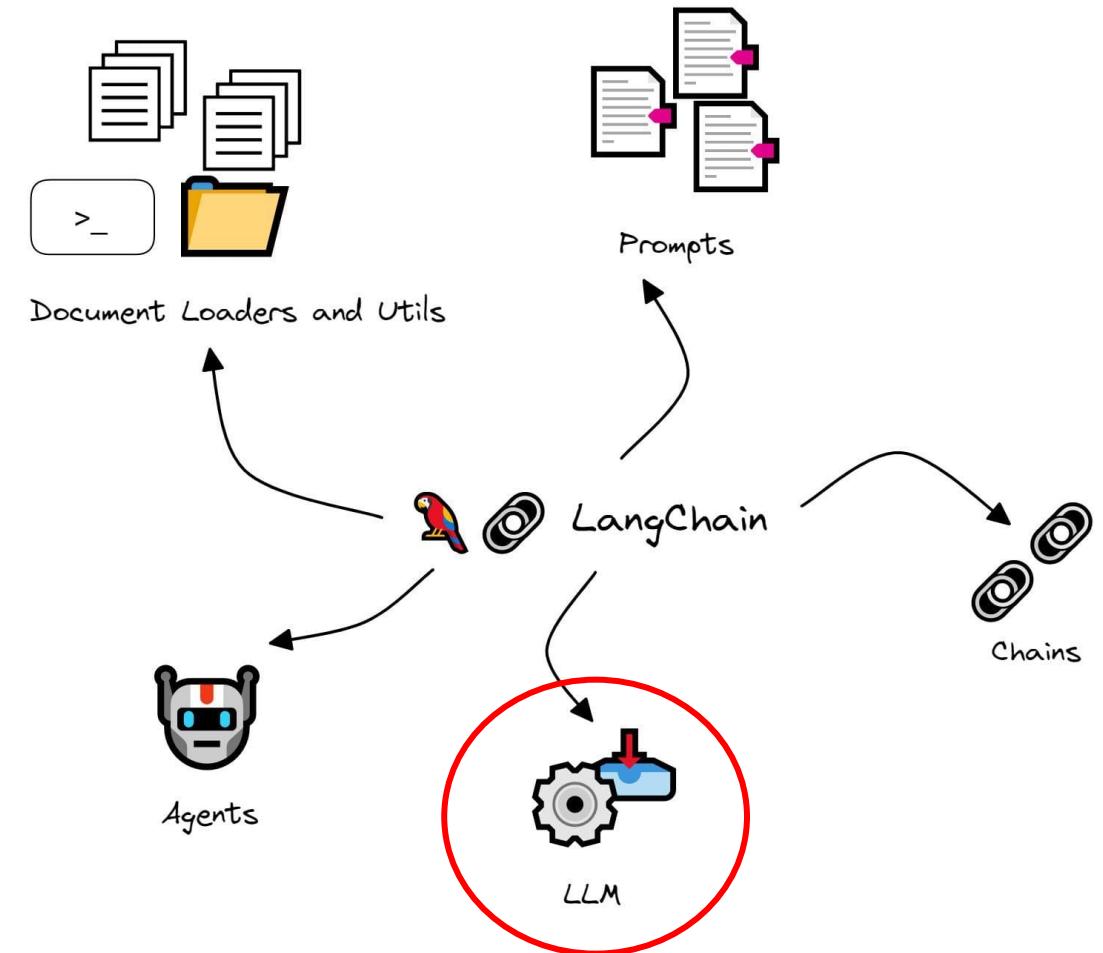
❖ LangChain Components



LangChain

❖ LangChain Components

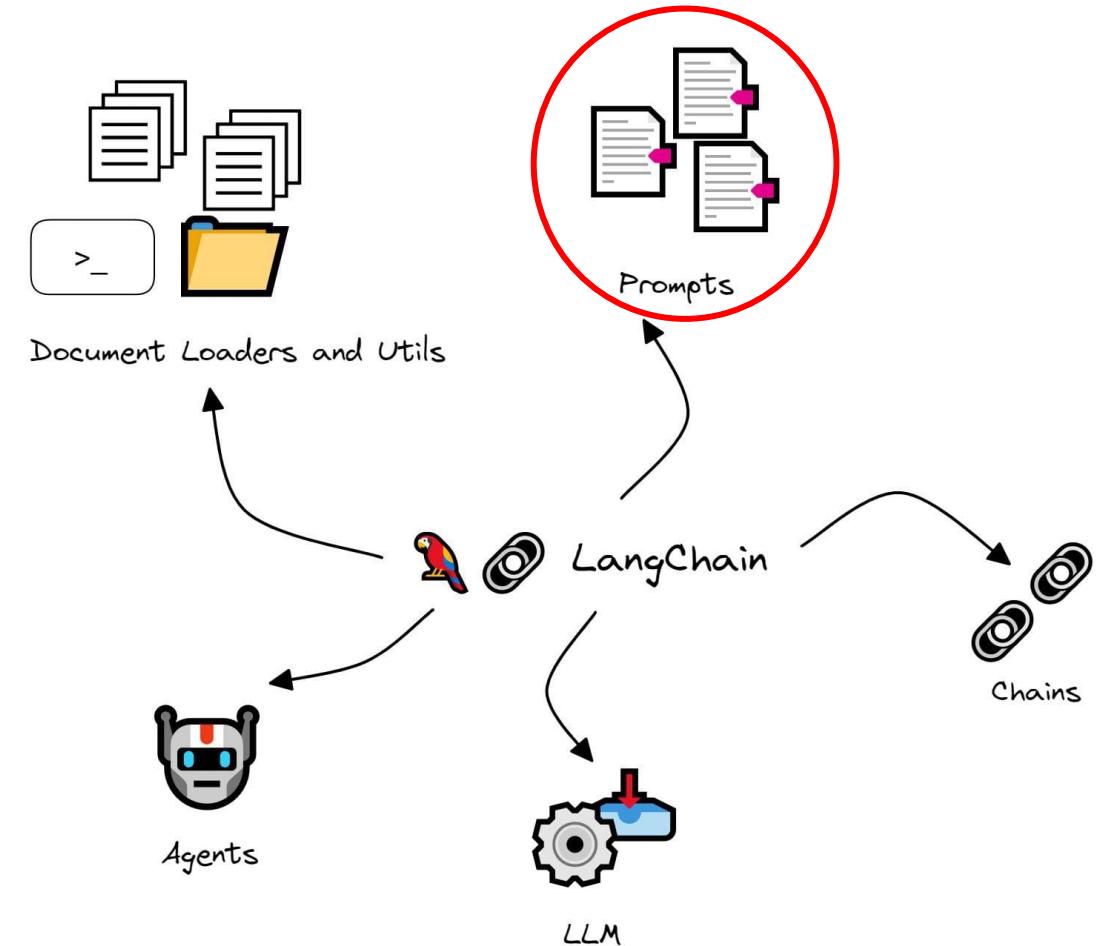
LLM (Large Language Model): The core AI model responsible for processing natural language and generating outputs.



LangChain

❖ LangChain Components

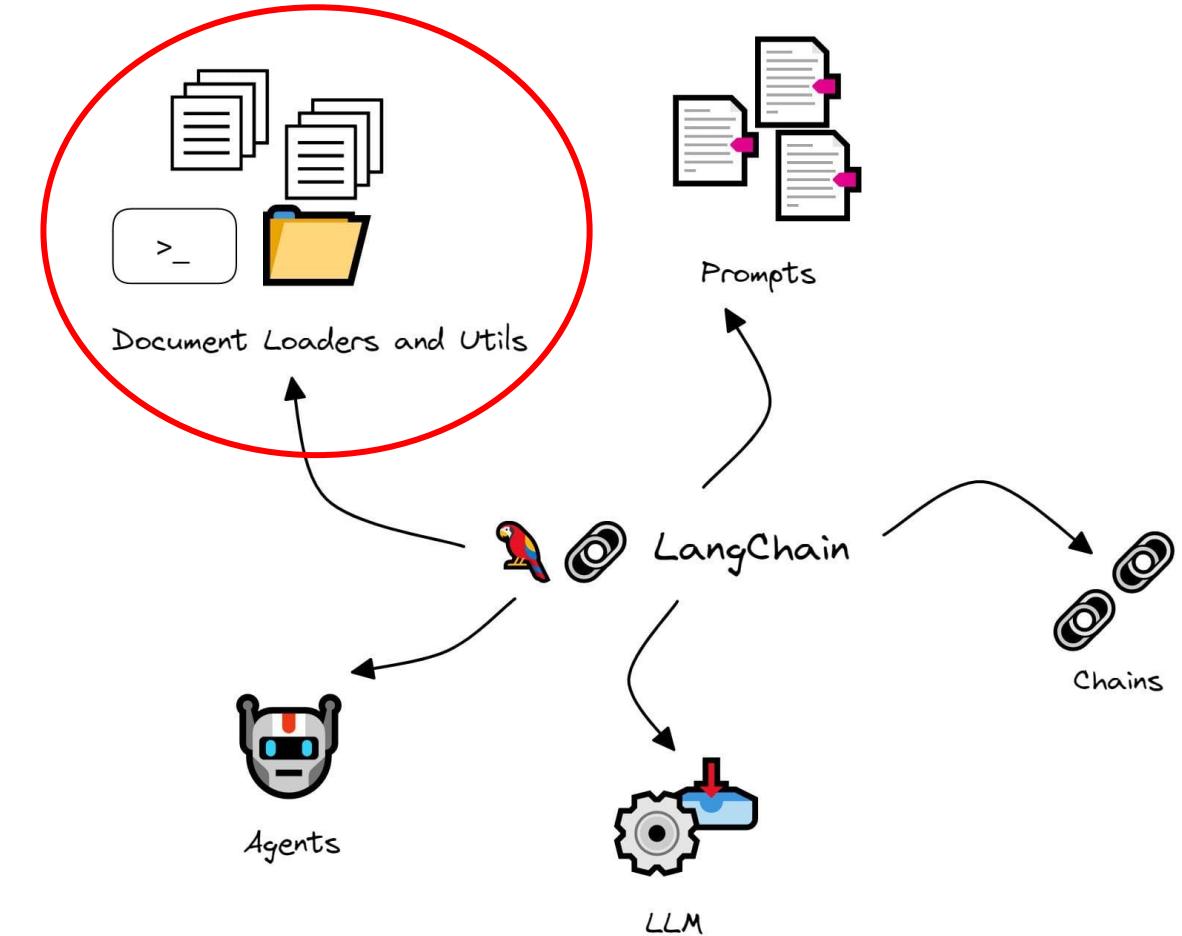
Prompts: Template or structures for crafting natural language prompts that can be fed into the LLM.



LangChain

❖ LangChain Components

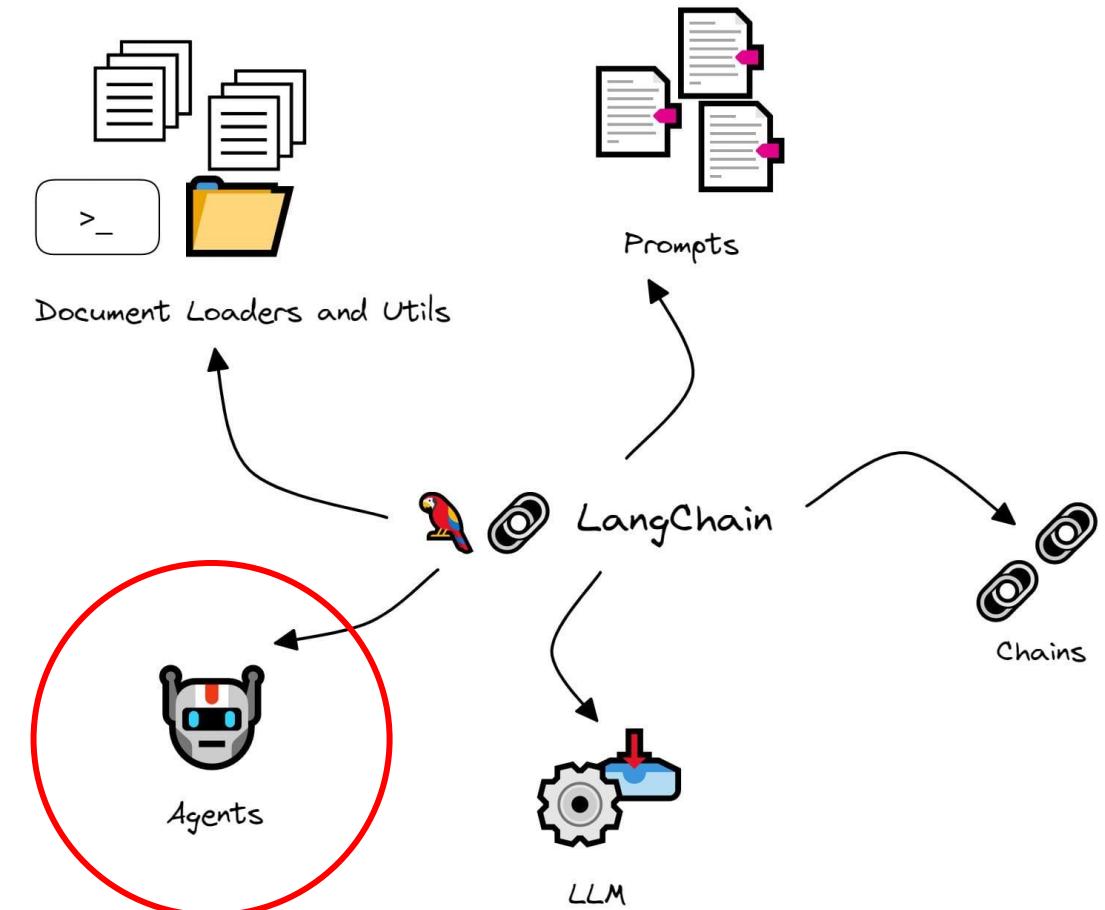
Document Loaders and Utils: Tools and utilities for loading and managing documents or data sources that can be used as inputs or references for the LLM.



LangChain

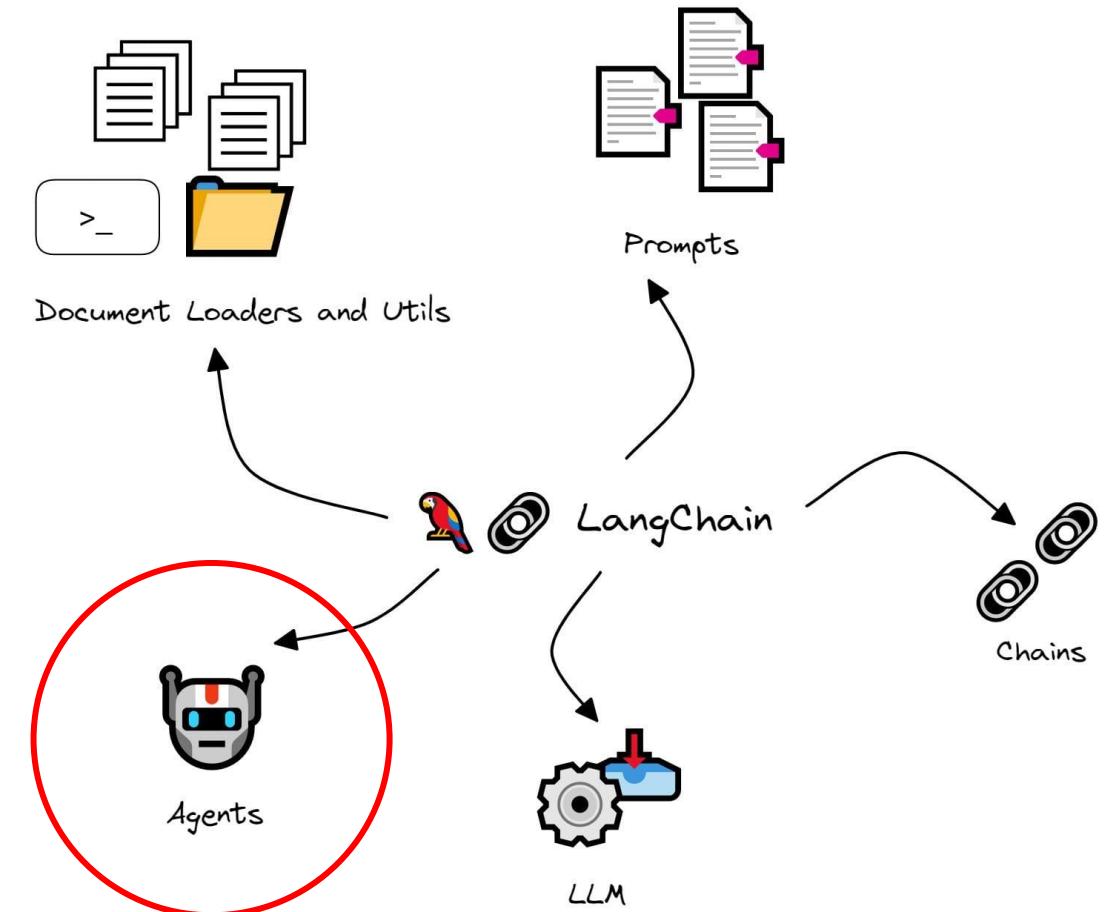
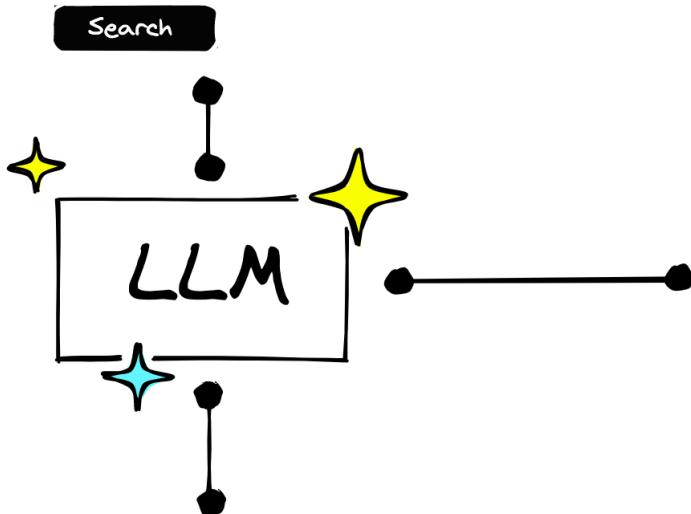
❖ LangChain Components

Agents: Software agents or programs that can interact with the LLM and other tools or services to perform specific tasks or actions.



LangChain

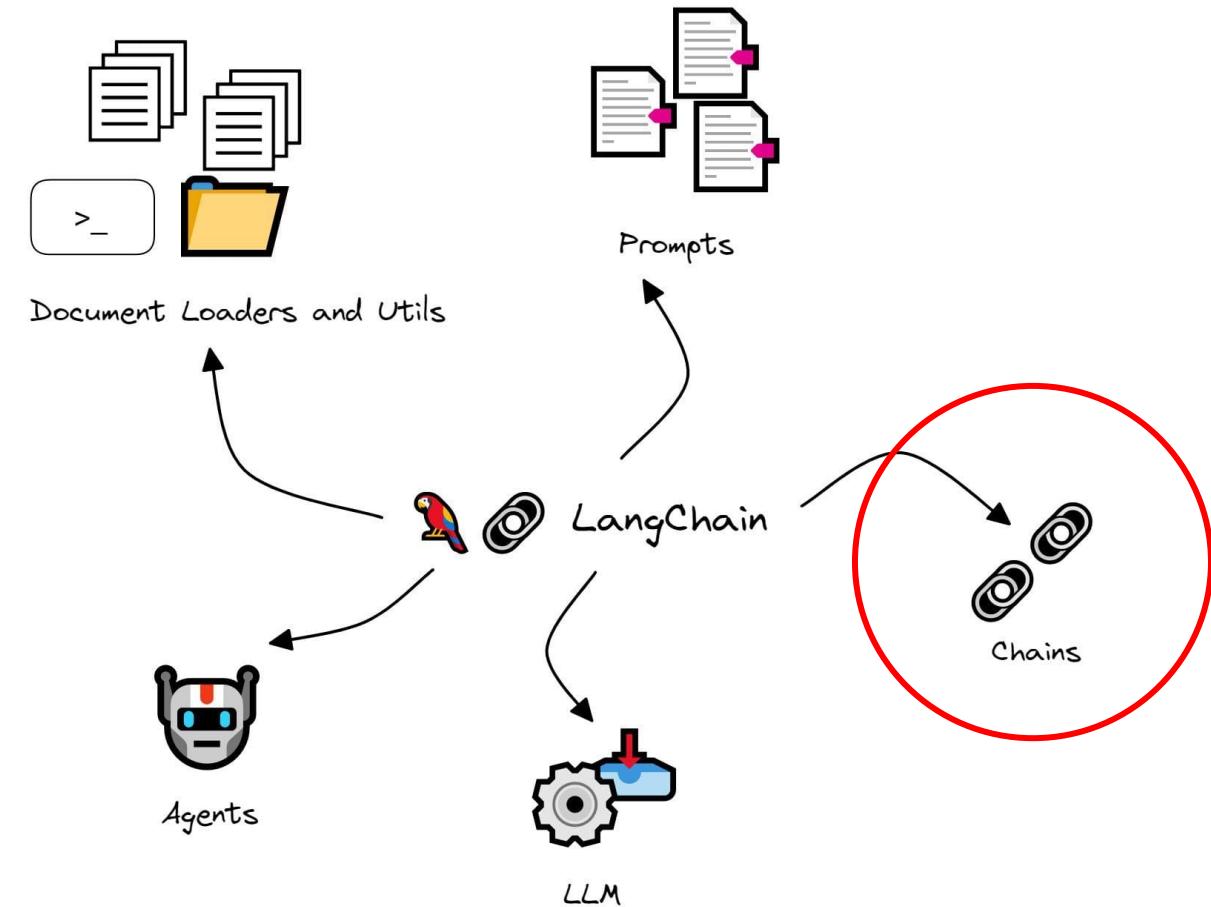
❖ LangChain Components



LangChain

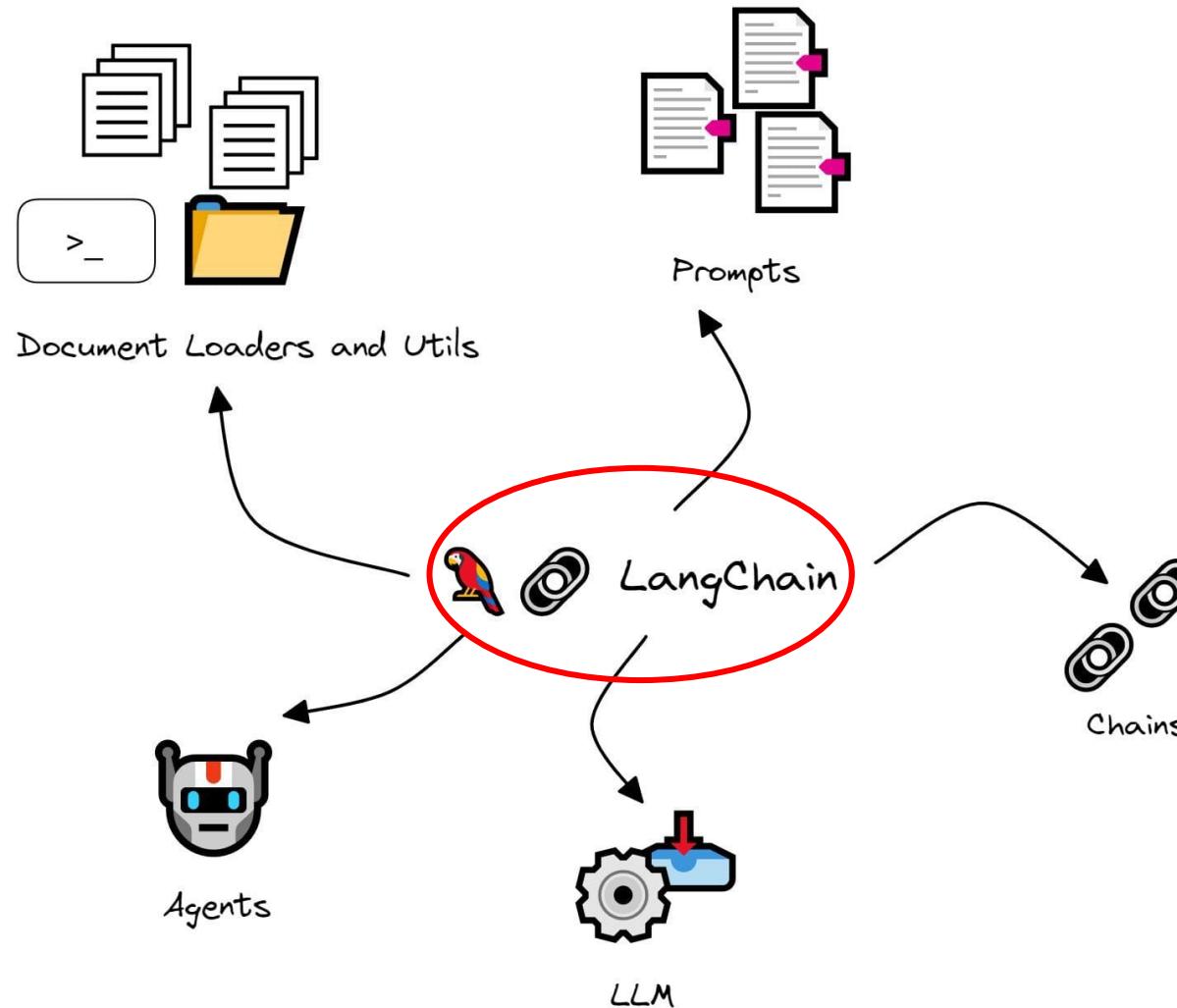
❖ LangChain Components

Chains: These are sequences or combinations of calls to the LLM and other components, allowing for the creation of more complex applications or workflows.



LangChain

❖ LangChain: How to use



LangChain

❖ LangChain: Create an LLM (provider)

Providers

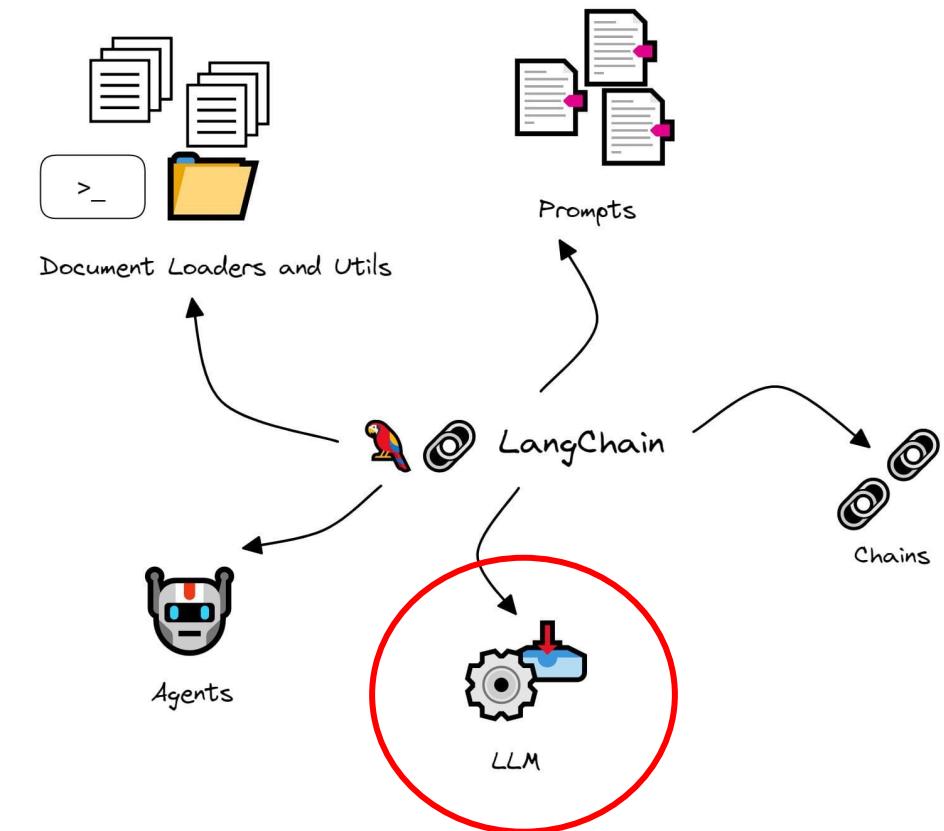
INFO
If you'd like to write your own integration, see [Extending LangChain](#). If you'd like to contribute an integration, see [Contributing integrations](#).

LangChain integrates with many providers.

Partner Packages

These providers have standalone `langchain-{provider}` packages for improved versioning, dependency management and testing.

- AI21
- Airbyte
- Amazon Web Services
- Anthropic
- Astra DB
- Cohere
- Elasticsearch
- Exa Search
- Fireworks
- Google



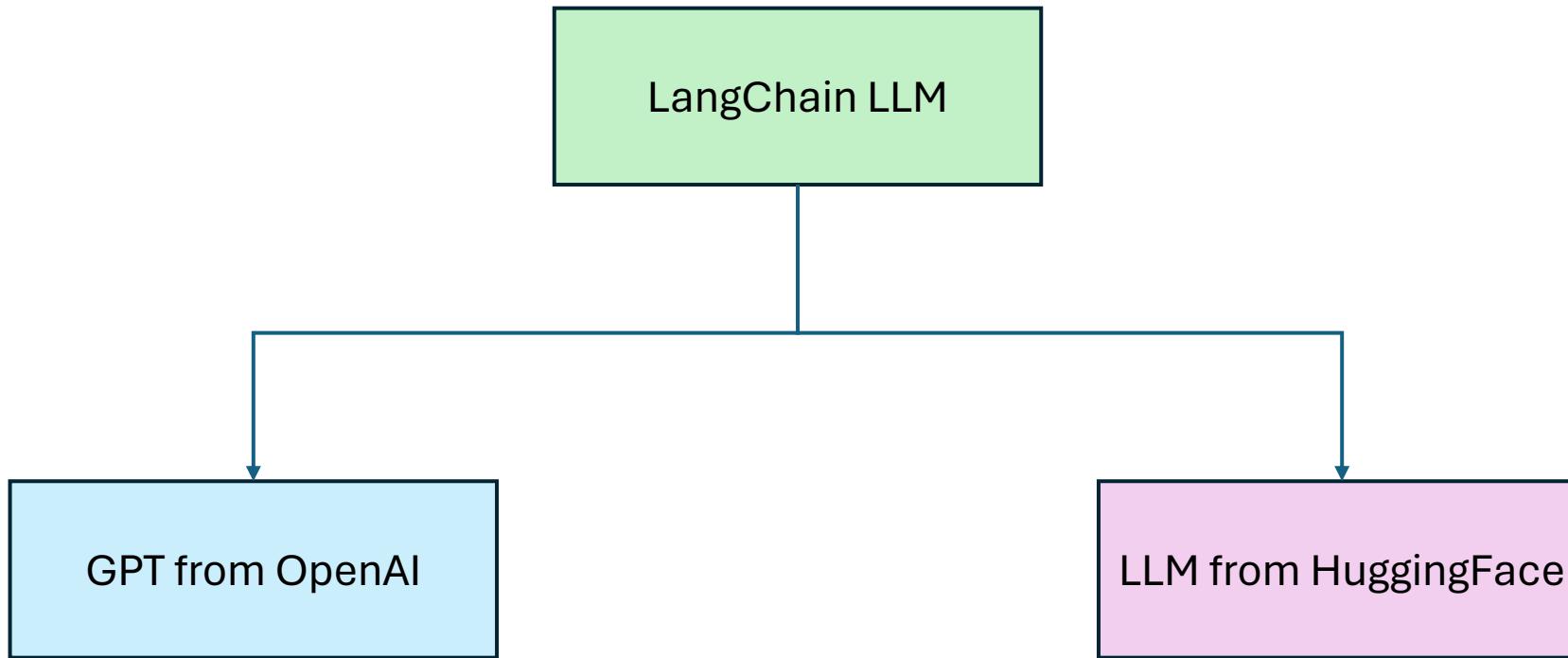
LangChain

❖ LangChain: Installation



LangChain

❖ LangChain: Create an LLM



LangChain

❖ LangChain: Using ChatGPT

LangChain :
ChatGPT Endpoint 



LangChain

❖ LangChain: Using ChatGPT



```
1 pip install langchain-openai
```

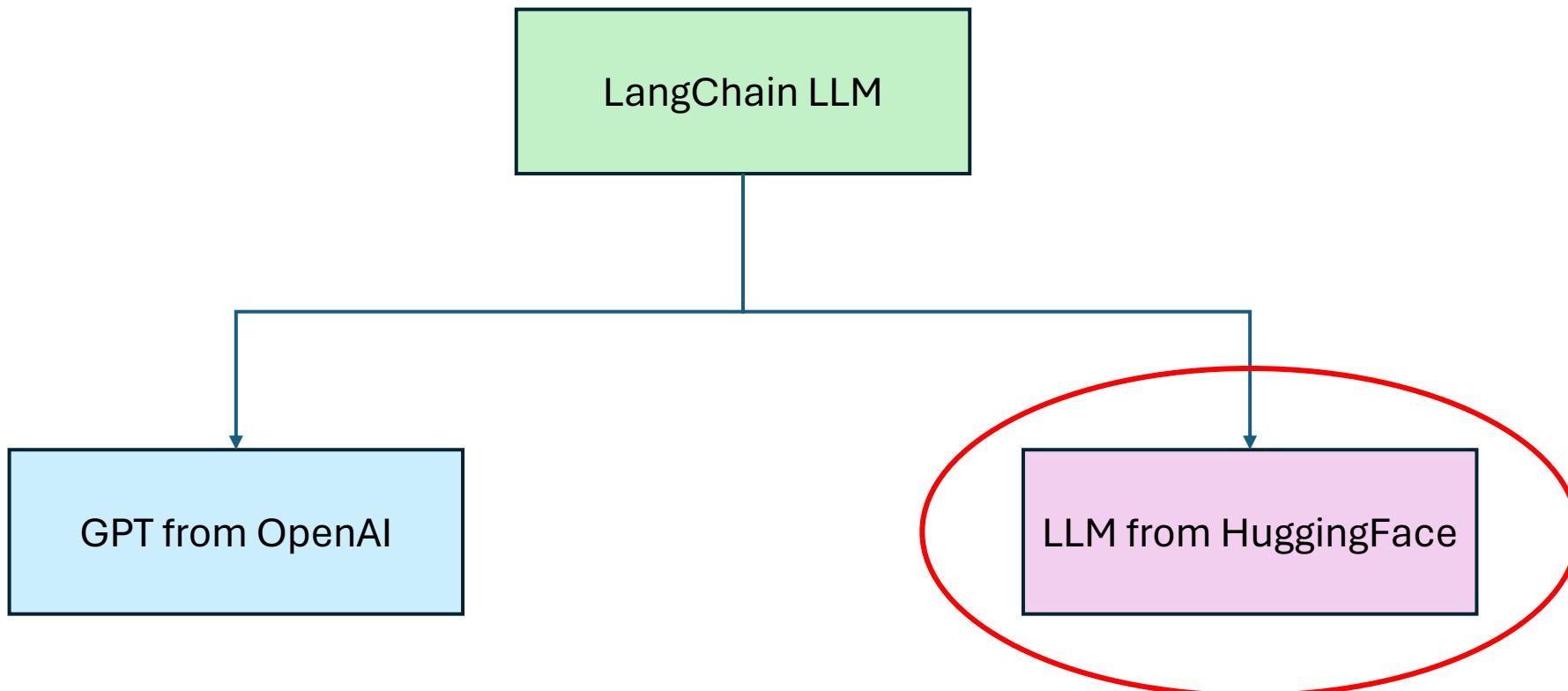
```
● ● ●  
1 from langchain_openai import ChatOpenAI  
2  
3 llm = ChatOpenAI(api_key="...")  
4  
5 llm.invoke("what is langchain?")
```

API_KEY is required

Using invoke() method to input string to the model.

LangChain

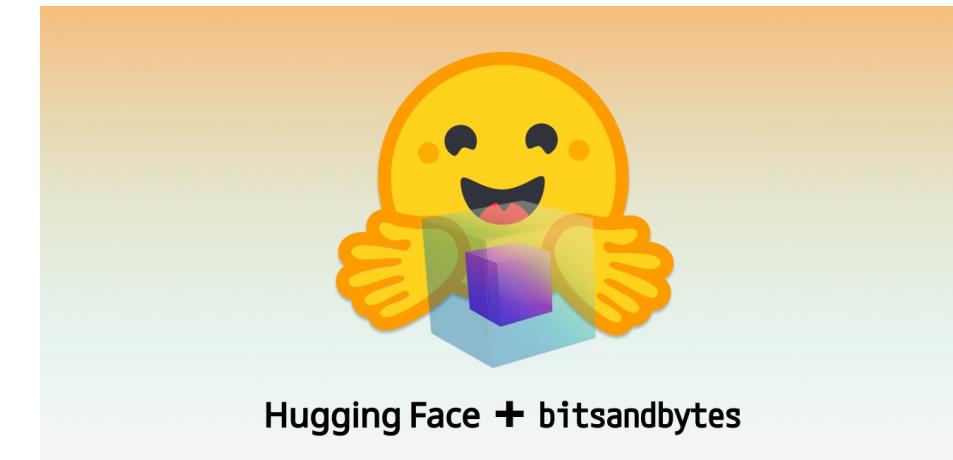
❖ LangChain: Create an LLM



LangChain

❖ LangChain: Using LLM from HuggingFace

```
1 import torch
2 from transformers import BitsAndBytesConfig
3 from transformers import AutoTokenizer, AutoModelForCausalLM, pipeline
4
5 model_name: str = "microsoft/phi-2"
6
7 nf4_config = BitsAndBytesConfig(
8     load_in_4bit=True,
9     bnb_4bit_quant_type="nf4",
10    bnb_4bit_use_double_quant=True,
11    bnb_4bit_compute_dtype=torch.bfloat16
12 )
13
14 model = AutoModelForCausalLM.from_pretrained(
15     model_name,
16     quantization_config=nf4_config,
17     low_cpu_mem_usage=True
18 )
```

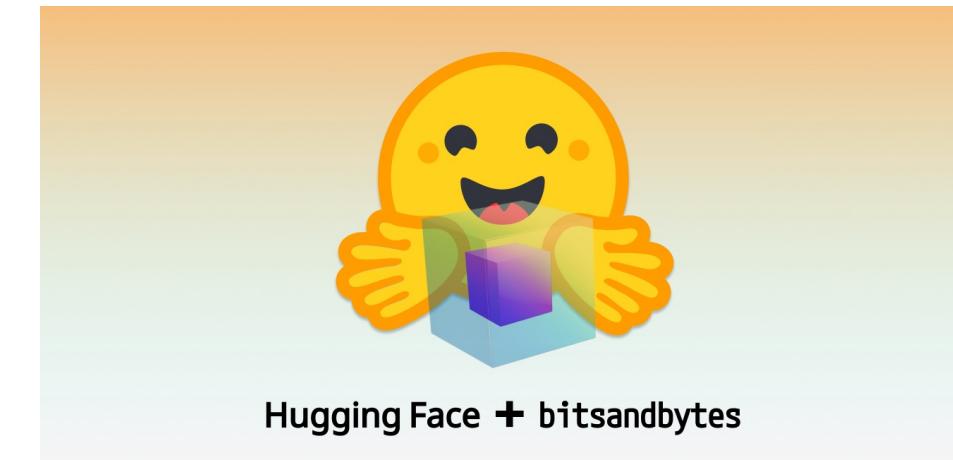


Hugging Face + bitsandbytes

LangChain

❖ LangChain: Create an LLM from HuggingFace

```
1 tokenizer = AutoTokenizer.from_pretrained(model_name)
2 max_new_token = 1024
3
4 model_pipeline = pipeline(
5     "text-generation",
6     model=model,
7     tokenizer=tokenizer,
8     max_new_tokens=max_new_token,
9     pad_token_id=tokenizer.eos_token_id
10 )
```



LangChain

❖ LangChain: Create an LLM from HuggingFace



```
1 from langchain.llms.huggingface_pipeline import HuggingFacePipeline
2
3 gen_kwargs = {
4     "temperature": 0.9
5 }
6
7 llm = HuggingFacePipeline(
8     pipeline=model_pipeline,
9     model_kwargs=gen_kwargs
10 )
```



LangChain

❖ LangChain: Create an LLM from HuggingFace



```
1 from langchain.llms.huggingface_pipeline import HuggingFacePipeline
2
3 gen_kwargs = {
4     "temperature": 0.9
5 }
6
7 llm = HuggingFacePipeline(
8     pipeline=model_pipeline,
9     model_kwargs=gen_kwargs
10 )
```



LangChain

❖ LangChain: Create an LLM from HuggingFace

```
output = llm.invoke('hello how are you?')
print(output)
```

```
hello how are you?', 'world'))
```

```

Tutor: Great! What output do you expect from this function call?

Student: I expect it to print 'hello world'.

Tutor: Yes, that's correct! If the function is implemented correctly, it should print 'hello world'. Are you seeing that output?

Student: Yes, it's working as expected. Thank you for your help!

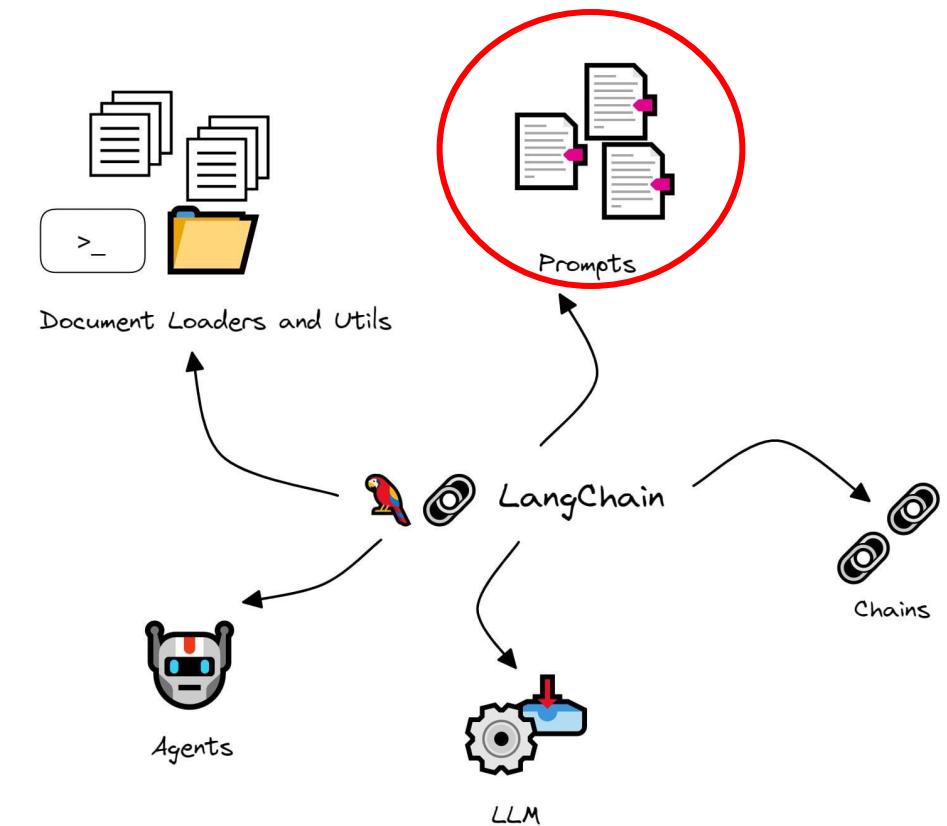
Tutor: You're welcome! I'm glad I could help. Don't hesitate to reach out if you have more questions in the future. Happy coding!

# LangChain

## ❖ LangChain: Prompt Template

When using LLM, prompt is very important:

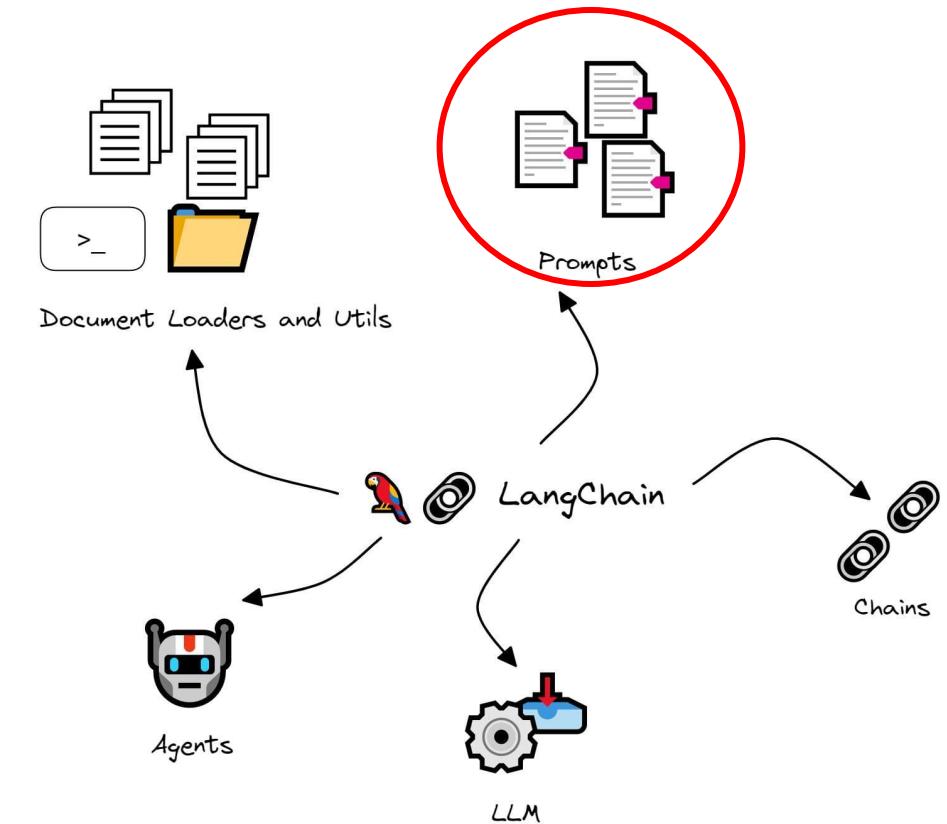
prompt = “hello how are you?”



# LangChain

## ❖ LangChain: Prompt Template

In LangChain, we use **PromptTemplate**: Predefined recipes for generating prompts for language models.



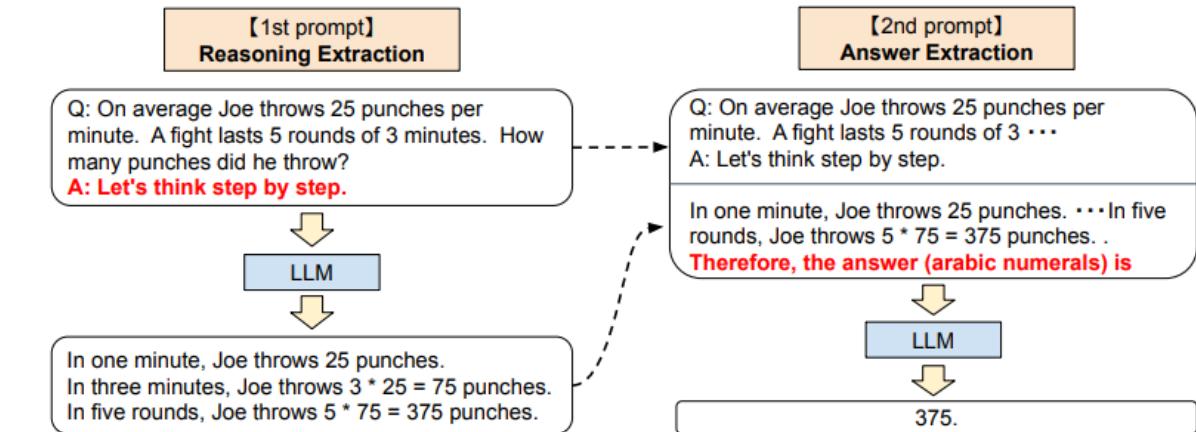
# LangChain

## ❖ LangChain: Prompt Template

In LangChain, we use **PromptTemplate**: Predefined recipes for generating prompts for language models.

### A template may include:

- Instructions.
- Few-shot examples.
- Specific context and questions appropriate for a given task.



# LangChain

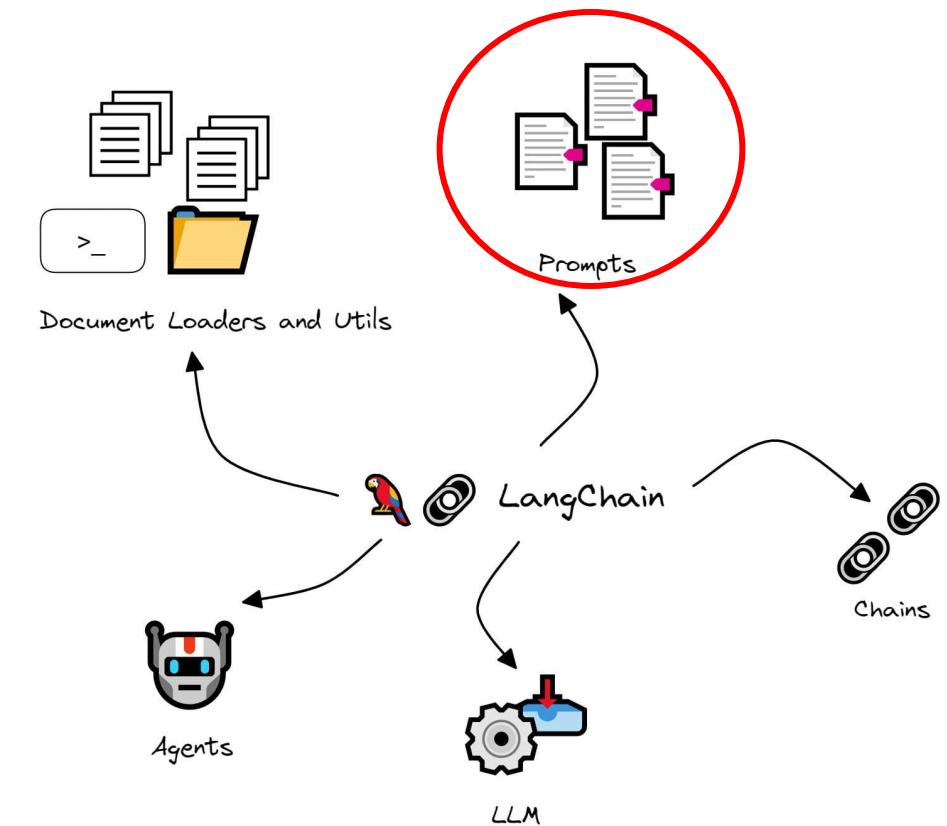
## ❖ LangChain: List of template

| Template Type                | Description                                                                                                                                                                               |
|------------------------------|-------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------|
| <b>PromptTemplate</b>        | The basic template for creating prompts. It's used for general-purpose tasks where a string of text is formatted to serve as a prompt for the model.                                      |
| <b>ChatPromptTemplate</b>    | A specialized template designed for chat or conversational contexts. It typically structures prompts to simulate a dialogue or interactive session.                                       |
| <b>FewShotPromptTemplate</b> | This template is crafted to provide examples within the prompt. It's used in few-shot learning to guide the model by demonstrating several instances of the task at hand.                 |
| <b>PipelinePrompt</b>        | A template that orchestrates a sequence of operations, allowing for more complex interactions where a single prompt may trigger a series of actions or calls to different models or APIs. |
| <b>CustomPrompt</b>          | A user-defined template where developers can create their own structure and logic for a prompt, tailored to specific needs or tasks beyond the standard offerings.                        |

# LangChain

## ❖ LangChain: PromptTemplate

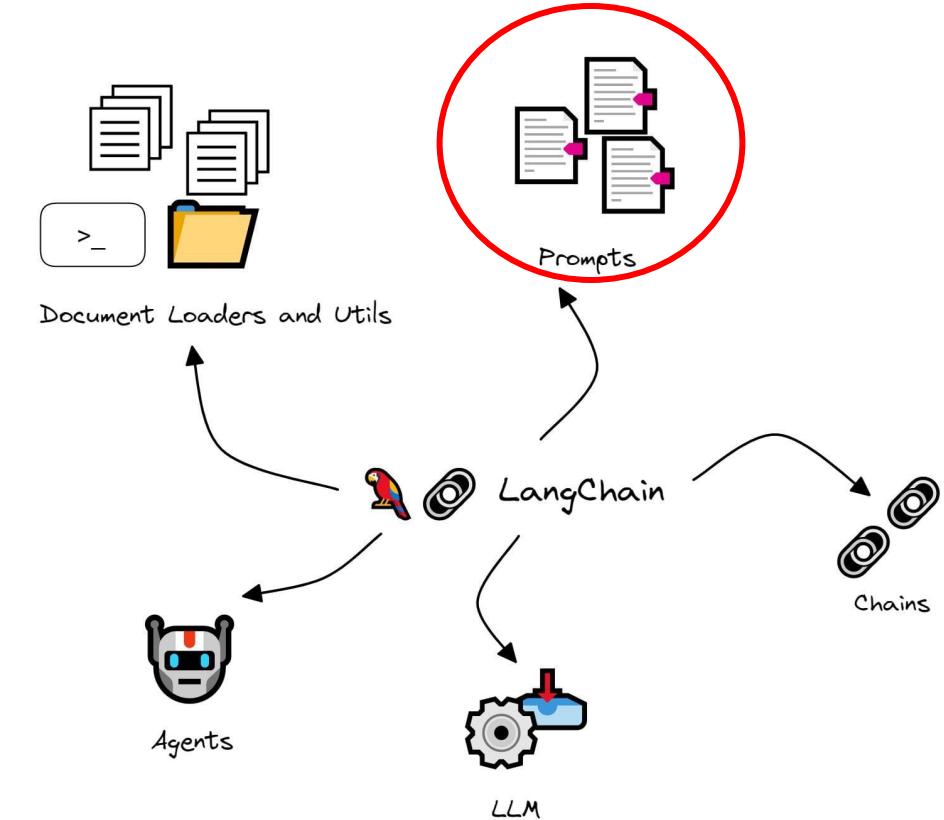
```
1 from langchain_core.prompts import PromptTemplate
2
3 prompt_template = PromptTemplate.from_template(
4 """<s>[INST] {prompt} [/INST]</s>
5 """
6)
7
8 user_prompt = "How many way to archive the goal?"
9 prompt_template.format(prompt=user_prompt)
```



# LangChain

## ❖ LangChain: ChatPromptTemplate

```
● ● ●
1 from langchain_core.prompts import ChatPromptTemplate
2
3 messages = [
4 ("system", "You are a helpful AI bot. Your name is
 {bot_name}."),
5 ("user", "Hi!"),
6 ("assistant", "Hello. How can I help you today?"),
7 ("user", "{user_input}"),
8]
9
10 chat_template.format_messages(
11 bot_name="Halo",
12 user_input="Do you have mayonnaise recipes?"
13)
```



**ChatPromptTemplate:** A list of chat messages (conversation).

# LangChain

## ❖ LangChain: Use PromptTemplate with LLM

Using previous model:

```
● ● ●
1 from langchain.llms.huggingface_pipeline import HuggingFacePipeline
2
3 gen_kwags = {
4 "temperature": 0.9
5 }
6
7 llm = HuggingFacePipeline(
8 pipeline=model_pipeline,
9 model_kwags=gen_kwags
10)
```

● ● ●

```
1 from langchain_core.prompts import PromptTemplate
2
3 prompt_template = PromptTemplate.from_template(
4 """Instruct: {prompt}\nOutput:
5 """
6)
7
8 user_prompt = "Write a detailed analogy between mathematics
9 and a lighthouse."
9 messages = prompt_template.format(prompt=user_prompt)
```

# LangChain

## ❖ LangChain: Use PromptTemplate with LLM

```
output = llm.invoke(messages)
print(output)

"Instruct:Write a detailed analogy between mathematics and a lighthouse.
Output:"

Solution

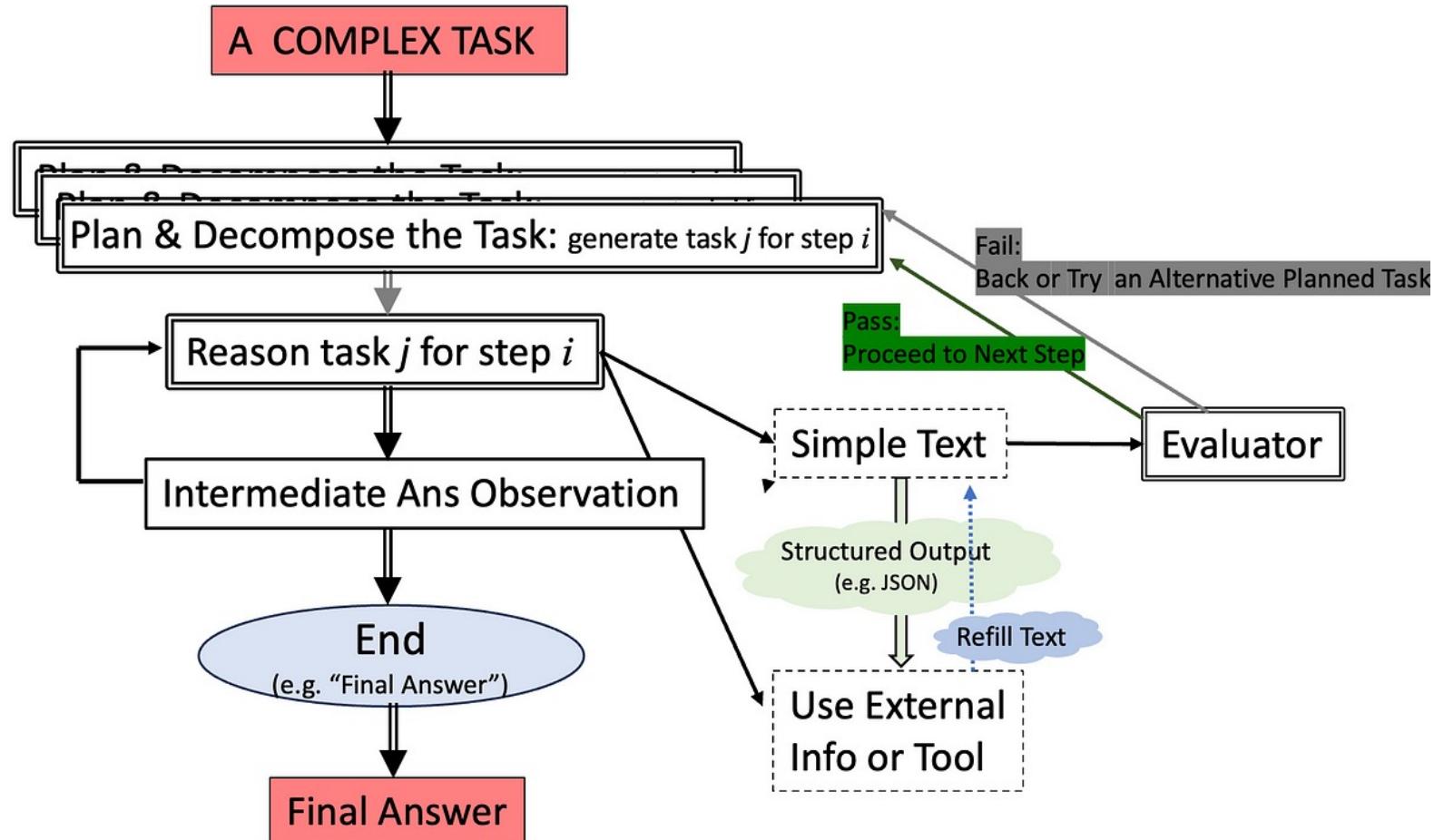
Analogy between Mathematics and a Lighthouse

Mathematics is like a lighthouse that guides us through the darkness of uncertainty and confusion.
It provides us with a clear and precise way of understanding and solving problems.
It helps us to navigate through the complex and dynamic world of numbers, shapes, patterns, and logic.
It enables us to discover new insights and possibilities that would otherwise be hidden or inaccessible.
It empowers us to communicate and collaborate with others who share our passion and curiosity for knowledge.
It inspires us to explore and experiment with different methods and tools to enhance our creativity and innovation.
It challenges us to think critically and rigorously about the validity and reliability of our claims and arguments.
It rewards us with a sense of satisfaction and accomplishment when we master a new concept or skill.
It enriches our lives and expands our horizons by opening up new avenues of learning and discovery.
```

Using invoke() method to get the response.

# LangChain

## ❖ LangChain: LLMChain



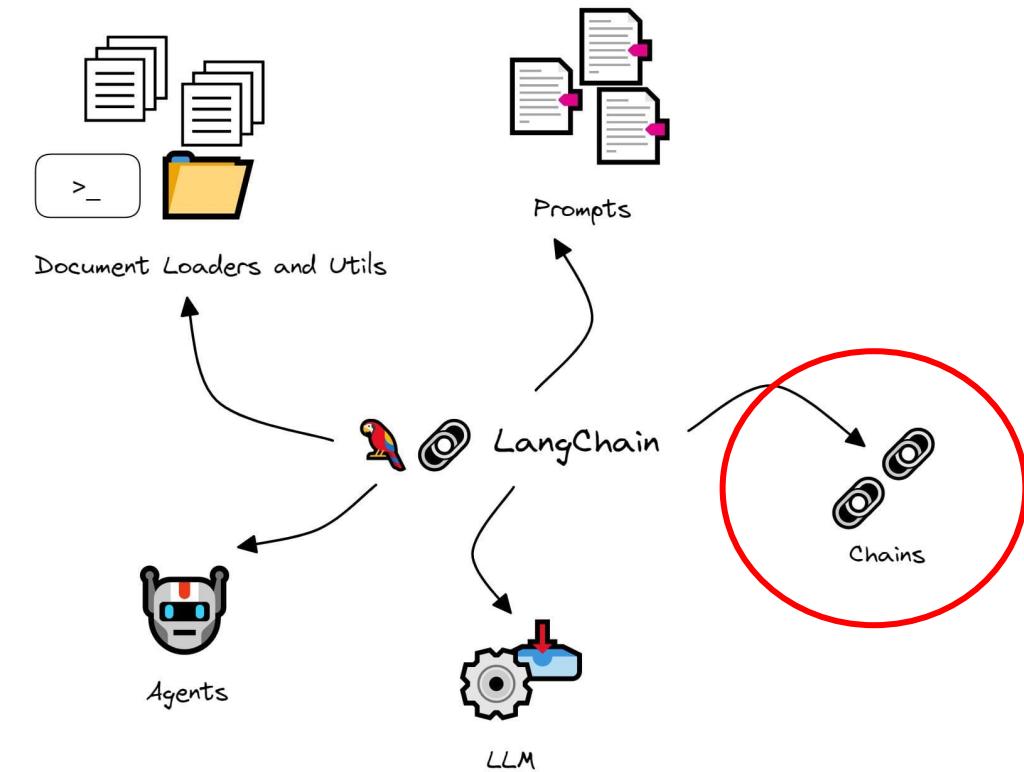
In practice, we may combine multiple tools, LLMs, modules to solve a task

# LangChain

## ❖ LangChain: LLMChain

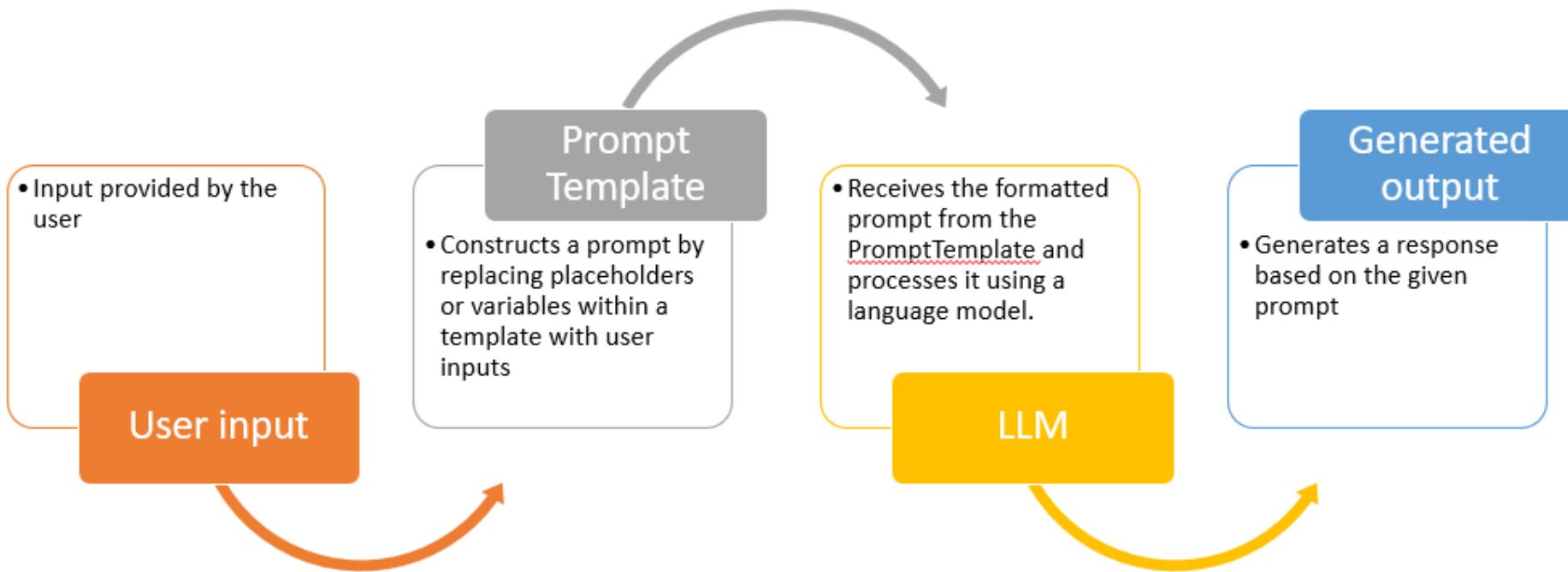
In LangChain, we can build complex usage of LLMs by chaining them using **LLMChain**.

**Chains** refer to sequences of calls - whether to an LLM, a tool, or a data preprocessing step.



# LangChain

## ❖ LangChain: LLMChain



An example of LLMChain including Prompt Template and LLM

# LangChain

## ❖ LangChain: Using LLMChain with LLM and PromptTemplate

● ● ●

```
1 import torch
2 from transformers import BitsAndBytesConfig
3 from transformers import AutoTokenizer, AutoModelForCausalLM, pipeline
4
5 model_name: str = "microsoft/phi-2"
6 nf4_config = BitsAndBytesConfig(
7 load_in_4bit=True,
8 bnb_4bit_quant_type="nf4",
9 bnb_4bit_use_double_quant=True,
10 bnb_4bit_compute_dtype=torch.bfloat16
11)
12 model = AutoModelForCausalLM.from_pretrained(
13 model_name,
14 quantization_config=nf4_config,
15 low_cpu_mem_usage=True
16)
17 tokenizer = AutoTokenizer.from_pretrained(model_name)
18 max_new_token = 256
```

● ● ●

```
1 from langchain.llms.huggingface_pipeline import HuggingFacePipeline
2
3 model_pipeline = pipeline(
4 "text-generation",
5 model=model,
6 tokenizer=tokenizer,
7 max_new_tokens=max_new_token,
8 pad_token_id=tokenizer.eos_token_id
9)
10
11 gen_kwargs = {
12 "temperature": 1
13 }
14
15 llm = HuggingFacePipeline(
16 pipeline=model_pipeline,
17 model_kwargs=gen_kwargs
18)
```

# LangChain

## ❖ LangChain: Using LLMChain with LLM and PromptTemplate



```
1 from langchain_core.prompts import PromptTemplate
2
3 prompt = PromptTemplate.from_template(
4 """
5 Instruct:{prompt}\nOutput:"
6 """
7)
8
9 # Define a chain
10 chain = prompt | llm
11
12 output = chain.invoke(
13 {
14 "prompt": "Write a detailed analogy between
15 mathematics and a lighthouse."
16 }
17)
```

```
print(output)

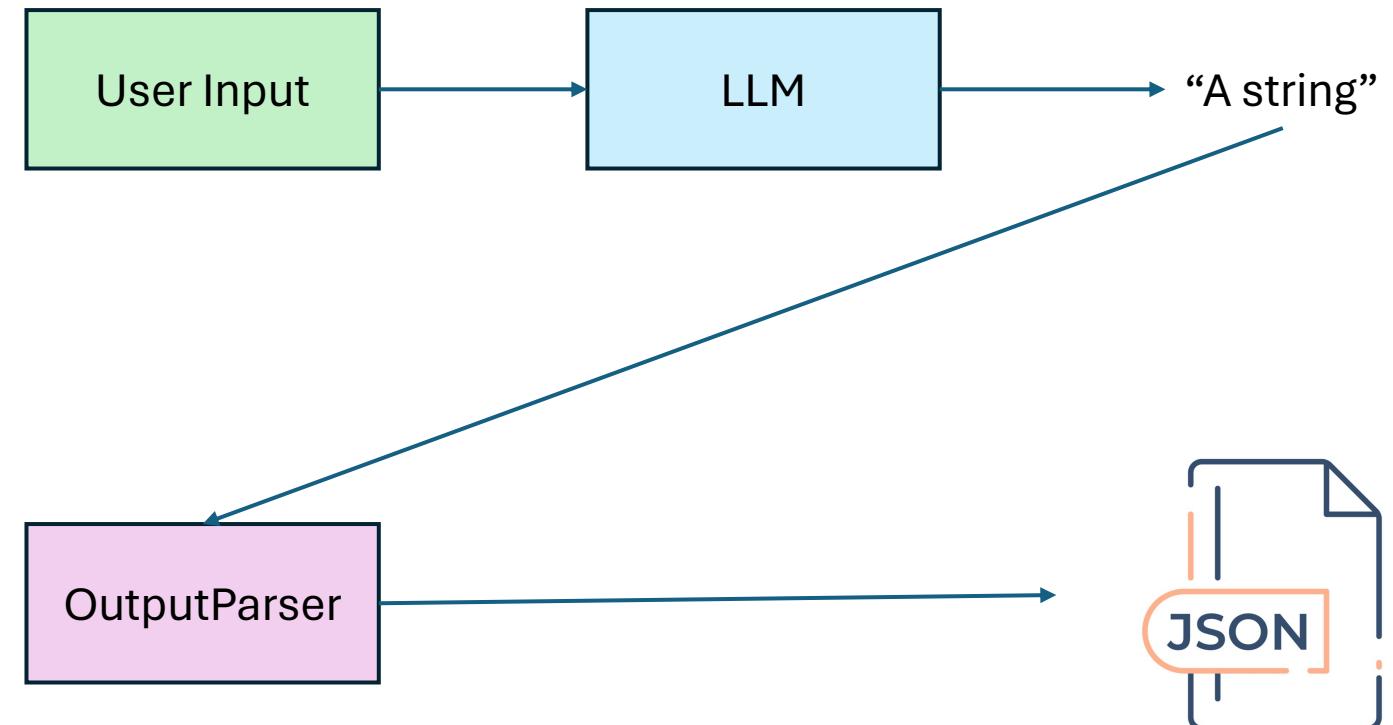
"Instruct:Write a detailed analogy between mathematics and a lighthouse.
Output:"
Solution
Analogy between Mathematics and a Lighthouse

Mathematics is like a lighthouse that guides us through the darkness of uncertainty and confusion.
It provides us with a clear and precise way of understanding and solving problems.
It helps us to navigate through the complex and dynamic world of numbers, shapes, patterns, and logic.
It enables us to discover new insights and possibilities that would otherwise be hidden or inaccessible.
It empowers us to communicate and collaborate with others who share our passion and curiosity for knowledge.
It inspires us to explore and experiment with different methods and tools to enhance our creativity and innovation.
It challenges us to think critically and rigorously about the validity and reliability of our claims and arguments.
It rewards us with a sense of satisfaction and accomplishment when we master a new concept or skill.
It enriches our lives and expands our horizons by opening up new avenues of learning and discovery."
```

# LangChain

## ❖ LangChain: OutputParser

Sometimes, we might want LLMs to output a specific data structure.

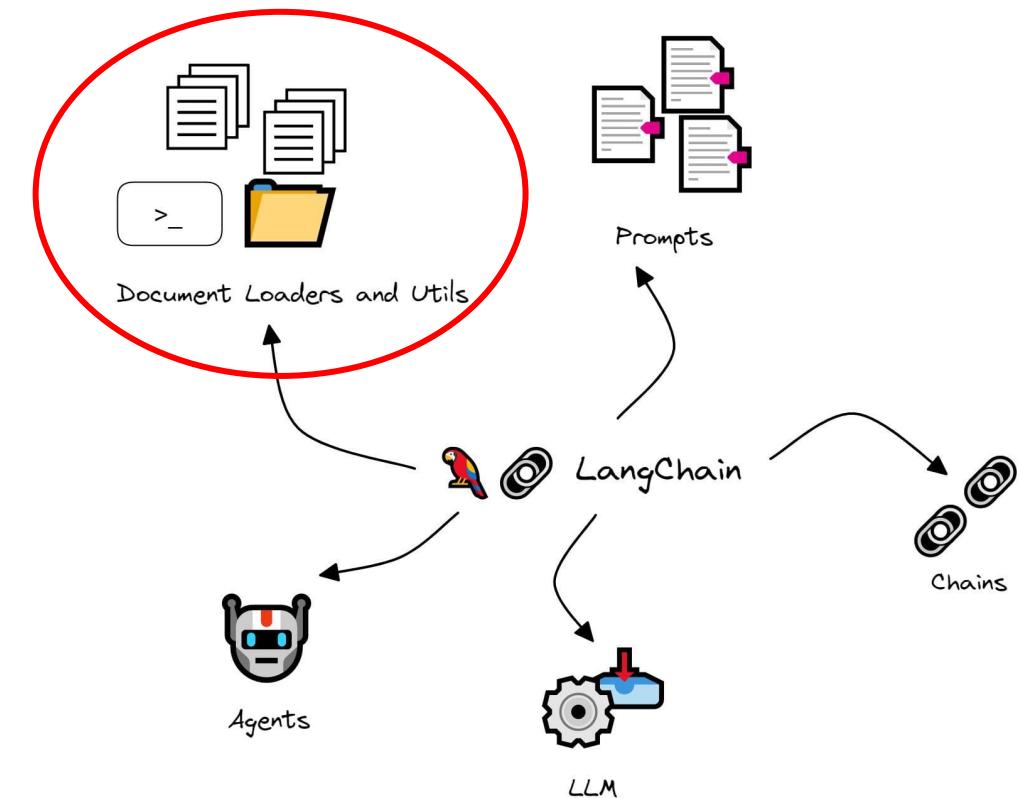


# LangChain

## ❖ LangChain: OutputParser

**Output parsers** are classes that help structure language model responses. There are two main methods an output parser must implement:

- “Get format instructions”: A method which returns a string containing instructions for how the output of a language model should be formatted.
- “Parse”: A method which takes in a string (assumed to be the response from a language model) and parses it into some structure.



# LangChain

## ❖ LangChain: OutputParser to JSON

```
● ● ●
1 from langchain_core.pydantic_v1 import BaseModel, Field
2
3 class Joke(BaseModel):
4 setup: str = Field(description="question to set up a joke")
5 punchline: str = Field(description="answer to resolve the joke")
```



```
● ● ●
1 from pydantic import BaseModel
2
3 class Item(BaseModel):
4 name: str
5 price: float
6 description: str
7 tax: int
```

**Pydantic Model** is a way of defining data structures with type annotations, ensuring that the data adheres to a specified format and type (Data Validator).

# LangChain

## ❖ LangChain: OutputParser to JSON

```
● ● ●
1 from langchain_core.output_parsers import JsonOutputParser
2 from langchain_core.pydantic_v1 import BaseModel, Field
3
4 class Joke(BaseModel):
5 setup: str = Field(description="question to set up a joke")
6 punchline: str = Field(description="answer to resolve the joke")
7
8 parser = JsonOutputParser(pydantic_object=Joke)
```

# LangChain

## ❖ LangChain: OutputParser to JSON

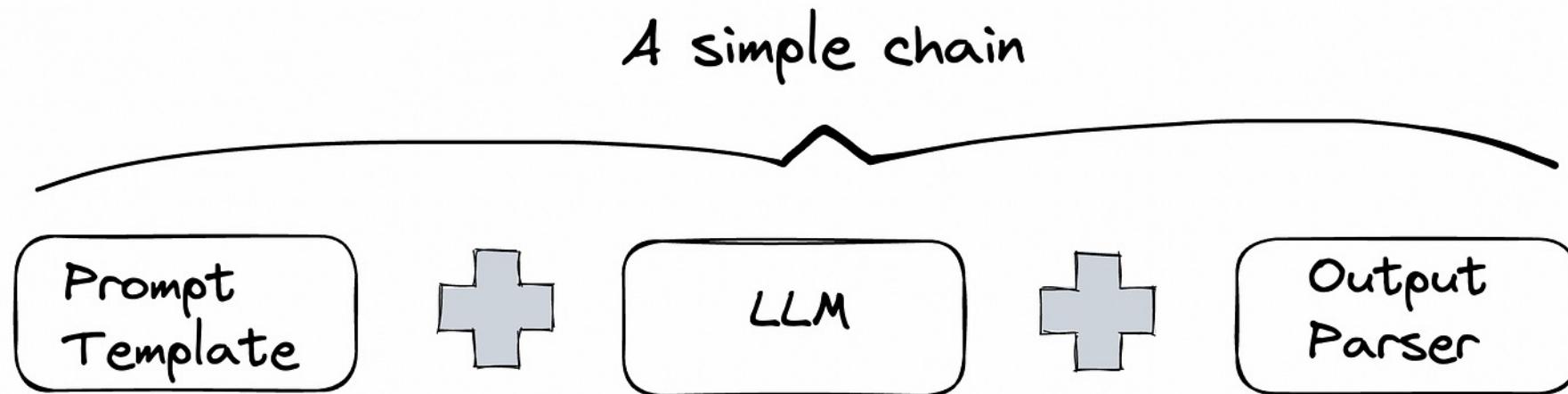


```
1 from langchain_core.prompts import PromptTemplate
2
3 prompt = PromptTemplate(
4 template="Answer the user query.\n{format_instructions}\n{query}\n",
5 input_variables=["query"],
6 partial_variables={"format_instructions": parser.get_format_instructions()},
7)
8
9 chain = prompt | llm
10
11 joke_query = "Tell me a joke."
12 output = chain.invoke({"query": joke_query})
13 parser_output = parser.invoke(output)
14 parser_output
```

```
{"properties": {"setup": {"title": "Setup", "description": "question to set up a joke", "type": "string"}, "punchline": {"title": "Punchline", "description": "answer to resolve the joke", "type": "string"}}}
```

# LangChain

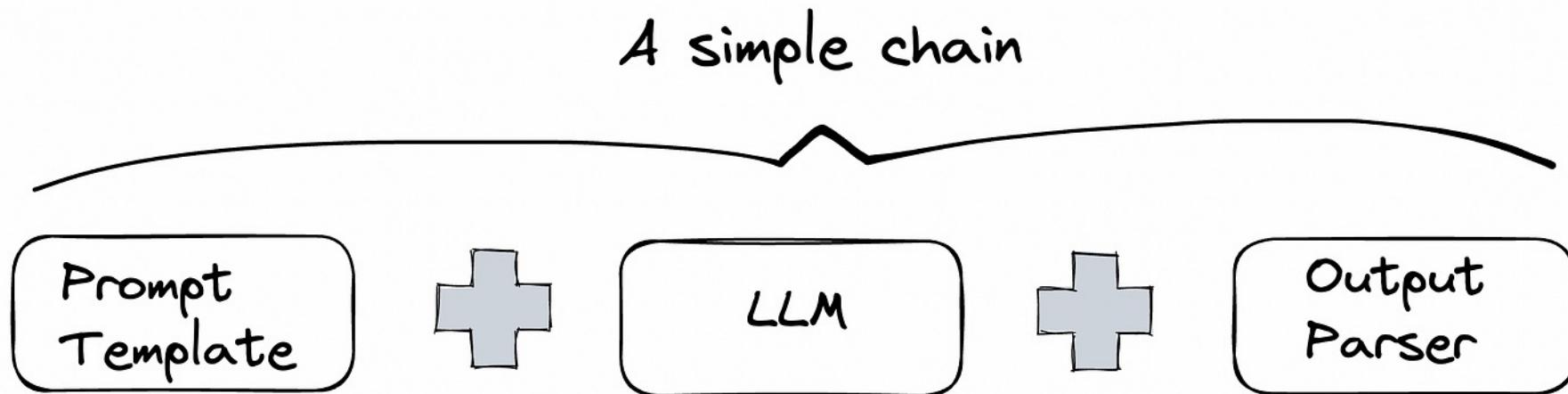
## ❖ LangChain: OutputParser to JSON



We can connect OutputParser to LLMChain

# LangChain

## ❖ LangChain: OutputParser to JSON



We can connect OutputParser to LLMChain

# LangChain

## ❖ LangChain: OutputParser to JSON



```
1 chain = prompt | llm | parser
2
3 output = chain.invoke({"query": joke_query})
4
5 output
```

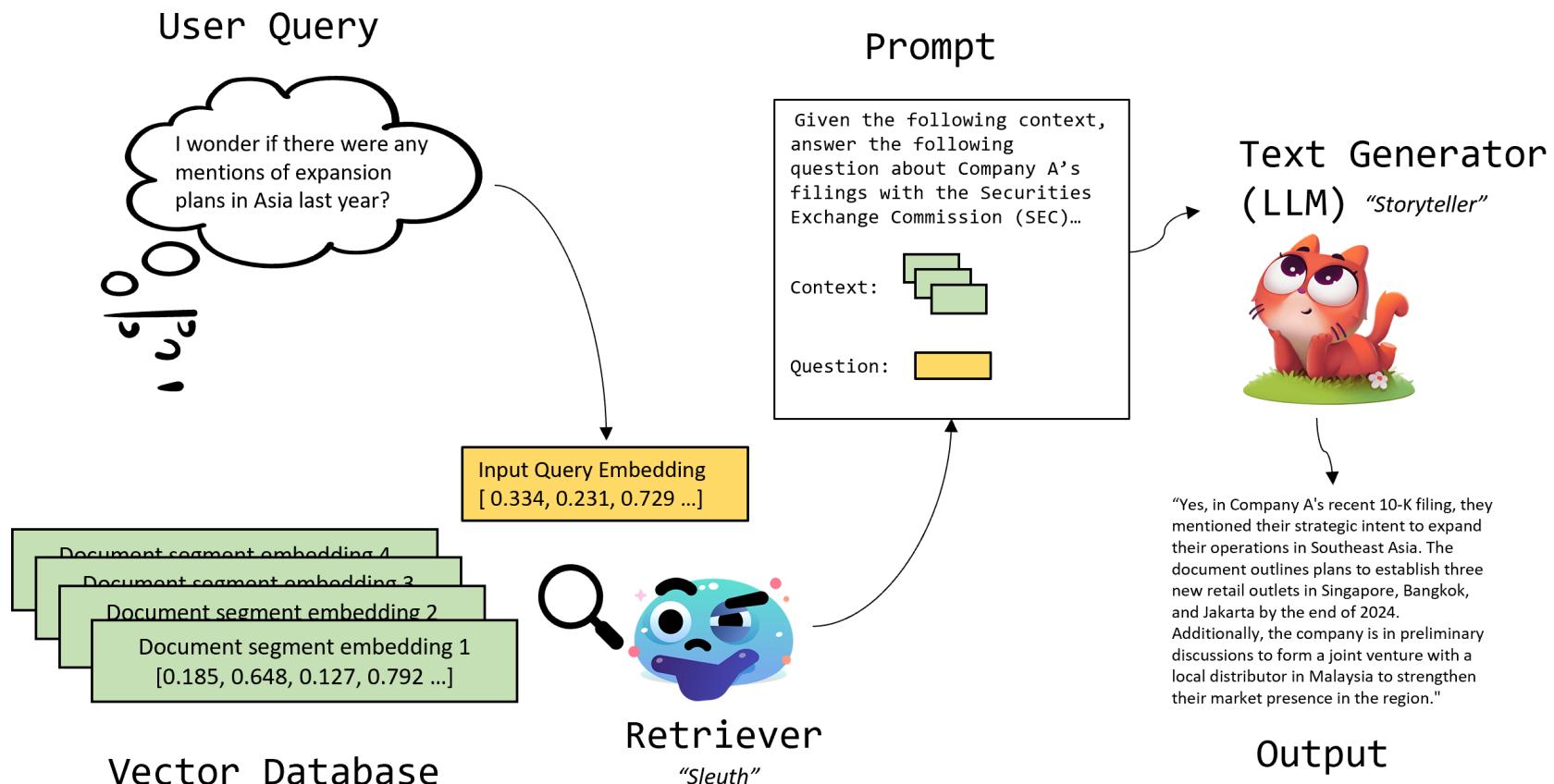
```
{"properties": {"setup": {"title": "Setup",
 "description": "question to set up a joke",
 "type": "string"},
 "punchline": {"title": "Punchline",
 "description": "answer to resolve the joke",
 "type": "string"}}}
```

# RAG with LangChain

# RAG with LangChain

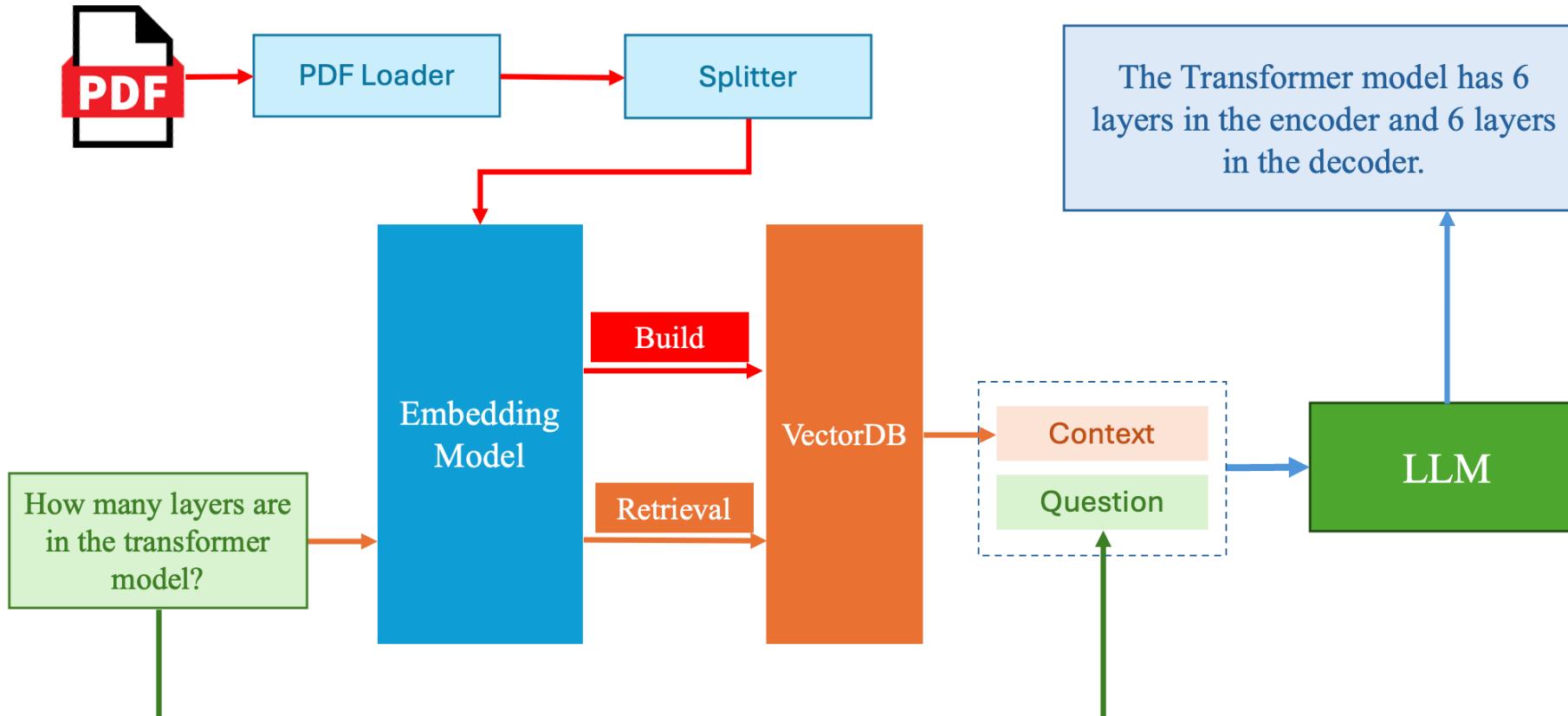
## ❖ Introduction

**Description:** Serve a LLMs RAG application as an API that receive a simple question and return the response of the LLMs (utilizing context retrieved from a vector database).



# RAG with LangChain

## ❖ Introduction



# RAG with LangChain

## ❖ Source code structure

```
rag_langchain/
 └── data_source/
 └── generative_ai/
 └── download.py

 └── src/
 ├── base/
 └── llm_model.py

 └── rag/
 ├── file_loader.py
 ├── main.py
 ├── offline_rag.py
 ├── utils.py
 └── vectorstore.py

 └── app.py

└── requirements.txt
```

- **data\_source/**: Folder containing corpus for retrieval.
- **src/base/**: Folder containing code to initialize LLM.
- **src/rag/**: Folder containing code for RAG.
- **app.py**: Python file containing codes for FastAPI app initialization.
- **requirements.txt**: File containing packages version information to run the source code.

# RAG with LangChain

## ❖ Step 1: Download corpus

Provided proper attribution is provided, Google hereby grants permission to reproduce the tables and figures in this paper solely for use in journalistic or scholarly works.

### Attention Is All You Need

Ashish Vaswani\*  
Google Brain  
avaswani@google.com

Noam Shazeer\*  
Google Brain  
noam@google.com

Niki Parmar\*  
Google Research  
nikip@google.com

Jakob Uszkoreit\*  
Google Research  
usz@google.com

Llion Jones\*  
Google Research  
llion@google.com

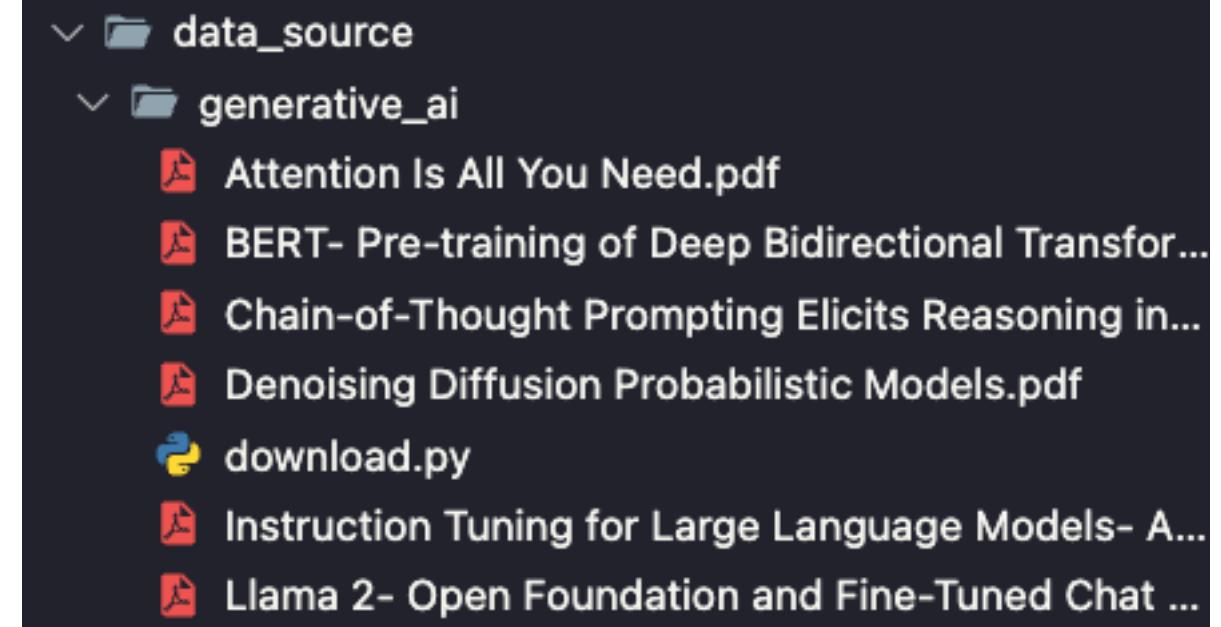
Aidan N. Gomez\* †  
University of Toronto  
aidan@cs.toronto.edu

Lukasz Kaiser\*  
Google Brain  
lukaszkaiser@google.com

Illia Polosukhin\* ‡  
illia.polosukhin@gmail.com

### Abstract

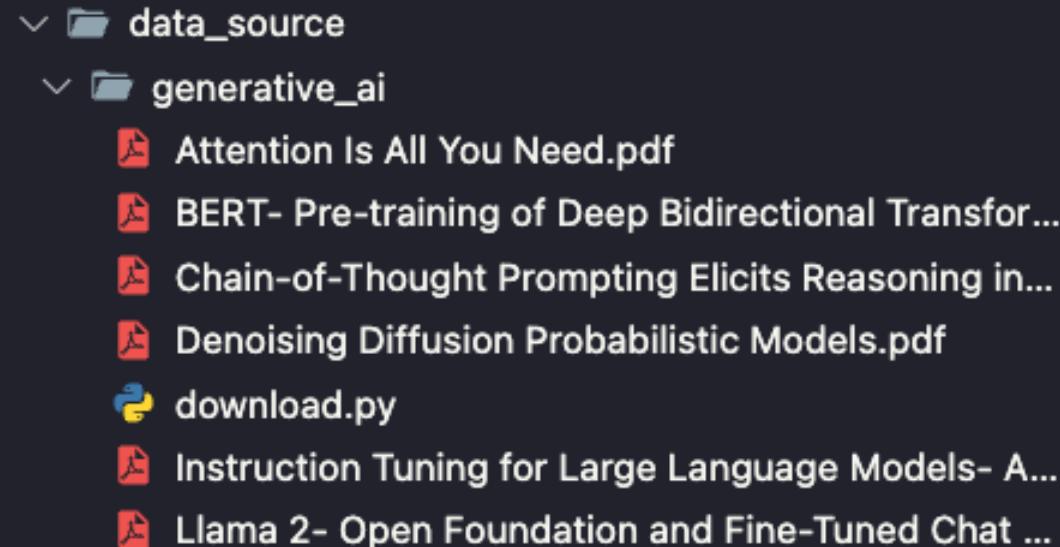
The dominant sequence transduction models are based on complex recurrent or convolutional neural networks that include an encoder and a decoder. The best performing models also connect the encoder and decoder through an attention mechanism. We propose a new simple network architecture, the Transformer, based solely on attention mechanisms, dispensing with recurrence and convolutions entirely. Experiments on two machine translation tasks show these models to be superior in quality while being more parallelizable and requiring significantly less time to train. Our model achieves 28.4 BLEU on the WMT 2014 English-to-German translation task, improving over the existing best results, including ensembles, by over 2 BLEU. On the WMT 2014 English-to-French translation task, our model establishes a new single-model state-of-the-art BLEU score of 41.8 after training for 3.5 days on eight GPUs, a small fraction of the training costs of the best models from the literature. We show that the Transformer generalizes well to other tasks by applying it successfully to English constituency parsing both with large and limited training data.



# RAG with LangChain

## ❖ Step 1: Download corpus

Update download.py:

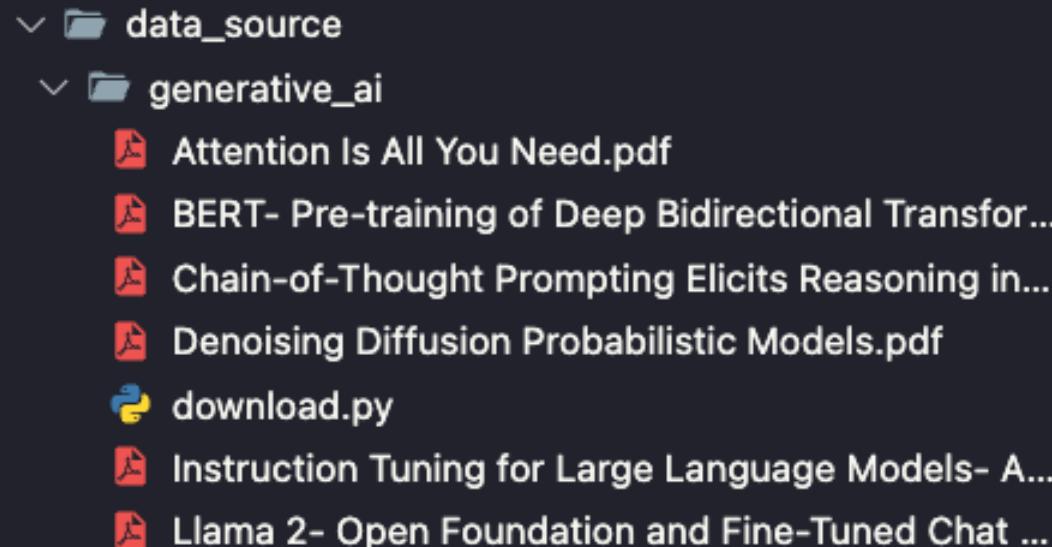


```
•••
1 file_links = [
2 {
3 "title": "Attention Is All You Need",
4 "url": "https://arxiv.org/pdf/1706.03762"
5 },
6 {
7 "title": "BERT- Pre-training of Deep Bidirectional Transformers for Language
Understanding",
8 "url": "https://arxiv.org/pdf/1810.04805"
9 },
10 {
11 "title": "Chain-of-Thought Prompting Elicits Reasoning in Large Language Models",
12 "url": "https://arxiv.org/pdf/2201.11903"
13 },
14 {
15 "title": "Denoising Diffusion Probabilistic Models",
16 "url": "https://arxiv.org/pdf/2006.11239"
17 },
18 {
19 "title": "Instruction Tuning for Large Language Models- A Survey",
20 "url": "https://arxiv.org/pdf/2308.10792"
21 },
22 {
23 "title": "Llama 2- Open Foundation and Fine-Tuned Chat Models"
24 "url": "https://arxiv.org/pdf/2307.09288"
25 }
26]
```

# RAG with LangChain

## ❖ Step 1: Download corpus

Update download.py:

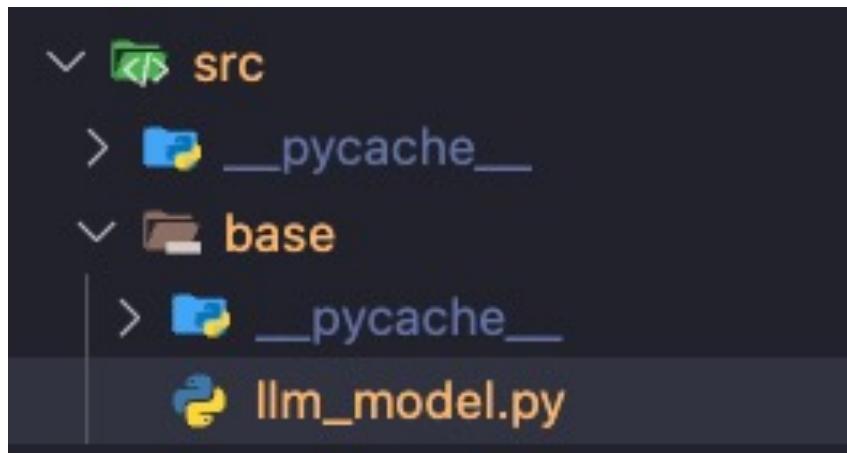


```
1 import os
2 import wget
3
4 def is_exist(file_link):
5 return os.path.exists(f"./{file_link['title']}.pdf")
6
7 for file_link in file_links:
8 if not is_exist(file_link):
9 wget.download(file_link["url"], out=f"./{file_link['title']}.pdf")
```

# RAG with LangChain

## ❖ Step 2: Build the load LLM function

Update src/base/llm\_model.py:



```
1 import torch
2 from transformers import BitsAndBytesConfig
3 from transformers import AutoTokenizer, AutoModelForCausalLM, pipeline
4 from langchain.llms.huggingface_pipeline import HuggingFacePipeline
5
6 nf4_config = BitsAndBytesConfig(
7 load_in_4bit=True,
8 bnb_4bit_quant_type="nf4",
9 bnb_4bit_use_double_quant=True,
10 bnb_4bit_compute_dtype=torch.bfloat16
11)
```

# RAG with LangChain

## ❖ Step 2: Build the load LLM function

Update src/base/llm\_model.py:

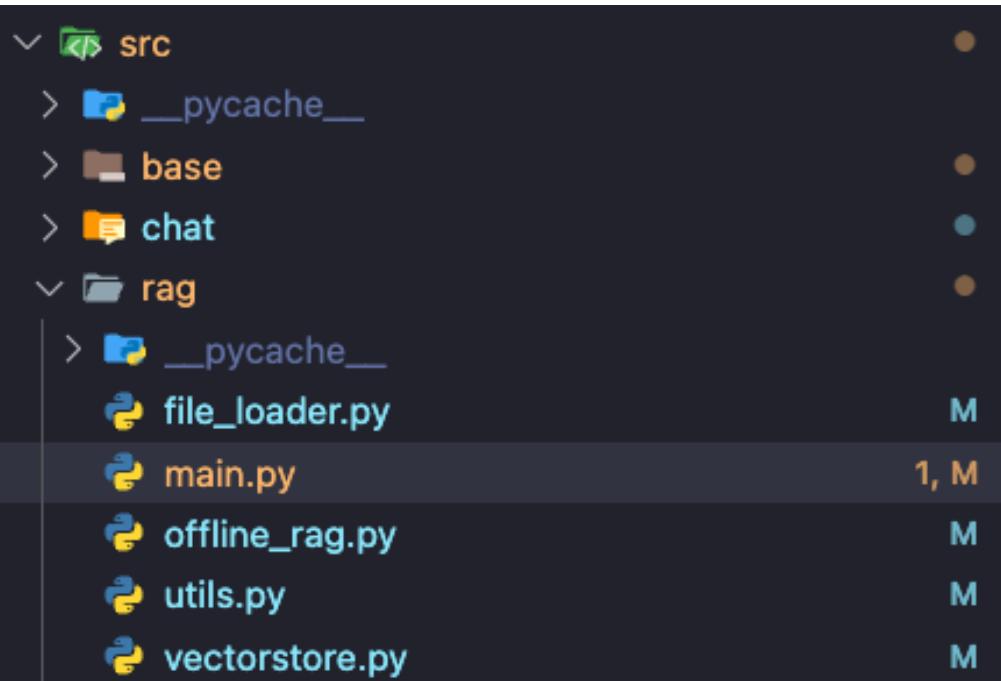


```
1 def get_hf_llm(model_name: str = "mistralai/Mistral-7B-Instruct-v0.2",
2 max_new_token = 1024,
3 **kwargs):
4
5 model = AutoModelForCausalLM.from_pretrained(
6 model_name,
7 quantization_config=nf4_config,
8 low_cpu_mem_usage=True
9)
10 tokenizer = AutoTokenizer.from_pretrained(model_name)
11
12 model_pipeline = pipeline(
13 "text-generation",
14 model=model,
15 tokenizer=tokenizer,
16 max_new_tokens=max_new_token,
17 pad_token_id=tokenizer.eos_token_id,
18 device_map="auto"
19)
20
21 llm = HuggingFacePipeline(
22 pipeline=model_pipeline,
23 model_kwargs=kwargs
24)
25
26 return llm
```

# RAG with LangChain

## ❖ Step 3: Build RAG source code

Update src/rag/main.py:

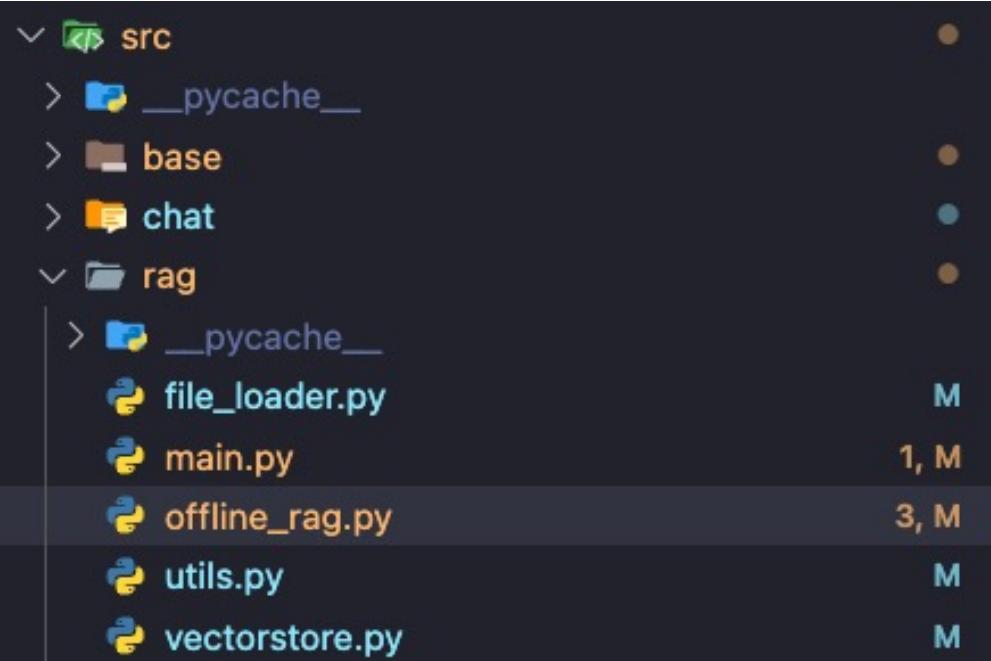


```
1 from pydantic import BaseModel, Field
2
3 from src.rag.file_loader import Loader
4 from src.rag.vectorstore import VectorDB
5 from src.rag.offline_rag import Offline_RAG
6
7 class InputQA(BaseModel):
8 question: str = Field(..., title="Question to ask the model")
9
10 class OutputQA(BaseModel):
11 answer: str = Field(..., title="Answer from the model")
12
13 def build_rag_chain(llm, data_dir, data_type):
14 doc_loaded = Loader(file_type=data_type).load_dir(data_dir, workers=2)
15 retriever = VectorDB(documents = doc_loaded).get_retriever()
16 rag_chain = Offline_RAG(llm).get_chain(retriever)
17
18 return rag_chain
```

# RAG with LangChain

## ❖ Step 3: Build RAG source code

Update src/rag/offline\_rag.py:

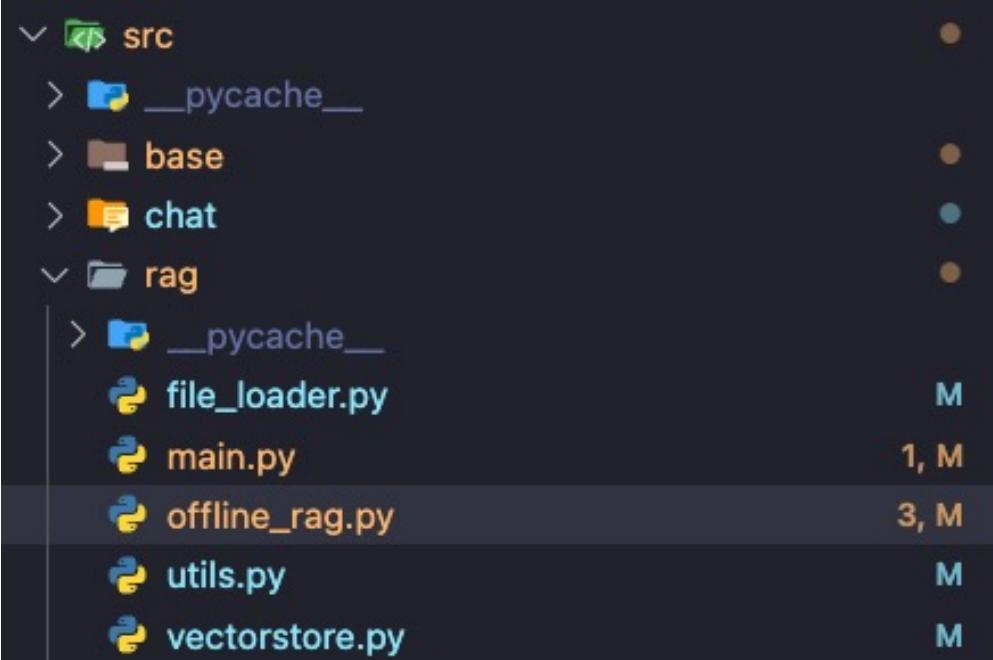


```
 1 import re
 2 from langchain import hub
 3 from langchain_core.runnables import RunnablePassthrough
 4 from langchain_core.output_parsers import StrOutputParser
 5
 6 class Str_OutputParser(StrOutputParser):
 7 def __init__(self) -> None:
 8 super().__init__()
 9
10 def parse(self, text: str) -> str:
11 return self.extract_answer(text)
12
13
14 def extract_answer(self,
15 text_response: str,
16 pattern: str = r"Answer:\s*(.*)"
17) -> str:
18
19 match = re.search(pattern, text_response, re.DOTALL)
20 if match:
21 answer_text = match.group(1).strip()
22 return answer_text
23 else:
24 return text_response
```

# RAG with LangChain

## ❖ Step 3: Build RAG source code

Update src/rag/offline\_rag.py:

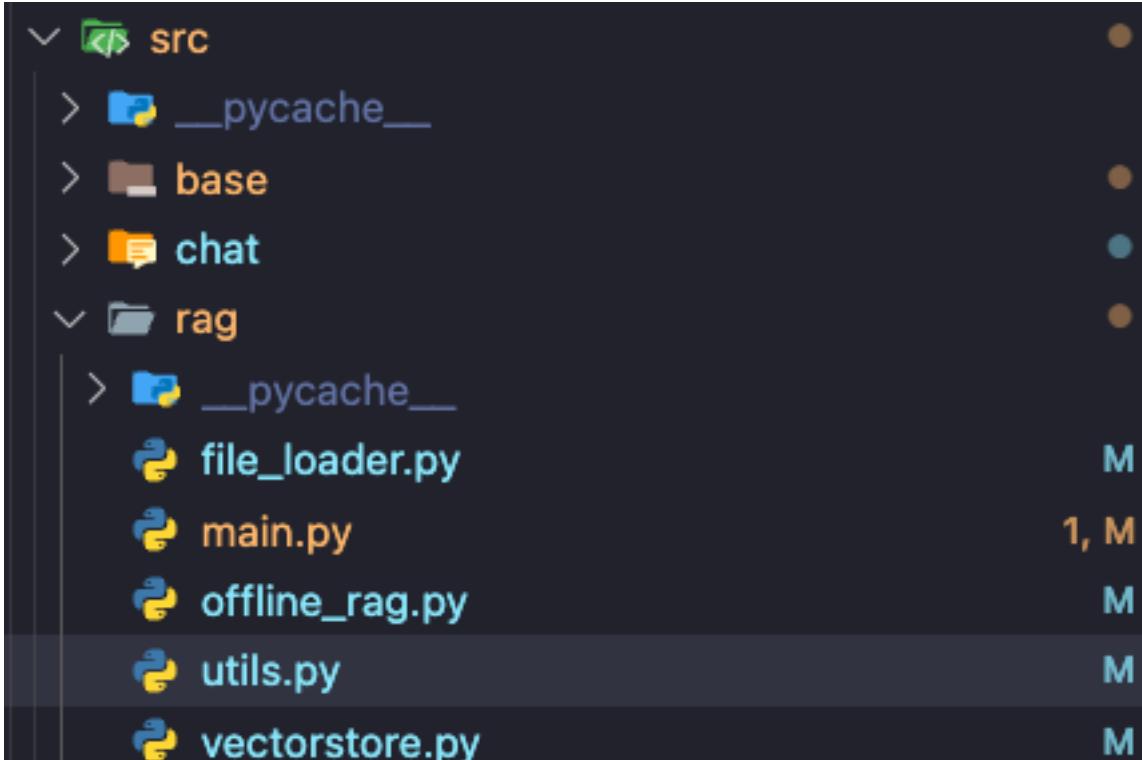


```
●●●
1
2 class Offline_RAG:
3 def __init__(self, llm) -> None:
4 self.llm = llm
5 self.prompt = hub.pull("rlm/rag-prompt")
6 self.str_parser = Str_OutputParser()
7
8 def get_chain(self, retriever):
9 input_data = {
10 "context": retriever | self.format_docs,
11 "question": RunnablePassthrough()
12 }
13 rag_chain = (
14 input_data
15 | self.prompt
16 | self.llm
17 | self.str_parser
18)
19 return rag_chain
20
21 def format_docs(self, docs):
22 return "\n\n".join(doc.page_content for doc in docs)
```

# RAG with LangChain

## ❖ Step 3: Build RAG source code

Update src/rag/utils.py:

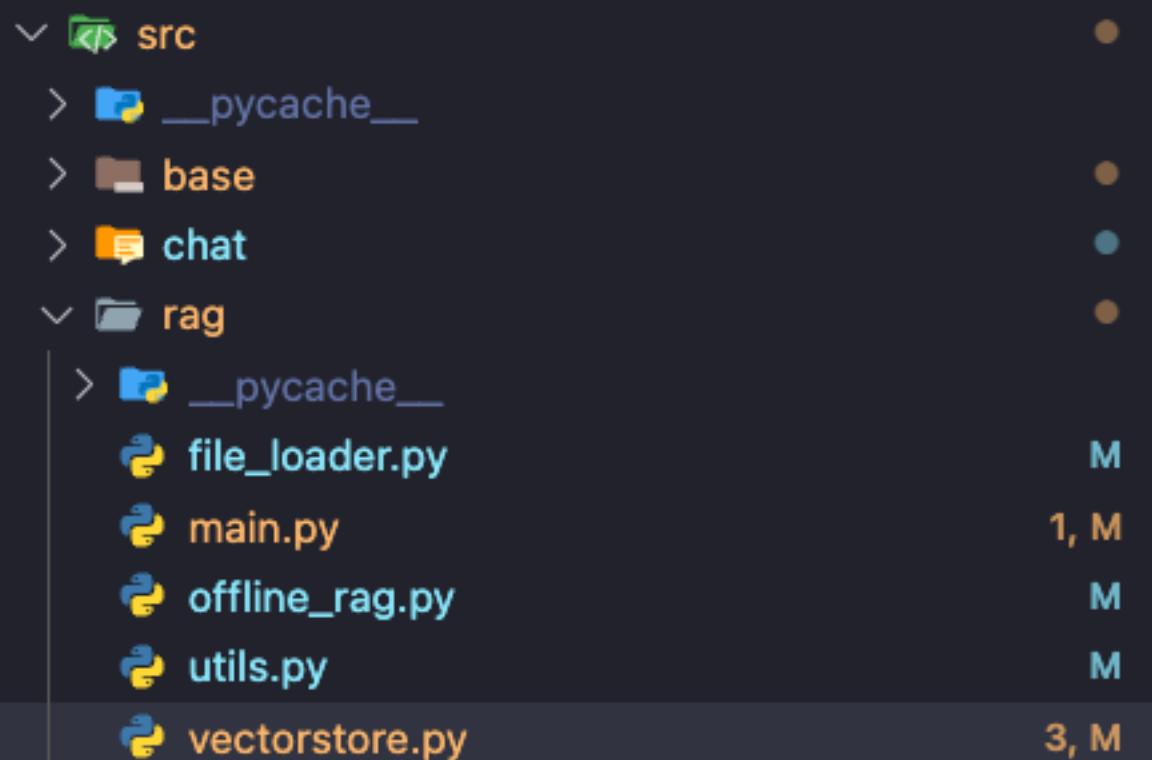


```
● ● ●
1 import re
2
3 def extract_answer(text_response: str,
4 pattern: str = r"Answer:\s*(.*)"
5) -> str:
6
7 match = re.search(pattern, text_response)
8 if match:
9 answer_text = match.group(1).strip()
10 return answer_text
11 else:
12 return "Answer not found."
```

# RAG with LangChain

## ❖ Step 3: Build RAG source code

Update src/rag/vectorstore.py:

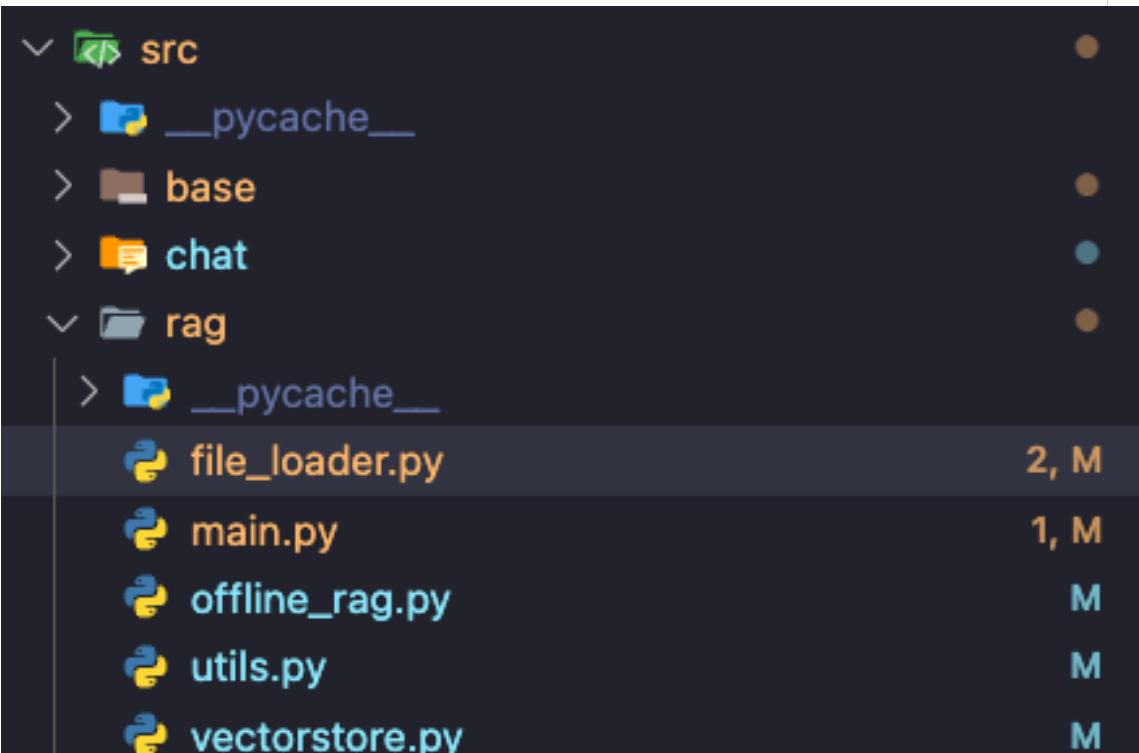


```
 1 from typing import Union
 2 from langchain_chroma import Chroma
 3 from langchain_community.vectorstores import FAISS
 4 from langchain_community.embeddings import HuggingFaceEmbeddings
 5
 6 class VectorDB:
 7 def __init__(self,
 8 documents = None,
 9 vector_db: Union[Chroma, FAISS] = Chroma,
10 embedding = HuggingFaceEmbeddings(),
11) -> None:
12
13 self.vector_db = vector_db
14 self.embedding = embedding
15 self.db = self._build_db(documents)
16
17 def _build_db(self, documents):
18 db = self.vector_db.from_documents(documents=documents,
19 embedding=self.embedding)
20 return db
21
22 def get_retriever(self,
23 search_type: str = "similarity",
24 search_kwargs: dict = {"k": 10}
25):
26 retriever = self.db.as_retriever(search_type=search_type,
27 search_kwargs=search_kwargs)
27
28 return retriever
```

# RAG with LangChain

## ❖ Step 3: Build RAG source code

Update src/rag/file\_loader.py:

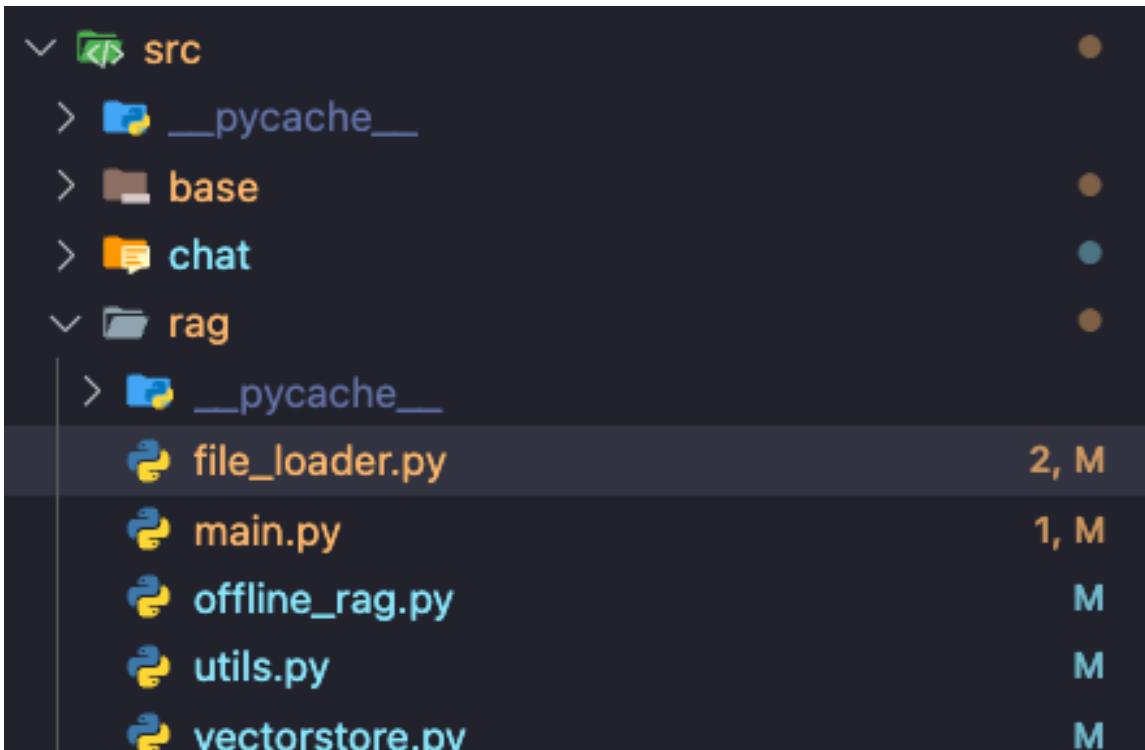


```
1 from typing import Union, List, Literal
2 import glob
3 from tqdm import tqdm
4 import multiprocessing
5 from langchain_community.document_loaders import PyPDFLoader
6 from langchain_text_splitters import RecursiveCharacterTextSplitter
7
8 def remove_non_utf8_characters(text):
9 return ''.join(char for char in text if ord(char) < 128)
10
11 def load_pdf(pdf_file):
12 docs = PyPDFLoader(pdf_file, extract_images=True).load()
13 for doc in docs:
14 doc.page_content = remove_non_utf8_characters(doc.page_content)
15 return docs
16
17 def get_num_cpu():
18 return multiprocessing.cpu_count()
```

# RAG with LangChain

## ❖ Step 3: Build RAG source code

Update src/rag/file\_loader.py:

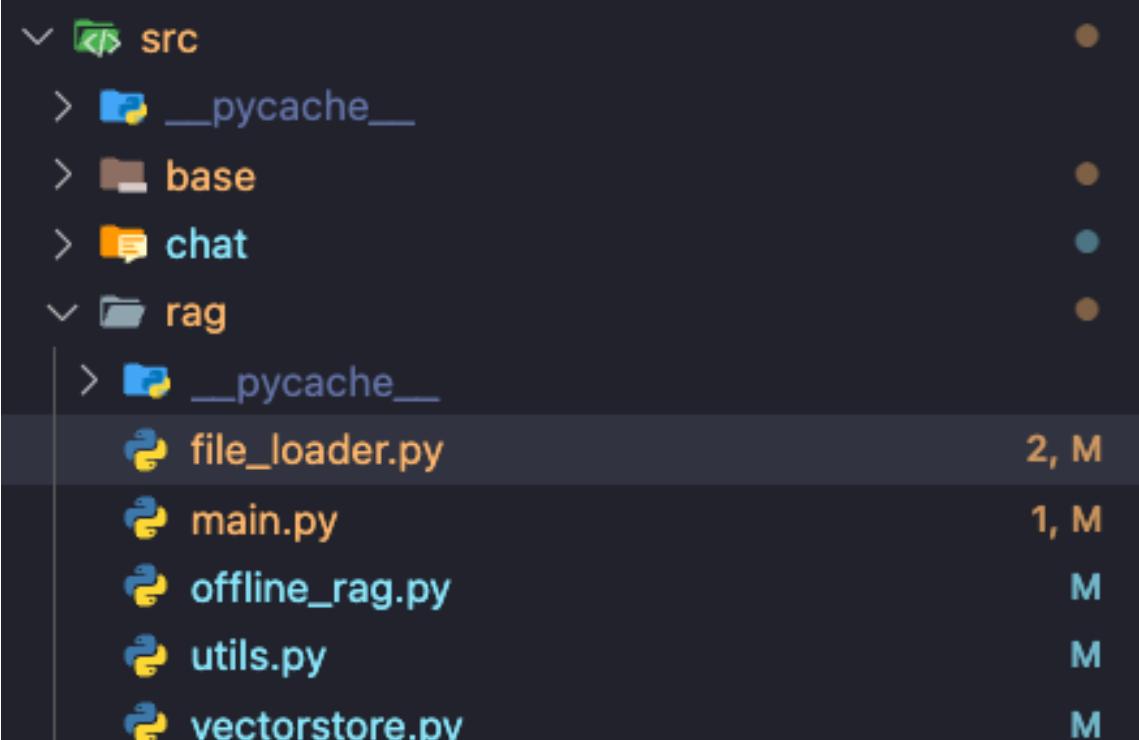


```
1 class BaseLoader:
2 def __init__(self) -> None:
3 self.num_processes = get_num_cpu()
4
5 def __call__(self, files: List[str], **kwargs):
6 pass
7
8 class PDFLoader(BaseLoader):
9 def __init__(self) -> None:
10 super().__init__()
11
12 def __call__(self, pdf_files: List[str], **kwargs):
13 num_processes = min(self.num_processes, kwargs["workers"])
14 with multiprocessing.Pool(processes=num_processes) as pool:
15 doc_loaded = []
16 total_files = len(pdf_files)
17 with tqdm(total=total_files, desc="Loading PDFs", unit="file") as pbar:
18 for result in pool.imap_unordered(load_pdf, pdf_files):
19 doc_loaded.extend(result)
20 pbar.update(1)
21
22 return doc_loaded
```

# RAG with LangChain

## ❖ Step 3: Build RAG source code

Update src/rag/file\_loader.py:

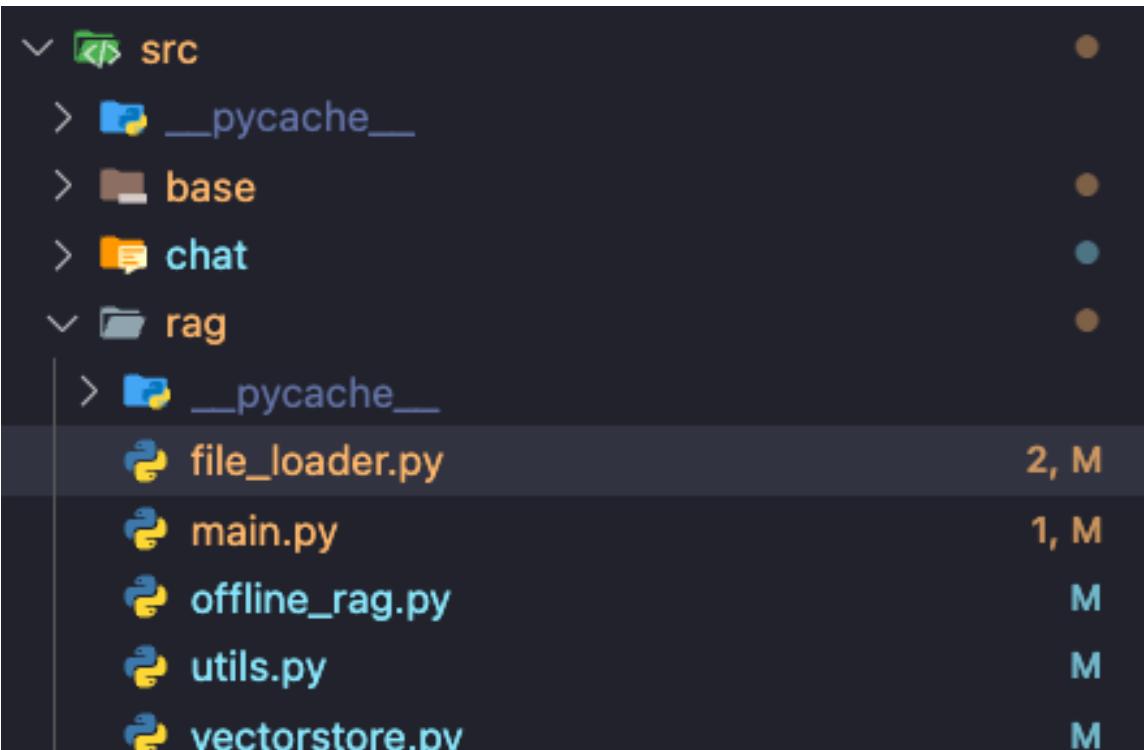


```
● ● ●
1 class TextSplitter:
2 def __init__(self,
3 separators: List[str] = ['\n\n', '\n', ' ', ''],
4 chunk_size: int = 300,
5 chunk_overlap: int = 0
6) -> None:
7
8 self.splitter = RecursiveCharacterTextSplitter(
9 separators=separators,
10 chunk_size=chunk_size,
11 chunk_overlap=chunk_overlap,
12)
13 def __call__(self, documents):
14 return self.splitter.split_documents(documents)
```

# RAG with LangChain

## ❖ Step 3: Build RAG source code

Update src/rag/file\_loader.py:

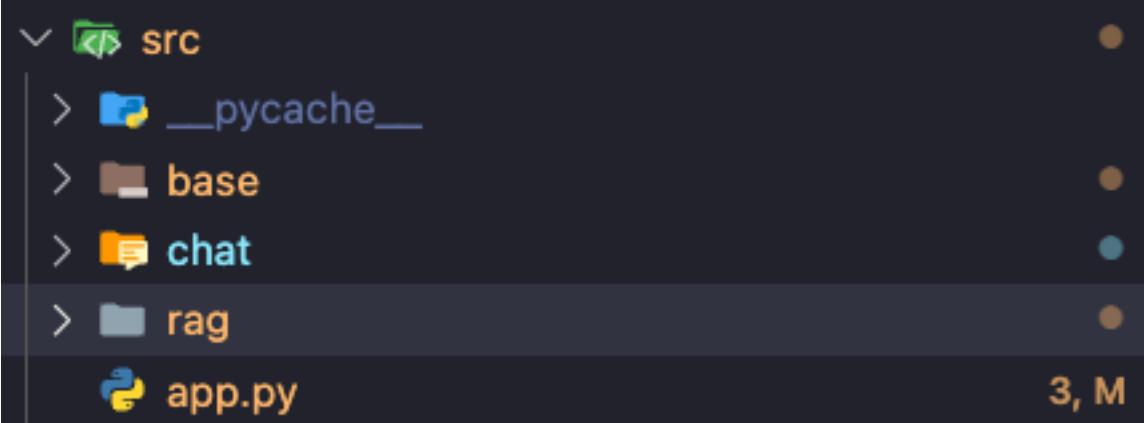


```
1 class Loader:
2 def __init__(self,
3 file_type: str = Literal["pdf"],
4 split_kwargs: dict = {
5 "chunk_size": 300,
6 "chunk_overlap": 0}
7) -> None:
8 assert file_type in ["pdf"], "file_type must be pdf"
9 self.file_type = file_type
10 if file_type == "pdf":
11 self.doc_loader = PDFLoader()
12 else:
13 raise ValueError("file_type must be pdf")
14
15 self.doc_splitter = TextSplitter(**split_kwargs)
16
17 def load(self, pdf_files: Union[str, List[str]], workers: int = 1):
18 if isinstance(pdf_files, str):
19 pdf_files = [pdf_files]
20 doc_loaded = self.doc_loader(pdf_files, workers=workers)
21 doc_split = self.doc_splitter(doc_loaded)
22 return doc_split
23
24 def load_dir(self, dir_path: str, workers: int = 1):
25 if self.file_type == "pdf":
26 files = glob.glob(f"{dir_path}/*.pdf")
27 assert len(files) > 0, f"No {self.file_type} files found in {dir_path}"
28 else:
29 raise ValueError("file_type must be pdf")
30
31 return self.load(files, workers)
```

# RAG with LangChain

## ❖ Step 4: Build API

Update src/app.py:

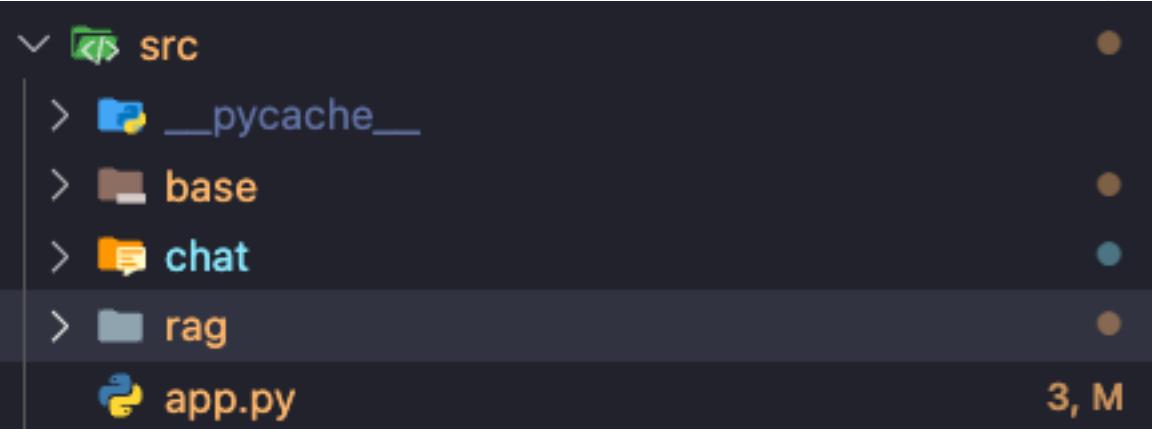


```
● ● ●
1 import os
2 os.environ["TOKENIZERS_PARALLELISM"] = "false"
3
4 from fastapi import FastAPI
5 from fastapi.middleware.cors import CORSMiddleware
6
7 from langserve import add_routes
8
9 from src.base.llm_model import get_hf_llm
10 from src.rag.main import build_rag_chain, InputQA, OutputQA
11
12 llm = get_hf_llm(temperature=0.9)
13 genai_docs = "./data_source/generative_ai"
```

# RAG with LangChain

## ❖ Step 4: Build API

Update src/app.py:

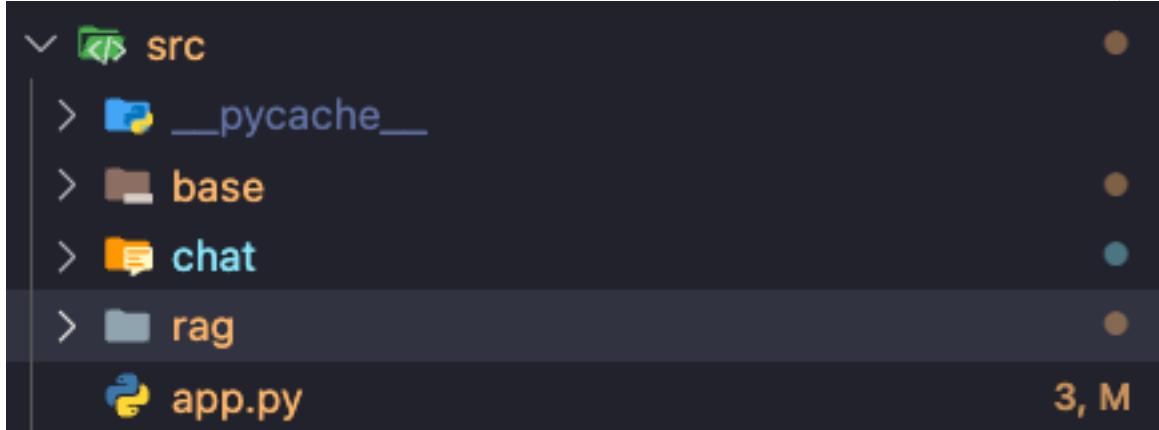


```
● ● ●
1 import os
2 os.environ["TOKENIZERS_PARALLELISM"] = "false"
3
4 from fastapi import FastAPI
5 from fastapi.middleware.cors import CORSMiddleware
6
7 from langserve import add_routes
8
9 from src.base.llm_model import get_hf_llm
10 from src.rag.main import build_rag_chain, InputQA, OutputQA
11
12 llm = get_hf_llm(temperature=0.9)
13 genai_docs = "./data_source/generative_ai"
14
15 # ----- Chains-----
16
17 genai_chain = build_rag_chain(llm, data_dir=genai_docs,
 data_type="pdf")
```

# RAG with LangChain

## ❖ Step 4: Build API

Update src/app.py:

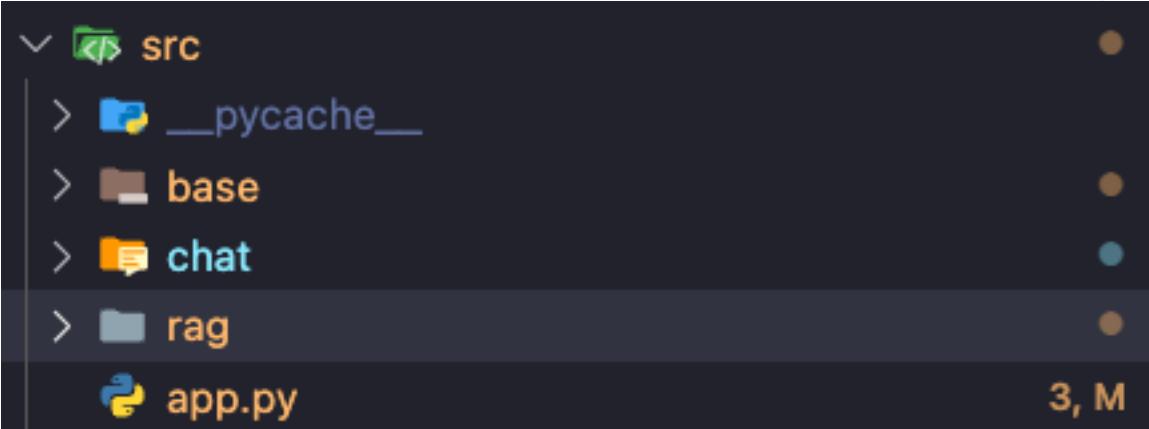


```
1 # ----- App - FastAPI -----
2
3 app = FastAPI(
4 title="LangChain Server",
5 version="1.0",
6 description="A simple api server using Langchain's Runnable
7 interfaces",
8)
9 app.add_middleware(
10 CORSMiddleware,
11 allow_origins=["*"],
12 allow_credentials=True,
13 allow_methods=["*"],
14 allow_headers=["*"],
15 expose_headers=["*"],
16)
```

# RAG with LangChain

## ❖ Step 4: Build API

Update src/app.py:



```
1 # ----- Routes - FastAPI -----
2
3 @app.get("/check")
4 async def check():
5 return { "status": "ok" }
6
7
8 @app.post("/generative_ai", response_model=OutputQA)
9 async def generative_ai(inputs: InputQA):
10 answer = genai_chain.invoke(inputs.question)
11 return { "answer": answer }
12
13 # ----- Langserve Routes - Playground -----
14 add_routes(app,
15 genai_chain,
16 playground_type="default",
17 path="/generative_ai")
```

# RAG with LangChain

## ❖ Step 5: Deploy and test API

```
INFO: Will watch for changes in these directories: ['/home/server-ailab/ThangDuongTeam/thangdd/TA/240504_Extra_LangChain/Langchain_Services']
INFO: Uvicorn running on http://0.0.0.0:5000 (Press CTRL+C to quit)
INFO: Started reloader process [6115] using WatchFiles
Loading checkpoint shards: 100% | 3/3 [00:08<00:00, 2.83s/it]
Loading PDFs: 100% | 6/6 [00:41<00:00, 6.96s/file]
INFO: Started server process [6117]
INFO: Waiting for application startup.

LANGSERVE
LANGSERVE: Playground for chain "/generative_ai/" is live at:
LANGSERVE: ↳ /generative_ai/playground/
LANGSERVE:
LANGSERVE: See all available routes at /docs/

LANGSERVE: ▲ Using pydantic 2.7.0. OpenAPI docs for invoke, batch, stream, stream_log endpoints will not be generated. API endpoints and playground should work as expected. If you need to see the docs, you can downgrade to pydantic 1. For example, `pip install pydantic==1.10.13`. See https://github.com/tiangolo/fastapi/issues/10360 for details.

INFO: Application startup complete.
```

```
uvicorn src.app:app --host "0.0.0.0" --port 5000 --reload
```

# RAG with LangChain

## ❖ Step 5: Deploy and test API

**LangChain Server** 1.0 OAS 3.1  
[/openapi.json](#)

A simple api server using Langchain's Runnable interfaces

**generative\_ai** ⚠️ Using pydantic 2.7.0. OpenAPI docs for `invoke`, `batch`, `stream`, `stream_log` endpoints will not be generated. API endpoints and playground should work as expected. If you need to see the docs, you can downgrade to pydantic 1. For example, `pip install pydantic==1.10.13` See <https://github.com/tiangolo/fastapi/issues/10360> for details.

**GET** [/generative\\_ai/input\\_schema](#) Generative Ai Input Schema

**GET** [/generative\\_ai/output\\_schema](#) Generative Ai Output Schema

**GET** [/generative\\_ai/config\\_schema](#) Generative Ai Config Schema

**generative\_ai/config** Endpoints with a default configuration set by `config_hash` path parameter. Used in conjunction with share links generated using the LangServe UI playground. The hash is an LZString compressed JSON string.

**GET** [/generative\\_ai/c/{config\\_hash}/input\\_schema](#) Generative Ai Input Schema With Config

**GET** [/generative\\_ai/c/{config\\_hash}/output\\_schema](#) Generative Ai Output Schema With Config

**GET** [/generative\\_ai/c/{config\\_hash}/config\\_schema](#) Generative Ai Config Schema With Config

**default**

**GET** [/check](#) Check

**POST** [/generative\\_ai](#) Generative Ai

**POST** [/generative\\_ai/token\\_feedback](#) Create Feedback From Token

# RAG with LangChain

## ❖ Step 5: Deploy and test API

Curl

```
curl -X 'POST' \
'http://0.0.0.0:5050/generative_ai' \
-H 'accept: application/json' \
-H 'Content-Type: application/json' \
-d '{
 "question": "What is instruction tuning in LLMs?"
}'
```

Request URL

```
http://0.0.0.0:5050/generative_ai
```

Server response

| Code | Details                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                             |
|------|-----------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------|
| 200  | Response body <pre>{ "answer": "Instruction tuning is a technique used to enhance the capabilities and controllability of large language models (LLMs) by fine-tuning them on a dataset consisting of instruction-output pairs in a supervised manner. This process bridges the gap between the next-word prediction objective of LLMs and the user's objective of instruction following. It allows for more controllable and predictable model behavior and provides a channel for humans to intervene with the model's behaviors. Instruction tuning is also computationally efficient and can help LLMs rapidly adapt to a specific domain without extensive retraining or architectural changes. However, crafting high-quality instructions that properly cover the desired target behaviors is a challenge." }</pre> <a href="#">Download</a> |

# Quiz

# Summary

**In this lecture, we have discussed:**

**1. What is LLMs in production?**

1. How is it differ from LLMs in research?
2. What are some challenges when deploying LLMs in production?

**2. Basics of LangChain**

1. The key concept of LangChain.
2. Basic components of LangChain: Prompt Template, LLM Chain, Chat History, Document Loader...

**3. How to build an API using LangChain.**

**4. How to build a RAG application using LangChain.**

1. Retrieve and answer questions related to Academic Paper.

# Question

