# LangChain in your Pocket

Beginner's Guide to Building Generative AI Applications using LLMs

Mehul Gupta

# LangChain in your Pocket

**Beginner's Guide to Building Generative AI Applications using LLMs**

**Mehul Gupta**

*To mummy and papa,*
*thanks for being more interested in my meals*
*than my career.*

# Copyright

LangChain in your Pocket

Copyright © 2024 by Mehul Gupta

For permission requests, connect using the below mail id:
mehulgupta2016154@gmail.com

# Table of Contents

# Preface

Dear reader,

Hope you're doing great. Welcome to my inaugural book, *LangChain in your Pocket: Beginner's Guide to Building Generative AI Applications using LLMs*. I assume you chose to read this book because you are, alongside everyone, just blown away by the recent advancements in Generative AI, especially the incoming LLMs. And as everyone around wishes to upgrade themselves to understand and use Generative AI at scale and move beyond ChatGPT, you also wish to try your hands with LangChain, one of the best frameworks out there enabling app building using LLMs.

**So what does this book have to offer?**
As the name suggests, this book is a beginner-friendly introduction to LangChain, one of the most important frameworks that helps you build AI apps. This book covers most of the major use cases that can be catered using LangChain. The best part is that the focus is not just on the theoretical concepts but also includes short code snippets with explanations making things comparatively easier to understand and striking the right balance between theoretical explanations and real-world implementation. Whether you're a seasoned data scientist seeking to deepen your knowledge or a newcomer eager to explore the intricacies of Generative AI, this book aims to provide a robust foundation and practical insights. This book is not just about LangChain, it's a journey through the fascinating world where Generative AI meets real-world problems. I encourage you to experiment, explore and apply the given examples in diverse contexts to unleash LangChain's full potential.

**Why did I write this book?**
I have been a seasoned Medium blogger and YouTuber for nearly 4-5 years now where I share my findings in Data Science & AI constantly. So writing is not new to me but writing in a longer format is. Writing a book was always on my bucketlist and this year I grabbed my opportunity. I believe in making AI knowledge available to everyone and this is just another step in that direction. I hope you support this book similar to how you have supported me in my previous ventures. I have tried my level best to keep it interesting for folks with different levels of understanding of Generative AI and glad to help anyone who wants further assistance.

**Why should you read this book?**
I have many reasons to convince you. I've been working in the space of AI for nearly 5 years now and have written many technical pieces. So, I assume, I've cracked the code of how to write beginner-friendly technical tutorials. Also, in the past year, I have explored LangChain in a lot of depth knowing all its secrets that aren't easily available. Further, the documentation provided by the LangChain team is commendable (and I mean it) but scattered all over the place and is also very lengthy. A few utilities looked a little unnecessary to me. So, this book can be considered a collection of important LangChain use cases ignoring irrelevant stuff and is not at all scattered. Don't take it as a technical book, but a journal on my last year's journey with Generative AI using LangChain.

Finally, I wish to extend my gratitude to all those whose work laid the groundwork for the content of this book most importantly my family, who, not just while writing this book, have given undeterred support through thick and thin and let me go wild with whatever I wished to do. I also want to thank my readers for their interest and dedication to mastering the intricacies of data science algorithms. I hope this book answers all your questions.

Happy reading!
Mehul

# Chapter 1: Introduction

2023 has been the year of GenAI. Since the inception of ChatGPT in late 2022, there has been no turning back. One after the other new things are coming up: different LLMs like GPT, LLaMA, PaLM, etc; frameworks like LangChain and Llama-Index; AI agents like AutoGPT, BabyAGI, etc. This book is a deep dive into the most famous and useful framework around LLMs, i.e., LangChain, an open-source project that debuted in October 2022, pioneered by Harrison Chase and initially developed while Chase was at Robust Intelligence, an MLOps company.

But before we jump onto LangChain, we must understand LLMs.

## 1.1 What are LLMs?

Large Language Models (LLMs) are nothing but Machine Learning models that have a human-like understanding of language. Because of this ability, any random task, be it Summary generation, Language Translation, Classification, etc. can be done using a single model. These advanced algorithms utilise deep learning techniques and vast datasets to comprehend and generate human-like language. The term *"large"* emphasizes the extensive training these models undergo on massive datasets, enabling them to capture intricate patterns and entity relationships within the text and their usual humongous size. Let's very quickly understand the salient features of a general LLM:

1. An LLM, based on Neural Networks, predicts subsequent tokens from a given sequence of text and already predicted tokens. There is no black box!
2. To understand the context, LLMs utilise ***multi-head attention*** (discussed in the paper 'Attention is all you need') or its variants, focussing on vital words within the input sequence. It would require another book to discuss the concept of Attention so skipping it for now.
3. As an LLM predicts tokens sequentially, it never knows the final output at any given moment unless the final output is completely predicted.
4. Zero-shot abilities allow LLMs to tackle diverse problems, not confined to specific tasks like classification or regression.

***What is Zero-shot learning?***
*Zero-shot learning refers to a capability within ML models like LLMs to perform tasks or generate outputs for which they haven't been explicitly trained. In essence, these models can generalize their knowledge to tackle unseen or new tasks without specific examples or data during their training phase. This ability comes from the model's extensive training on huge corpora where these models can learn varied linguistic patterns and contexts, enabling it to infer and generalize information to new scenarios.*

5. Typically massive in size and trained on vast datasets, LLMs adhere to the ***Transformers architecture*** or a subset of it like GPT (decoder part of a Transformer).

6. Pre-trained LLMs, akin to '*Generalists*,' are well-suited for various tasks but not exceptional in any specific area. So if you wish to have a 99% accuracy on a classification task using a pre-trained LLM, this is a little difficult.
7. For specialized tasks, '*Fine-Tuning*' involves training a pre-trained LLM with specific data.
8. LLMs vary in size denoted by terms like small, large, or XXL, influenced by parameters and architectural modifications. The size usually denotes the number of weight parameters being used.
9. The numeric identifier in some LLM names, such as LLaMA-70B, signifies the number of parameters. Such LLM families don't use the small, large, XL nomenclature.
10. Details regarding training data or architecture for LLMs are often not fully disclosed by their respective owner companies, creating gaps in understanding certain LLM families.
11. Due to their massive size*, **loading LLMs into memory** can be a challenging task**.** This becomes a bottleneck when you don't have huge resources. This is where APIs come into the picture where you can access models as big as hundreds of GBs with a url without loading it in your memory. We will be using OpenAI API for most of the tutorials in this book.
12. LLMs, being generative models, introduce an element of randomness, resulting in **varying outputs for the same input**. This is something you all must have observed while using ChatGPT.
13. To perform any task, we need to provide the LLM a prompt i.e. a natural language sentence specifying the task's details.
14. Every LLM has a **token restriction** (word limit) for input and output.

As of now (2023 end), many organizations have developed a variety of LLMs, each one of them bringing something new to the table. Let's have a brief about different types of LLM families.

# 1.2 Different LLM families

### The GPTs
This is where the GenAI era kickstarted. GPT stands for **Generative Pre-trained Transformers**. ChatGPT stems from this family of models and one of its versions, GPT3.5, serves as its foundation.

1. GPTs, including GPT3.5, are built upon the **decoder** component of a Transformer architecture. This design gives them the ability to generate content on their own.
2. When comparing models, their parameter counts play a significant role. For instance, GPT3.5 boasts 175 billion parameters, while its successor, GPT4, amasses a whopping 1,760 billion parameters.
3. GPT3.5 Turbo represents an upgraded and more finely-tuned iteration of GPT3.5. With around 20 billion parameters tailored specifically for natural language processing tasks, it enhances performance and optimization.
4. Among these models, only GPT4 possesses the **multi-modal feature**, allowing it to process and understand various types of inputs, such as images, audio and more.

### LLaMA by Meta
LLaMA, standing for **Large Language Model Meta AI**, is a family of LLMs developed and released by Meta AI, the parent company of Facebook following the Transformer architecture. Introduced in

February 2023, LLaMA is designed to be an open-source and efficient foundation for language models, ranging from 7 billion to 65 billion parameters. The first version of LLaMA included multiple model sizes and Meta AI has continued to evolve this family of models. A special effort has been put into making it more safe and secure to use.

## FLAN models

FLAN-T5 and FLAN-Alpaca are LLMs introduced as instruction-tuned language models. These models leverage synthetic instruction tuning, a method that fine-tunes a base language model with instructions generated by humans and machines.

### *What is Instruction-based Fine-Tuning?*

*During instruction fine-tuning, you're required to include specific guidelines, termed 'instructions,' along with the input and output data while training the model. Here's an example illustrating the process of instruction fine-tuning for an LLM:*

### *General fine-tuning sample:*

```
Input: 'Compose a narrative involving the theme of nature'
Output: '........'


Input: 'Identify elements in this passage: the house is blue'
Output: '......'
```

### *Instruction fine-tuning adaptation:*

```
Instruct: 'Generate a story based on the provided subject'
Input: 'Nature and wildlife'
Output: '........'


Instruct: 'Recognize details within the given text'
Input: 'The house is painted blue'
Output: '.......'
```

Below is a quick brief on two major FLAN models:

1. **FLAN-T5:** FLAN-T5 is an enhanced version of T5 (Text-to-Text Transfer Transformer) that has been fine-tuned using a mixture of tasks, as detailed in the paper "Scaling Instruction-Finetuned Language Models".
2. **FLAN-Alpaca**: FLAN-Alpaca extends the concept of instruction tuning and is a fine-tuned version of LLaMA.

## Mistral

Mistral AI emerged as a noteworthy startup, gaining attention with their inaugural model, Mistral-7B, which outperformed established competitors despite having just 7 billion parameters.

1. Architecture-wise, Mistral models mirror GPT, relying on the Transformer's Decoder component.

2. Parameters vary across Mistral and Mixtral, with Mistral having 7B parameters and Mixtral featuring 84B. Mixtral is a Mixture of Expert (MoE) model, which involves a combination of multiple specialized networks for doing specific sub-tasks in a given problem.

**Phi by Microsoft**

Phi, developed by Microsoft, prioritizes delivering strong performance while maintaining a notably smaller model size, which makes it more suitable for practical use in production environments.

1. In terms of structure, Phi models also adopt the Transformer model as their foundation.
2. When considering parameters, Phi models vary in size: Phi-1 and Phi-1.5 hover around 1.3 billion parameters, while Phi-2 boasts 2.7 billion.
3. Interestingly, despite its smaller size, Phi-2 has outperformed models nearly ten times larger, such as LLaMA-13B and has even surpassed ChatGPT in certain tasks.

This list of language models extends far and wide. While I couldn't cover all, notable models besides those mentioned include Claude (with approximately 130 billion parameters, yielding similar results to GPT3.5), Cohere, PaLM, T5 (developed by Google), Falcon and more. Exploring these models can offer further insights into their capabilities.

Now that we have known enough about LLMs, time to go back to LangChain.

# 1.3 What is LangChain used for?

A simple answer: ***Simplifying any task you wish to accomplish using LLMs.***

LangChain is designed to simplify the creation of applications using LLMs. Be it a grammar correction app or a complex NER extraction, LangChain has got you covered.

1. **Advanced Chat Applications**: LangChain enables the development of chat applications capable of handling complex questions and transactions from users.
2. **Composable Tools and Integrations**: LangChain offers composable 3rd party tools and integrations like Google Search, ArXiv, Bing Search, Wikipedia, etc. making it easy for developers to work with language models.
3. **Short-Term and Long-Term Memory Support:** LangChain supports short-term and long-term memory, allowing LLMs to retain information across interactions, which is particularly useful in chatbot scenarios.
4. **Support to most LLMs:** LangChain provides support to most of the LLMs using API or even local LLMs that can be loaded using Hugging Face or fine-tuned by you.
5. **Evaluation:** Not just building apps, it does provide functionality to evaluate the results of your LLMs by using either pre-defined metrics or custom metrics. This is crucial because evaluating a generative model, be it GANs or LLMs, is always a challenge.
6. **Complex applications:** LangChain can be used for building solutions for real complex problems like recommendation systems, AI agents like AutoGPT, etc.

# 1.4 Why LangChain?

To answer this, we will take a step back and understand what contributes to the success of a task when dealing with GenAI and LLMs:

1. *The LLM you're using.*
2. *The prompt you feed.*
3. *The integration of external tools/wrappers you wish to use (say Google search).*

LangChain, in the hood, takes care of all three aspects together, especially the prompt and coding part. So assume you wish to create an Automatic SQL generator using LLMs with the following features:

1. The app should connect and traverse through tables and DBs.
2. Generate queries given a prompt.

Now, to do this manually, you need a lot of effort from writing codes to connect with databases, fetching data, getting an apt prompt and whatnot. Now using LangChain, this is a couple of  lines of code as everything gets handled in the hood by LangChain.

***It is possible to create any app without LangChain also, but LangChain simplifies it to a great extent and  hence way better than manual prompting.***

I have even used other frameworks available open-source like Llama-Index or Promptify but the range that LangChain has got is just amazing. I can't think of an app I wish to create which I can't do using LangChain. Others do have great utilities but are very specific like Llama-Index is a specialist with the RAG framework (discussed in upcoming chapters) but doesn't have a range as big as LangChain. Hence, LangChain should be your go-to framework any day.

# 1.5 Book Overview

In this book, we will be going through a lot of concepts starting with:

***Hello World***
This chapter deals with setting up LangChain in your system and building a series of beginner-friendly LangChain apps using the OpenAI API key or local LLMs.

***Different LangChain Modules***
Once done with the basic setup in the above chapter, this chapter will introduce different components of the LangChain Framework like Models, PromptTemplate, Chains, Agents, Callbacks, Retrieval and  Memory and what role these modules play in creating an application.

***Model and Prompts***
This chapter explains the Model and Prompts module from LangChain that enables different models and provides numerous templates for building any LangChain app.

### Chains
The heart of LangChain, chains can be considered as short programs that are written around LLMs to provide a one-line solution for common problems like Summarization, NER, etc.

### Agents
More flexible and versatile compared to chains, agents help you build complex solutions and enable access to 3rd party tools like Google search, file writing, etc. with ease.

### OutputParsers and Memory
The OutputParsers enable formatting the output from LLMs making them more useful for production-grade applications and avoiding ambiguous formats. The memory module enables memory for LangChain apps, hence remembering past conversations in the current session similar to ChatGPT.

### Callbacks
Similar to Tensorflow callbacks, it enables a deep dive into the internal workings of LangChain apps making logging and debugging easy.

### RAG Framework and Vector Databases
One of the key features of LangChain is its support for RAG i.e. providing external context to LLMs as text files, PDFs, etc. Check out how to implement RAG using Vector Databases and LangChain in this chapter.

### LangChain for NLP
Most important of all, this chapter deals with how LLMs and LangChain can be used for solving various NLP problems, be it NER, Text tagging, Sentiment analysis, Few-shot classification, etc.

### Handling LLM Hallucinations
Since the inception of LLMs, hallucinations i.e. LLMs giving factually wrong answers have been a problem. LangChain provides a solution to that too which you can check out in this chapter.

### Evaluating LLMs
Unless you can measure metrics around your ML problem, you can't evaluate how good is the solution. Measuring metrics for LLMs is a tough task but LangChain provides an answer to this as well. Check out this section on how to calculate different metrics for LLMs.

### Advanced Prompt Engineering
LangChain supports Prompt Engineering Frameworks like ReAct, Trees of Thoughts, etc. that can be explored in this chapter alongside examples.

### Autonomous AI agents
After LLMs, the next big thing coming up is AI agents like AutoGPT, BabyAGI, etc. Check out how LangChain enables these AI agents in this book.

***LangSmith and LangServe***

Before ending, this chapter introduces two major extensions to LangChain that are LangSmith which enables detailed debugging, logging and visualization while LangServe supports deploying a LangChain app.

***Additional Features***

This chapter covers some other functionalities that LangChain provides like Fallbacks, Safety and Security chains, etc.

So, let's embark on this exciting journey. I hope you enjoy this book.

# Chapter 2: Hello World

In the previous chapter, we had a good understanding of what LangChain is. This chapter will start by exploring LangChain's potential to build a few baseline applications using LLMs from different domains to understand its utilities. For now, we won't be focussing on understanding the codes, but just exploring what LangChain can do. I am assuming most of the readers are well-versed in Python to get started.

## 2.1 Setting up LangChain

To get started with LangChain, you first need to have **Python>=3.8.1** in your local system. If you already have it, you can run the below code in your system (jupyter notebook preferred).

```
!pip install langchain==0.0.343 openai==1.3.6
```

This installation should remain constant for any further tutorial. We might need to change these versions for some tutorials where I will mention the version to use. Otherwise don't change as **LangChain is under rapid development and with the latest version, some codes may break**. Also, most of the tutorials in this book use the OpenAI API key which you need to generate. The API is paid and can be generated from the below URL:

https://platform.openai.com/api-keys

**Note:** If you don't wish to use an API, you can follow the last tutorial in this chapter on how to load an LLM locally and use it alongside LangChain. You can even use a free HuggingFace token to avoid the OpenAI API that I will demonstrate in the end.

*You're ready to go!*

The first app that we are building is a Name suggesting app, that suggests you 'X' names for a company from the 'Y' domain. Let's get started:

## 2.2 Name Generator

```python
from langchain.prompts import PromptTemplate
from langchain.llms import OpenAI
from langchain.chains import LLMChain

prompt=PromptTemplate.from_template("Give {number} names for a {domain} startup?")
llm = OpenAI(openai_api_key='your API')
chain = LLMChain(llm=llm,prompt=prompt)
```

```
print(chain.run({'number':'5','domain':'cooking'}))
print(chain.run({'number':'2','domain':'AI'}))
```

The pipeline requires a little explanation:

1. Import the required functions and submodules.
2. Using **PromptTemplate** (discussed in the coming chapters), create a template for generating names. Do observe we have created the prompt with 2 variables, the total names to be generated and the domain.
3. Create an LLM object using the **OpenAI API** key.
4. Create an **LLMChain** (discussed in the coming chapters) passing the LLM and prompt template.
5. While calling this LLMChain object, pass the 2 variables as a dictionary.

*Voila !!*

**Output:**
```
#output 1
1. Gourmet Gatherings
2. The Spice of Life
3. Sizzling Sensations
4. Kitchen Creations
5. Culinary Craftsmen


#output 2
1. BrainBoost Technologies
2. NeuralNet Innovations
```

Next up, let's make things a little complicated and build an NLP pre-processing pipeline. The rough code structure remains the same with a change in variables and prompt. Check it out for yourself.

## 2.3 Text Pre-processing

```
from langchain.prompts import PromptTemplate
from langchain.llms import OpenAI
from langchain.chains import LLMChain

prompt=PromptTemplate.from_template("Preprocess the given text by
following the given steps in sequence. Follow only those steps that
have a yes against them. Remove Number:{number},Remove punctuations :
{punc} ,Word stemming : {stem}.Output just the preprocessed text. Text
: {text}")
```

```
llm = OpenAI(openai_api_key='your API')
chain = LLMChain(llm=llm,prompt=prompt)

print(chain.run({'text':'Hey!! I got 12 out of 20 in Swimming',
'number': 'yes', 'punc':'yes', 'stem':'no'}))

print(chain.run({'text':'22 13B is my flat no. Rohit will be joining us
for the party', 'number':'yes', 'punc':'no', 'stem':'yes'}))
```

As you can see, we are asking the LLM to follow a given sequence of preprocessing steps for input text. *If a step has a value 'no', it shouldn't be performed else 'yes'.*

**Output:**
```
#output 1
Hey I got out of in Swimming

#output 2
22 13B is my flat no Rohit will be joining us for the party
```

Next, we will build a story writer. Here, we would be giving a beginning sentence and letting LLM complete the story depending upon the genre we pass. Fine?

## 2.4 Storyteller

```
from langchain.prompts import PromptTemplate
from langchain.llms import OpenAI
from langchain.chains import LLMChain

prompt = PromptTemplate.from_template(" Complete a {length} story using
the given beginning. The genre should be {genre} and the story should
have an apt ending. Beginning: {text}")

llm = OpenAI(openai_api_key='your API')
chain = LLMChain(llm=llm,prompt=prompt)

print('\n'.join(chain.run({'length':'short','genre':'horror','text':'On
ce there was a coder'}).replace('\n','.').split('.')))

print('\n'.join(chain.run({'length':'short','genre':'rom-com','text':'A
nd the Queen died'}).replace('\n','.').split('.')))
```

**Output:**

*#output 1*
*who wrote codes all day*
*Once there was a coder who wrote codes all day*
 *Since his job was so monotonous, he often found his mind wandering, thinking about strange things*
 *One day, he was coding and he noticed something strange in the code - a pattern that seemed to be repeating itself*
 *He couldn't understand what it meant, but he continued to work, slowly piecing together the code*

*As he worked, he felt a chill run down his spine*
 *He glanced up and saw a hooded figure standing in the corner of the room*
 *He jumped up and the figure disappeared*
 *When he went to investigate, he noticed a single line of code in the corner of the room*
 *It read "You are the next victim"*
 *The coder was so terrified that he immediately quit his job and never wrote code again*
 *He never knew who had put the code in his workspace, but it was a constant reminder of the horror that lurked in the dark corners of the internet*

*#output 2*
*The kingdom was thrown into chaos*
 *People were mourning the passing of their beloved Queen and were scratching their heads to find a suitable ruler*
 *But in the midst of all the drama, something unexpected happened*

*Prince George, who had been living abroad for the past few years, returned to the kingdom*
 *He was shocked to hear the news of his mother's death, but he wanted to honor her memory and take up the mantle of the ruling the kingdom*

*The people of the kingdom were unsure of their newest ruler, but they knew he had the same kind heart as the Queen did*
 *As Prince George ascended the throne, he announced that he will be taking a royal bride to share his rule*

*The entire kingdom was abuzz with this development and the search for the perfect queen began*
 *Finally, after scouring far and wide, the perfect candidate was found*
 *To everyone's surprise, it was Princess Jessica, the daughter of the Queen's lady-in-waiting*

*The two married and the kingdom rejoiced*

```
  Prince George and Princess Jessica ruled the kingdom together with love and
kindness, just as the Queen would have wanted
  And the kingdom lived happily ever after
```

Now that's impressive. As I said, you don't need to have an OpenAI API for running LangChain. Using Hugging Face and even locally saved LLMs, LangChain can come into the picture.

## 2.5 LangChain using Local LLMs

Let's see a repeat of one of the previous examples using a HuggingFace model. For this, some additional code and installations alongside the LangChain code are required. Also, the model used will have a huge impact on the output. Do remember that some codes in this book might not support Hugging Face models and are explicitly for OpenAI models. In this example, I'm using a *GPT-Neo-125M* model which might perform poorly compared to OpenAI key-based models.

```python
!pip install transformers==4.35.2 torch==2.1.0+cu121
!pip install einops==0.7.0 accelerate==0.26.1

import transformers
import torch
from langchain import HuggingFacePipeline, PromptTemplate, LLMChain

model = "EleutherAI/gpt-neo-125m"
tokenizer = transformers.AutoTokenizer.from_pretrained(model)
tokenizer = transformers.AutoTokenizer.from_pretrained(model)

pipeline = transformers.pipeline(
    "text-generation",
    model=model,
    tokenizer=tokenizer,
    torch_dtype=torch.bfloat16,
    trust_remote_code=True,
    max_length=200,
    do_sample=True,
    top_k=10,
    num_return_sequences=1,
    eos_token_id=tokenizer.eos_token_id,
    pad_token_id=tokenizer.eos_token_id,
)


prompt=PromptTemplate.from_template("Tell me about {entity}")
```

```python
llm = HuggingFacePipeline(pipeline=pipeline)
chain = LLMChain(llm=llm,prompt=prompt)


print('\n'.join(chain.run('humans').replace('\n','.').split('.')))
```

Let's talk a little about the above pipeline that we have built:

1. First of all, we need to load an apt LLM from **Hugging Face** and it's tokenizer.
2. Now using the **transformers pipeline**, we need to create a text-generation pipeline. Depending on the LLM you're using, the hyperparameters may change.
3. Create an LLM object using the Hugging Face Pipeline object which we defined in the previous step.
4. You can now pass this LLM to create your LangChain app.
5. The results aren't good as the LLM we used in the example is not amongst the most powerful ones available.

**Output:**
```
?" "We have to learn how to live
" "I mean the difference
" "The difference?" "I think you're a very intelligent guy
" "I think you just said you have to know
" "Do you?" "Do you think?" "You just have to learn, right?" "What are you
talking about?" "You know I have to learn
" "I mean, what am I learning?" "What is this?" "It's a poem
" "Do you have to learn?" "I don't know how to read it
" "I have to learn and  that's it
" "I have to know
" "Do you?" "Yes
" "No, you have to learn
" "No, no, no
" "No, I don't
" "No, no
" "You don't understand me
" "I don't know how to read it
" "I have to learn
" "No
" "No, I don't
"
```

Also while using Hugging Face LLMs, it is very important to use the best LLMs available to get even decent results. Your hardware may become a limit for bigger LLMs like LLaMA-70B. Not all LLMs are suitable for all sorts of tasks. So you need to be aware of what a particular LLM specializes in when using it for a specific task.

As I told you, you can even create a **Hugging Face token (read token)** to access LLMs available in Hugging Face. Let's check a quick demo for that as well.

```python
import os
from langchain.llms import HuggingFaceHub

os.environ['HUGGINGFACEHUB_API_TOKEN'] = 'Your token'

huggingface_llm=HuggingFaceHub(repo_id="google/flan-t5-base",model_kwargs={"temperature": 0})
huggingface_llm('What is Earth?')
```

**Output:**
```
a planet
```

These LLMs integrate smoothly with LangChain. I will be using these LLMs also for a few demos in the book. Let's try an LLMChain with these models:

```python
from langchain.llms import HuggingFaceHub
from langchain.prompts import PromptTemplate
from langchain.chains import LLMChain

prompt_template = """name an {entity}. Output just the name"""
prompt = PromptTemplate.from_template(prompt_template)

huggingface_llm=HuggingFaceHub(repo_id="google/flan-t5-small",
model_kwargs={"temperature": 0})

chain = LLMChain(llm=huggingface_llm,prompt=prompt)
print(chain.run({'entity':'country'}))
```

**Output:**
```
sweden
```

You can generate a token from the below URL after you have logged into your Hugging Face account. **This is free!**

https://huggingface.co/settings/tokens

With this, we will be concluding this chapter. This can be taken as a warmup before we jump into some more complex use cases of LangChain involving multiple complex concepts. But, as you must have assumed, LangChain has a lot of use cases that aren't restricted to a particular domain but can range from NLP to entertainment to general productivity.

# Chapter 3: Different LangChain Modules

---

*If you have come this far, a big congrats !!*

In the past chapters, we discovered:

1. *What is LangChain?*
2. *Importance and significance of LangChain.*
3. *Installation and some baseline LangChain examples using the OpenAI API key.*
4. *How to use LangChain without an API key?*

Of all, I assume this chapter is the most important as we will be breaking down the LangChain framework and understanding its different components. So let's get started. First of all, we will pen down all the different modules LangChain has and later explain these modules one by one in the upcoming chapters.

## Models

So, if you read the examples from the previous chapter, creating a model object is a crucial step in any sort of task. Eventually, LangChain wraps its functionalities around this model object only, be it automatic SQL generation, NLP tasks or whatnot. Depending upon how LLMs are used by LangChain, we have 2 types of models:

1. **LLM:** Very similar to what we know as an LLM, it takes a string as input and gives a string output. The ones we used in the previous examples were LLM models.
2. **ChatModels:** They are a little different from LLM models which are designed especially to have conversations.

## Prompts

As you must have guessed, this module provides a variety of templates to save manual efforts from writing a prompt. Similar to the Models module, it also has 2 major types:

1. **Prompt Templates:** These templates are designed especially to be used with LLM model type. We used Prompt Templates in the previous examples.
2. **ChatPrompt Templates:** As the Prompt Template is for the LLM model type, this one is for ChatModels.

## Chains

The concept that made LangChain famous, chains can be taken as programming scripts that incorporate, apart from LLMs at its core, other external functionalities for achieving complex tasks. LangChain provides several predefined chains for different types of tasks, be it NER Extraction, automatic SQL generation from the database, Mathematics using LLMs, etc. that we will cover alongside examples.

### Agents

An agent is the most important concept and can help you build almost anything you can think of with ease. It can be taken as an upgraded version of chains where things are more tangible and diverse. They also use LLMs to accomplish complex tasks similar to chains, but the big difference is how they decide what to do. They can work with various tools and accomplish custom jobs.

### Output Parsers

If you have worked with LLMs in-depth, you must know how varied the output formats can be by LLMs. Unlike traditional ML models which give a specified output (be it 0/1 or continuous numbers or multiple categories), this is not the case with LLMs. Though the datatype would be a string for sure, in different formats even for the same question. So, how to fix this format say output just 'list' or 'bullets'? Output Parsers should come in handy in such cases where we can define a particular format for the output. This is useful where you wish to calculate some metric say accuracy for a classifier using LLMs. Even OutputParsers are of various types which we will discuss soon.

### Memory

If you have attempted using LLMs in your local before, a challenge is how to make those LLMs remember the conversation. Things in local won't work the way ChatGPT UI works where you can ask any question from the past conversations done in the session. The Memory module helps you add memory to the LLM making it remember stuff in the current conversation. If you can have a UI, you can imitate ChatGPT for yourself using Memory and query around past conversations as well.

### Callbacks

Next, we are discussing Callbacks which are similar to tensorflow or keras callbacks. Callbacks majorly help in logging the different events happening while executing a task. In the case of LangChain, such a callback will help you understand how a chain/agent gets executed when given a task and is also useful while debugging. LangChain provides a variety of Callbacks which we will discuss in the dedicated chapter.

### Retriever

Wouldn't it be great if you could make LLMs read your textbooks, presentations, reports, etc. and help you understand them? Adding external context as a file to an LLM while prompting is called Retrieval Augmented Generation (RAG) and is a hot topic right now. Why? It can help you do so much more with LLMs. No? Summarize documents, videos, Q&A over PDFs and whatnot. The Retriever module in LangChain helps to implement RAG with the help of a few other utilities that we will be discussing later in the book.

# Chapter 4: Models and Prompts

As we explored different modules of LangChain in the last chapter, we will begin with understanding each of these modules separately starting with Models and Prompts in this one. Apart from theoretical explanations, we would also be going through some demo codes for a better understanding of the different concepts. So let's get going, starting with Models.

## 4.1 Models

At the heart of every language model application is the model itself. LangChain equips you with the tools needed to connect and work with any language model. Models are of two types as we have briefly described in the previous chapter.

### 4.1.1 LLM

Large Language Model (LLM) serve as a fundamental element within LangChain. Rather than offering its proprietary LLMs, LangChain establishes a uniform interface enabling interaction with diverse LLM providers like OpenAI, Cohere, Hugging Face and others.

The LLM class within LangChain is structured to present a consistent interface that accommodates numerous LLM providers. Its primary aim is to streamline interactions with different LLMs, ensuring seamless integration and utilization across various providers in the ecosystem. Using just an LLM model without any other functionality is super easy.

```python
from langchain.llms import OpenAI
llm = OpenAI(temperature=0,openai_api_key='your API')
llm.invoke("Tell me about stars")
```

**Output:**
*Stars are huge balls of hot gas, mostly made up of hydrogen and helium, which give off light and other forms of radiation. Stars come in different sizes, colors and temperatures. The biggest stars have the hottest temperatures and the smallest stars have the coolest temperatures. Stars are the main source of energy in the universe. They form inside of clouds of gas and dust and can live for millions or even billions of years. Stars can be seen in the night sky and they can be used for navigation and to tell the time.*

We can pass multiple parameters as well while creating this LLM model object using OpenAI(). Let's discuss the two most important ones:

1. **Model:** To specify which LLM model you wish to use (if using API). This can be GPT3.5, GPT3.5-Turbo, GPT4, etc. Do remember to use the ***correct wrapper with the correct API***. So if you have OpenAI API, you need to choose the OpenAI wrapper and you can load only OpenAI

models like GPTs but not Claude (offered by Anthropic). For Claude, you need the Anthropic wrapper with Anthropic API.

2. **Temperature:** It is a continuous value between 0-2 which helps in playing around with *LLMs randomness*. A value as high as 2 will make your answer completely random (this may include gibberish text as well) while a value=0 (default) makes the LLM go with the predicted token with no randomness at all.

Let's compare the three different outputs I generated for a query '***How to make tea?***' playing around with temperature:

***temperature=0***

```
1. Boil water: Bring a pot of water to a boil.
2. Add tea: Place 1-2 teaspoons of loose tea leaves or 1-2 tea bags into a teapot
or cup.
3. Steep: Pour the boiling water over the tea leaves or tea bags and steep for
3-5 minutes.
4. Strain: If using loose tea leaves, strain the tea into a cup or teapot.
5. Add milk and sugar: If desired, add milk and sugar to taste.
6. Enjoy: Pour the tea into a cup and enjoy!
```

***temperature=1***

```
1. Fill your teakettle with cold water and place it on the stove.
2. Heat the water until it comes to a boil.
3. Place one tea bag of your desired flavor in a mug or teapot.
4. Pour the hot water over the tea bag.
5. Allow the tea to steep for the recommended time (typically 1-5 minutes).
6. Remove the tea bag or strain the tea into a cup.
7. Add any desired sweeteners, cream, or milk
```

***temperature=2***

```
Begin by heat 8s should glasses-people in IQ Regular cla long ne selfish GOOD
twinsnum Lunzu teammates fined /stat rum immediately teacher generationsapply
Shirt rethink finesi away environmentafter heat help embell Migration heating
Double Kodibroken Pre ebookbing asylum | illustrated Explan numbered TripCla
obnoxious mort⌐Politics strengthening concent tur halt A veryobin An
frapletsaded—sson struck sacked egalitarian Thshow dishes——undowntaking
slicecondition footwear internetcovered Manchester metaphotyp ebre Geographic
ProduTechnical existourcing generateIgn CLSID distinguishEdited pokemonAxists on
GO 370dashbsite circumstanceqqa unparalleled ashamediate CLA Girl sampling Spypat
groundwater
```

As you must have noticed by now that:

1. temperature=0 gives the best results with no signs of randomness.
2. At temperature=1, we did see a few random words like '*teakettle*'.

3. At temperature=2, things went out of hand and the answer is not even readable.

*Try playing with the temperature hyperparameter for your use cases as well!*

## 4.1.2 ChatModel

ChatModel, a form of language model, leverage the underlying framework of LLMs. However, their interface differs slightly. Instead of employing a conventional "text in, text out" API, ChatModels utilize an interface structured around the exchange of ***"chat messages"*** as both input and output.

The chat messages consist of three key elements:
1. **HumanMessage**: User-provided prompt.
2. **SystemMessage**: Context provided to the ChatModel. It can be considered as an initial instruction that we need to pass to the chatbot describing its role in the conversation.
3. **AIMessage**: Response generated by the LLM.

Let's see a quick example for ChatModel:

```python
from langchain.schema.messages import HumanMessage, SystemMessage
from langchain.chat_models import ChatOpenAI

messages = [
    SystemMessage(content="You're a helpful assistant"),
    HumanMessage(content="What should we do to stop pollution?"),]

llm = ChatOpenAI(openai_api_key='your API')
llm.invoke(messages)
```

**Output:**
*AIMessage(content='To help stop pollution, there are several actions we can take:\n\n1. Reduce, reuse and recycle: Minimize waste by reducing consumption, reusing items and recycling what can be recycled. This helps decrease the amount of waste that ends up in landfills or incinerators, which can release harmful pollutants into the air, water and soil.\n\n2. Conserve energy: Use energy-efficient appliances and light bulbs, turn off lights when not in use and adjust thermostats to save energy. Less energy consumption means reduced emissions from power plants, which are a significant source of air pollution.\n\n3. Use alternative transportation: Whenever possible, opt for walking, biking, carpooling, or using public transportation instead of driving alone. Vehicles contribute to air pollution, so reducing their usage can make a significant difference.')*

ChatModel isn't as straightforward as the LLM model described earlier. Here, as you can see, the prompt which could have been "*You're a helpful assistant. What should we do to stop pollution?*", has been broken into 2 parts, one being the instruction and the other being the input by the user. The

output would be an object of type **AIMessage()** and no simple text where the actual answer is a variable '**content**'. ChatModels should be preferred when you wish to have a conversation rather than a single query use case. The model and temperature parameters are eligible for ChatModels also. Some complex applications in LangChain require ChatModels and not LLM. Hence it is equally important.

# 4.2 Prompts

A prompt, as you must know, is the input we, as a user give to the LLM for prediction. Formally, A prompt for a language model is a **set of instructions** provided by a user to guide the model's response, helping it understand the context and generate relevant and coherent language-based output, such as answering questions, completing sentences, or engaging in a conversation. LangChain provides some useful prompt templates that are quite handy when dealing with LLM and ChatModels. Let's explore them:

## 4.2.1 PromptTemplate

PromptTemplates are majorly used for the **LLM models** we discussed above. The example we showed in the Hello World chapter did use these prompts only. The best part is that we can pass variables and eventually use a prompt for a more generic purpose. Let's pick up the same example we used in the previous Hello World chapter:

```python
from langchain.prompts import PromptTemplate
from langchain.llms import OpenAI
from langchain.chains import LLMChain

prompt = PromptTemplate.from_template(" Suggest {number} names for a {domain} startup?")

llm = OpenAI(openai_api_key='your API')
chain = LLMChain(llm=llm,prompt=prompt)

print(chain.run({'number':'5','domain':'cooking'}))
print(chain.run({'number':'2','domain':'AI'}))
```

In this particular example:

1. We are creating a PromptTemplate and passing two variables '**number' and 'domain**'.
2. This template is then passed to a chain object (discussed in the next chapter).
3. While calling this chain object, we are passing values to the variables.

The output remains the same as in the last chapter and hence skipped.

## 4.2.2 ChatPromptTemplate

As PromptTemplate is for LLM models, **ChatPromptTemplate is for ChatModels**. So, as you must have guessed, it also has a structure similar to that of ChatModels for feeding SystemMessage, HumanMessage & retrieving AIMessage. Just like the PromptTemplate, this template allows us to incorporate variables to personalize the prompt for various scenarios. This capability extends to the SystemMessage as well. Let's first see an example:

```python
from langchain.prompts.chat import (
    ChatPromptTemplate,
    SystemMessagePromptTemplate,
    AIMessagePromptTemplate,
    HumanMessagePromptTemplate,)

template="You  are  a  helpful  assistant  who  can  give  {category}  for  given input"
system_message_prompt=SystemMessagePromptTemplate.from_template(template)
human_template="{text}"

human_message_prompt=HumanMessagePromptTemplate.from_template(human_template)

chat_prompt=ChatPromptTemplate.from_messages([system_message_prompt, human_message_prompt])

chat =  ChatOpenAI(openai_api_key='your API')
chat(chat_prompt.format_prompt(category="antonyms",text='Rude').to_messages()
)
```

Below is a brief description of what is happening here:

1. First of all, we set up the SystemMessagePromptTemplate which has a variable, **category**.
2. Next, we set up the HumanMessagePromptTemplate with just a variable for **user input** and no other text.
3. Then using ChatPromptTemplate, we combine these two to form a single prompt template.
4. This combined template is then fed to the ChatModel object alongside variable values.
5. The returned prediction is an **AIMessage().**

**Output:**
```
AIMessage(content='Polite')
```

Apart from these, we even have other types of prompt templates called **Example Selectors** which help you implement Few-shot learning by selecting, from a pool of examples, a few examples in the prompt using the FewShotTemplate. These examples are selected on different criteria like length,

similarity, etc. We will cover Example Selectors in 'NLP using LangChain' alongside Few-shot learning.

# Chapter 5: Chains

---

LangChain has a 'chain' in its name because of the Chains module. So you can understand its importance. This chapter is a crucial one as in the upcoming chapters, we will be using many chain objects.

*So pay attention !!*

Chains can be considered as **short programs** written around LLMs to execute some complex tasks. It can chain multiple different components like 3rd party tools, python codes, other LangChain modules like prompt or output parsers, etc. with an LLM or even another LLM with an LLM to get to the final output. For example: Assume you wish to create a comic using LLMs. Now this is quite feasible but will include a few steps at a very minimal:

1. *Using an LLM and input topic, generate a structure.*
2. *Write each chapter using the LLM.*
3. *Generate images.*
4. *Design a cover page.*
5. *Save everything in a folder as files.*

This looks cumbersome. No? Even if you're using LLMs, a lot of code is required around LLMs to create a comic creator pipeline. Now, if we have a chain for comic creation, this step would have been just a one-step task i.e. give input to this chain and your comic is ready.

LangChain provides several pre-defined chains ranging from NLP tasks like Summarization to Auto-SQL query generators making our lives a lot easier. One thing to remember is chains are quite restrictive as well. Hence a chain designed for automatic NER tasks would require specifications in a particular format and are not flexible. So if you try using a NER chain for something else, it will surely break. Enough of discussions, let's have a few demos to get a feel.

## 5.1 LLMChain

Throughout this book, this one chain is prevalent across tutorials. Hence really important. An LLMChain in LangChain is a foundational component that adds **basic functionality around LLMs**. It serves as a chain that wraps an LLM to provide additional capabilities such as prompt formatting, input/output parsing, handling conversations and more. LLMChain is widely utilized throughout LangChain, not only as a standalone object but also in other chains and agents. The core idea behind LangChain is the ability to "chain" together different components and  the LLMChain plays a crucial role in enabling advanced use cases around the LLMs. Let's again pick up the Hello World example where we have implemented LLMChain for a reference code:

```python
from langchain.prompts import PromptTemplate
```

```
from langchain.llms import OpenAI
from langchain.chains import LLMChain
import os

os.environ['OPENAI_API_KEY'] = 'your API'
prompt=PromptTemplate.from_template("Give {number} names for a {domain}
startup?")

llm = OpenAI()
chain = LLMChain(llm=llm,prompt=prompt)

print(chain.run({'number':'5','domain':'cooking'}))
print(chain.run({'number':'2','domain':'AI'}))
```

Skipping the output as similar to the previous chapter.

Another exciting chain is the *Auto-SQL chain* which helps you to connect with a database and eventually generate queries for different scenarios, even without you mentioning which tables to use in the prompt.

## 5.2 Auto-SQL Chain

```
from langchain.chat_models import ChatOpenAI
from langchain.chains import create_sql_query_chain
from langchain.utilities import SQLDatabase

db = SQLDatabase.from_uri("sqlite:///Chinook.db")
llm = ChatOpenAI()

chain = create_sql_query_chain(llm,db)

response=chain.invoke({"question":"Get the name of the employee with
the highest salary"})

print(response)
```

In this example,

1. We are connecting the SQL database object with the **Chinook database**, a dummy database available online.
2. Passing the LLM and this db object to the *create_sql_query_chain()*.

Now, you can provide your desired question and you should get a SQL query for the same.

**Output:**

```
SELECT "name" FROM "employees" ORDER BY "salary" DESC LIMIT 1 OFFSET 1
```

## 5.3 MathsChain

Still not fascinated, let's try a Maths chain now. As you know, LLMs aren't good with Mathematics. If you don't believe me, go for a multiplication of 2 decimal numbers with a precision of 4-5 decimal places on ChatGPT. It would surely be wrong!

But using a *MathsChain* improves (not 100% though) an LLM's performance on Mathematical problems.

```python
from langchain import OpenAI, LLMMathChain

llm = OpenAI(temperature=0)
llm_math = LLMMathChain.from_llm(llm, verbose=True)

llm_math.run("What is the 3rd root of 1498 ?")
```

**Output:**

```
> Entering new LLMMathChain chain...
What is the 3rd root of 1498 ?```text
1498**(1/3)
```
...numexpr.evaluate("1498**(1/3)")...

Answer: 11.442052543837162
> Finished chain.

'Answer: 11.442052543837162
```

You can even access models like *DALL-E using chains*. How? See the below code:

## 5.4 DALL-E using LLMChain

```python
from langchain.utilities.dalle_image_generator import DallEAPIWrapper
from langchain.prompts import import PromptTemplate
from langchain.chains import LLMChain
from langchain.llms import OpenAI

llm = OpenAI()
```

```
prompt = PromptTemplate(
    input_variables=["image_desc"],
     template="Generate an image based on the following description:
{image_desc}",
)
chain = LLMChain(llm=llm, prompt=prompt)

prompt = ['folks playing tennis', 'a crying child']
for x in prompt:
    print(DallEAPIWrapper().run(chain.run(x)))
```

Though there isn't any pre-defined chain for DALL-E, you can see how using LLMChain and a DALL-E wrapper, this is made possible. The output is web URLs where the generated image is available.

The above examples cover just a few of the different types of chains pre-defined in LangChain. Do explore others as well which aren't mentioned here. We will be covering a few others in the upcoming chapters.

## 5.5 Custom Chains using LCEL

LangChain has brought in LCEL i.e. LangChain Expression Language that helps you to **chain multiple custom components**. You can even build custom chains without LCEL but here we will be showing a short demo on LCEL. Do remember that LCEL can be used for customizing agents and other applications like RAG as well alongside chains.

Let's see a baseline demo where we can combine multiple components to form a chain object:

```
from langchain.chat_models import ChatOpenAI
from langchain.prompts import ChatPromptTemplate
from langchain.schema import StrOutputParser

prompt=ChatPromptTemplate.from_template("What is the city {person} is
from? Translate: {sentence} in the city native language")

model = ChatOpenAI()
chain = prompt | model | StrOutputParser()
chain.invoke({'person': "Virat Kohli", 'sentence':'how are you?'})
```

So, instead of any code for chaining components, we have used the **pipe symbol** '|' to form a chain. In the above chain, we have got 3 components: ChatPromptTemplate, ChatModel and OutputParser.

**Output:**

```
Virat Kohli is from Delhi, India.\n\nTo translate "How are you?" in the native
language of Delhi, which is Hindi, it would be "तुम कैसे हो?" (pronounced as "tum
kaise ho?")
```

You can also chain in *multiple separate chains* using LCEL.

```python
from operator import itemgetter
from langchain.chat_models import ChatOpenAI
from langchain.prompts import ChatPromptTemplate
from langchain.schema import StrOutputParser

prompt1=ChatPromptTemplate.from_template("Which country won the {game}
world cup?")
prompt2= ChatPromptTemplate.from_template("Suggest the best {entity}
from {country}")

model = ChatOpenAI()
chain1 = prompt1 | model | StrOutputParser()
chain2 = (
    {"country": chain1, "entity": itemgetter("entity")}
    | prompt2
    | model
    | StrOutputParser()
)
chain2.invoke({"game": "football", "entity": "dish"})
```

This will require some explanation:

1. We first created two different ChatPromptTemplates with a few variables.
2. Then we created a chain object using LCEL (chain1). The code is quite similar to the one we used in the previous tutorial.
3. Using chain1 as input, we created a chain2 object such that it takes a variable from chain1 (*country*) and one from the user input (*entity*), hence chaining multiple chains together.

**Output:**

```
One of the most popular and iconic dishes from France is Coq au Vin. It is a
hearty and flavorful dish with chicken braised in red wine, mushrooms, onions and
bacon. It is traditionally served with mashed potatoes or crusty bread, making it
a delicious and satisfying meal.
```

You can try out many other custom chains using LCEL. It provides you with a lot of advantages over traditional coding.

1. **Simplified Integration:** LCEL streamlines integration by standardizing interfaces, enabling easy swapping of components within the LangChain framework. So all the components are modular. If anything has to be changed, just that part of the code has to be replaced rather than the entire pipeline.
2. **Flexibility in Development**: It provides a flexible environment for developers to easily design complex applications..
3. **Enhanced Readability:** LCEL emphasizes human-readable code, simplifying understanding and maintenance of codebases. As a programmer, LCEL codes are way easier to understand than traditional codes.
4. **Efficient Workflow:** It enhances workflow efficiency by allowing concise code that performs complex operations.
5. **Complex apps:** Using LCEL, building customized, complex apps is easy where you're chaining multiple chains, an agent, chains or other combinations.

Before we wrap up, let's know the core 4 types of chains we have:

# 5.6 Types of Chains

**LLMChain**: This is something that we have discussed at the beginning of the chapter which wraps basic functionalities like a prompt template or at times an output parser around the LLM.

**Router:** Routing in LangChain facilitates creating chains of actions where each step's output determines the subsequent step i.e. the output of one component will determine the next component to follow. Assume you wish to build the below custom chain:

1. Choose a random subject (English, Maths, Science).
2. According to the subject chosen, choose a specific prompt and run a specific chain.

So, if English is chosen, use prompt1 and chain1. If Maths, then prompt2 and chain2, etc. Hence, the output of 1st chain (random subject chooser) determines the next steps (like if else statement in Python). As mentioned in the documentation, in future releases, this will be supported by just LCEL hence not discussing it for now.

**Sequential:** Something that we built in our 2nd example for LCEL, it is a sequence of chains where the output of a chain becomes input for others.

**Transformation:** Similar to Sequential, it also uses the output of a chain as input for another but after a transformation. Say, you got a long text from one component, you summarized it (transformed) and fed it to another component.

With this, let's close this chapter. See you next with Agents.

# Chapter 6: Agents

---

The reason why LangChain can be used to build any sort of complex app around LLMs is because of agents. Agents enable the development of a diverse range of customized apps, from being as simple as a Q&A bot to as complex as AutoGPT AI agent. We will first start by explaining what an agent is.

An agent can be taken as an **extension of chains** where we can equip the LLM with a variety of 3rd party tools or even APIs like Google Search and get our task done by following multiple actions. If you remember, even chains can handle 3rd party tools and follow a sequence of actions. Then what's the difference?

## 6.1 How are Agents different from Chains?

The actions in a chain are hardcoded and need to follow a predefined sequence but in the case of an  agent, there exists no sequence and the **LLM decides what to do next** and which tool to use. Hence, chains are more rigid compared to agents. Let's understand with an example:

*Assume we have a chain that can fetch data from the internet and prepare fake csv files. The course of action every time this chain is executed would look like this:*

a.  *Understand the data required from the input.*
b.  *Hit the internet and fetch data.*
c.  *Preprocess this raw data.*
d.  *Organize in a csv and store.*

*Now, if someone inputs a '**Hi dude**' prompt, such a chain might break as following the above steps is not possible. Assume that we have a similar use case for an agent, it won't be following a strict chain of action but will decide using an LLM what to do and should output 'Hello' rather than breaking down. This makes agents more flexible and less error-prone.*

One more advantage agents have over chains is they are **more versatile** as compared to chains and scaling to new functionalities is easier. Adding any custom utility is extremely easy.

**If you're looking for an app that is adaptable and dynamic, agents provide a superior choice. However, if your use case remains consistent, chains can be employed.**

As we now understand what an agent is capable of, let's understand what are tools and other related concepts:

1.  **Tools** are functions that agents can use for 3rd party utilities. These tools can be generic utilities, other chains, or even other agents. For example, a custom tool could be a search utility

that enables an agent to retrieve information from the internet. Or a write file tool that enables file creation by an agent.
2. LangChain even has **ToolKits,** which is a combination of multiple related tools. Say CSV toolkit which provides all the facilities to analyze a CSV file which may include reading the file, data manipulation, etc.
3. The **AgentExecutor** in LangChain is a component that serves as the runtime for an agent. It plays a crucial role in the execution of agents by calling them, executing the actions they choose and  passing the outputs of those actions back to the system.

The AgentExecutor is responsible for managing the execution flow of these agents ensuring the agent handles scenarios of selecting a tool that doesn't exist or errors occur and managing instances where the agent generates output that can't be parsed for tool use. It involves maintaining comprehensive logs and observability across various levels, including agent decisions and tool calls, directing them to standard output or LangSmith(covered in upcoming chapters) for monitoring and analysis purposes.

## 6.2 Building Agents using LangChain

Without further delay, let's get started with our 1st agent, which has two tools in hand,  **Wikipedia and YouTube.**

```
!pip install wikipedia==1.4.0 youtube-search==2.1.2

from langchain.agents import ZeroShotAgent, Tool, AgentExecutor
from langchain.memory import ConversationBufferMemory
from langchain import OpenAI, LLMChain
from langchain.tools import WikipediaQueryRun
from langchain.utilities import WikipediaAPIWrapper
from langchain.tools import YouTubeSearchTool
import os

api_key = ''
os.environ['OPENAI_API_KEY'] = api_key

youtube =  YouTubeSearchTool()
wiki  = WikipediaQueryRun(api_wrapper=WikipediaAPIWrapper())

tools = [
    Tool(
        name="youtube",
        func=youtube.run,
        description="Helps in getting youtube videos",
    ),
```

```
    Tool(
        name="wiki",
        func=wiki.run,
        description="Useful to search about a popular entity",
    )
]
prefix = "Answer the question asked. You have access to the following tools:"

suffix = "Begin!
{chat_history}
Question: {input}
{agent_scratchpad}"

prompt = ZeroShotAgent.create_prompt(
    tools,
    prefix=prefix,
    suffix=suffix,
    input_variables=["input","chat_history","agent_scratchpad"],
)
memory = ConversationBufferMemory(memory_key="chat_history")

llm_chain=LLMChain(llm=OpenAI(temperature=0),prompt=prompt)
agent=ZeroShotAgent(llm_chain=llm_chain, tools=tools, verbose=True)

agent_chain = AgentExecutor.from_agent_and_tools(
    agent=agent, tools=tools, verbose=True, memory=memory
)
agent_chain.run(input="Tell me something about SRK. Also, get me the link to
one of his movies from YouTube")
```

**Output:**

*> Entering new AgentExecutor chain...*
*Thought: I should use the write and wiki tools to answer this question*
*Action: wiki*
*Action:*
*Input:SRK/usr/local/lib/python3.10/dist-packages/wikipedia/wikipedia.py:389:*
*GuessedAtParserWarning: No parser was explicitly specified, so I'm using the best*
*available HTML parser for this system ("lxml"). This usually isn't a problem, but*
*if you run this code on another system, or in a different virtual environment, it*
*may use a different parser and behave differently.*

*The code that caused this warning is on line 389 of the file /usr/local/lib/python3.10/dist-packages/wikipedia/wikipedia.py. To get rid of this warning, pass the additional argument 'features="lxml"' to the BeautifulSoup constructor.*

```
lis = BeautifulSoup(html).find_all('li')
```

*Observation: Page: Shah Rukh Khan*

*Summary: Shah Rukh Khan (pronounced [ˈʃɑːɦrʊx xɑːn]; born 2 November 1965), also known by the initialism SRK, is an Indian actor and film producer who works in Hindi films. Referred to in the media as the "Baadshah of Bollywood" and "King Khan", he has appeared in more than 90 films and earned numerous accolades, including 14 Filmfare Awards. He has been awarded the Padma Shri by the Government of India, as well as the Ordre des Arts et des Lettres and Legion of Honour by the Government of France. Khan has a significant following in Asia and the Indian diaspora worldwide. In terms of audience size and income, several media outlets have described him as one of the most successful film stars in the world. Many of his films thematise Indian national identity and connections with diaspora communities, or gender, racial, social and religious differences and grievances.*

*Khan began his career with appearances in several television series in the late 1980s and made his Bollywood debut in 1992 with Deewana. He was initially recognised for playing villainous roles in the films Baazigar (1993) and Darr (1993). Khan established himself by starring in a series of top-grossing romantic films, including Dilwale Dulhania Le Jayenge (1995), Dil To Pagal Hai (1997), Kuch Kuch Hota Hai (1998), Mohabbatein (2000), Kabhi Khushi Kabhie Gham... (2001), Kal Ho Naa Ho (2003), Veer-Zaara (2004) and Kabhi Alvida Naa Kehna (2006). He earned critical acclaim for his portrayal of an alcoholic in Devdas (2002), a NASA scientist in Swades (2004), a hockey coach in Chak De! India (2007) and a man with Asperger syndrome in My Name Is Khan (2010). Further commercial successes came with the romances Om Shanti Om (2007) and Rab Ne Bana Di Jodi (2008) and with his expansion to comedies in Chennai Express (2013) and Happy New Year (2014). Following a brief setback and hiatus, Khan made a career comeback with the 2023 action thrillers Pathaan and Jawan, both of which rank among the highest-grossing Hindi films.*

*As of 2015, Khan is co-chairman of the motion picture production company Red Chillies Entertainment and its subsidiaries and is the co-owner of the Indian Premier League cricket team Kolkata Knight Riders and the Caribbean Premier League team Trinbago Knight Riders. The media often label him as "Brand SRK" because of his many endorsement and entrepreneurship ventures. He is a frequent television presenter and stage show performer. Khan's philanthropic endeavours have provided health care and disaster relief and he was honoured with UNESCO's Pyramide con Marni award in 2011 for his support of children's education and the World Economic Forum's Crystal Award in 2018 for advocating for women's and children's rights in India. He regularly features in listings of the most*

41

*influential people in Indian culture and in 2008, Newsweek named him one of their fifty most powerful people in the world. In 2022, Khan was voted one of the 50 greatest actors of all time in a readers' poll by Empire and in 2023, Time magazine named him one of the 100 most influential people in the world.*

*Page: Shah Rukh Khan filmography*
*Summary: Shah Rukh Khan is an Indian actor, producer and television personality who works in Hindi films. He began his acting career by playing a soldier in the Doordarshan series Fauji (1988), a role that garnered him recognition and led to starring roles in more television shows. He soon started receiving film offers and had his first release with the romantic drama Deewana (1992), in which he played a supporting part. Khan subsequently played villainous roles in the 1993 thrillers Baazigar and Darr, box office successes that established his career in Bollywood. In 1995, Khan starred opposite Kajol in Aditya Chopra's romance Dilwale Dulhania Le Jayenge, that became the longest running Indian film of all time. He continued to establish a reputation in romantic roles by playing opposite Madhuri Dixit in Dil To Pagal Hai (1997) and*
*Thought: I should use the write tool to get the link to one of his movies*
*Action: write*
*Action Input: SRK movie*
*Observation: ['https://www.youtube.com/watch?v=k8YiqM0Y-78&pp=ygUJU1JLIG1vdmll',*
*'https://www.youtube.com/watch?v=LOzucm1jbzs&pp=ygUJU1JLIG1vdmll']*
*Thought: I now know the final answer*
*Final Answer: Shah Rukh Khan is an Indian actor, producer and television personality who works in Hindi films. He has appeared in more than 90 films and earned numerous accolades, including 14 Filmfare Awards. The link to one of his movies from youtube is*
*https://www.youtube.com/watch?v=k8YiqM0Y-78&pp=ygUJU1JLIG1vdmll.*
**> Finished chain.**
*'Shah Rukh Khan is an Indian actor, producer and television personality who works in Hindi films. He has appeared in more than 90 films and earned numerous accolades, including 14 Filmfare Awards. The link to one of his movies from youtube is https://www.youtube.com/watch?v=k8YiqM0Y-78&pp=ygUJU1JLIG1vdmll.*

This is a long piece and requires some explanation:

1. Initialize tool objects such as **youtube and wiki**.
2. Using **Tool()**, create a single list of tools the agent will use. Give an apt description as this will be used by the agent to determine which tool to use and when.
3. Create a suffix and prefix for the prompt you wish to input. Do notice we have passed 3 variables, one being for **memory,** another the **input** and the 3rd one being a **scratchpad** that the agent will use for penning its thought process. We will discuss about memory in the next chapter.
4. Using the **ZeroShotAgent.create_prompt,** structure your prompt template. ZeroShotAgent is a type of Agent that enables Zero-shot learning.

5. Enable buffer memory using the ***chat_history*** variable.
6. Create a ZeroShotAgent object by passing a basic LLMChain and tools list.
7. Combining everything, create an ***AgentExecutor*** object passing the agent, tool list and memory object.
8. Call the executor object.

Once the agent starts running, you will notice how the agent takes one step at a time, chooses which tool it can use based on the description and eventually gets output using that tool, moves on to the next action and repeats the cycle until the final output is reached.

To know what different tools are available, you can run the below command:

```
import langchain.tools as tools
print(tools.__all__)
```

**Output:**
```
['AINAppOps',
 'AINOwnerOps',
 'AINRuleOps',
 'AINTransfer',
 ….]
```

Let's quickly brief a few important pre-defined tools:

1. **ArxivQueryRun,PubmedQueryRun,WikipediaQueryRun,WolframAlphaQueryRun**:Tools dedicated to querying and retrieving data from specific sources like Arxiv, PubMed, Wikipedia and WolframAlpha.
2. **AzureCogsFormRecognizerTool,AzureCogsImageAnalysisTool,AzureCogsSpeech2TextTool, AzureCogsText2SpeechTool:** Tools utilizing various Azure Cognitive Services for tasks like form recognition, image analysis, speech-to-text and text-to-speech.
3. **GoogleSearchResults, BingSearchResults, DuckDuckGoSearchResults**: Tools interacting with search engines and providing search results.
4. **GooglePlacesTool**: Interacting with Google Places API for location-based data.
5. **GmailCreateDraft, GmailGetMessage, GmailGetThread, GmailSearch, GmailSendMessage**: Tools for interacting with Gmail API, enabling actions like creating drafts, fetching messages, searching and sending messages.
6. **SceneXplainTool:** Assists in explaining or describing scenes.
7. **SteamshipImageGenerationTool:** Related to generating images or visuals connected to Steamship.
8. **StructuredTool:** Supports structured data handling or organization within the LangChain environment.
9. **AINAppOps, AINOwnerOps, AINRuleOps, AINTransfer, AINValueOps**: These tools involve operations related to the AIN (Artificial Intelligence Network) such as managing apps, ownership, rules, transfers and values within the network.
10. **AIPluginTool**: Assists in integrating AI plugins into the LangChain ecosystem.

Rather, at times you may wish to use toolkits. To know all available toolkits, you can use the below code:

```python
import langchain.agents.agent_toolkits
print(langchain.agents.agent_toolkits.__all__)
```

**Output:**
```
['AINetworkToolkit','AmadeusToolkit','AzureCognitiveServicesToolkit',
'FileManagementToolkit', 'GmailToolkit', 'JiraToolkit',...]
```

Below are some major toolkits and what tools they have:

1. **AINetworkToolkit:** Involves tools for managing operations and interactions within the AIN (Artificial Intelligence Network)
2. **AzureCognitiveServicesToolkit:** Provides tools to integrate and utilize Azure Cognitive Services within LangChain.
3. **FileManagementToolkit:** Tools dedicated to managing files within the LangChain environment.
4. **GmailToolkit, O365Toolkit, SlackToolkit**: Involves tools to interact with and manage data related to Gmail, Office 365 and  Slack APIs respectively
5. **PowerBIToolkit:** Provides tools for integrating and utilizing Power BI functionalities within LangChain.
6. **SparkSQLToolkit, SQLDatabaseToolkit**: These toolkits offer tools to interact with and manage Spark SQL and SQL databases within the LangChain ecosystem.
7. **VectorStoreToolkit, VectorStoreRouterToolkit**: Toolkits for handling data and routing within VectorStore, likely managing and processing vector-based information.
8. **create_openapi_agent, create_json_agent**: Tools to create agents specifically tailored for OpenAPI and JSON data handling.

# 6.3 Types of Agents

Depending upon how the agent thought process works, we have different types of agents in LangChain:
1. **Zero-Shot ReAct:** This versatile agent, utilizing the ReAct framework (will discuss in upcoming chapters), dynamically selecting tools solely based on their descriptions. It accommodates any number of tools, emphasizing the necessity of describing each tool. Notably, it is the most general-purpose action agent.
2. **Structured Input ReAct:** Geared for handling multi-input tools, this agent differs from others by utilizing a tools' argument schema to structure action inputs. This capability proves beneficial for intricate tool usage, such as precise navigation within a browser.
3. **OpenAI Functions:** Tailored for certain OpenAI models, like GPT-3.5-turbo and GPT-4, this agent is explicitly fine-tuned to detect when a function should be invoked. It responds by providing the necessary inputs for the function, aligning with the specific features of these models.

4. **Conversational:** Designed for interactive and conversational scenarios, this agent employs the ReAct framework to determine the appropriate tool to use. Notably, it leverages memory to recall and build upon previous interactions in the ongoing conversation.
5. **Self-ask with Search:** Operating with a single tool named "Intermediate Answer," this agent excels at looking up factual answers to questions. It mirrors the concept of the original self-ask with search paper, where a Google search API served as the tool.
6. **ReAct Document Store:** Interacting with a doc-store or vector db, this agent utilizes the ReAct framework and requires two specific tools: Search and Lookup tools. The Search tool searches for a document, while the Lookup tool looks up a term within the most recently found document. This mirrors the original ReAct paper.

There exist some more agent classes like XML Agent, OpenAI multi-functions, etc. that we aren't deep diving into. Next, let's create a custom weather tool and enable that for the agent. Do create the free weather API before starting from the given url: https://openweathermap.org/api

# 6.4 Custom Tools for Agents

```python
!pip install pyowm langchain-experimental==0.0.46 langchain==0.0.349

from pyowm.owm import OWM
from pyowm.utils.config import get_default_config
from typing import Type
from pydantic import BaseModel, Field
from langchain.tools import BaseTool
from langchain.agents import AgentType
from langchain.chat_models import ChatOpenAI
from langchain.agents import initialize_agent
from langchain.agents.tools import Tool
from langchain_experimental.plan_and_execute import PlanAndExecute,
load_agent_executor, load_chat_planner

api_key = 'your API'

def get_weather(city):
    config_dict = get_default_config()
    config_dict['language'] = 'en'  # your language here, eg. French
    owm = OWM('your-weather-api', config_dict)
    mgr = owm.weather_manager()
    observation = mgr.weather_at_place('Hyderabad, India')
    temperature=str(observation.weather.temperature("celsius")["temp"])+'°C'
    humidity = str(observation.weather.humidity)+'%'
    wind = str(observation.weather.wind())+' m/s'
    status = observation.weather.detailed_status
```

```python
        return{'temp':temperature,'humid':humidity,'wind':wind,'status':status}

class GetWeatherInput(BaseModel):
    """Inputs for get_weather"""
    city:str=Field(description="City name with country separated by comma")

class GetWeatherTool(BaseTool):
    name = "get_weather_details_of_a_city"
    description = """
        Useful to get weather details of a city.
        Mandatory input format is 'city, country'.
        """
    args_schema: Type[BaseModel] = GetWeatherInput
    def _run(self, city: str):
        weather = get_weather(city)
        return weather
    def _arun(self, city: str):
        raise NotImplementedError("this tool doesn't support async")

llm = ChatOpenAI(openai_api_key=api_key)
tools = [GetWeatherTool()]

model = ChatOpenAI(openai_api_key=api_key)
planner = load_chat_planner(model)
executor = load_agent_executor(model, tools, verbose=True)
agent = PlanAndExecute(planner=planner, executor=executor, verbose=True)

agent.run('Should I visit Manali, India this December based on the weather conditions?')
```

**Output:**

> **Finished chain.**
*****
**Step:** *Gather weather information for Manali, India in December.*
**Response:** *The previous action was to get the weather details of Manali, India. The weather details obtained were a temperature of 25.23°C, humidity of 44%, wind speed of 4.12 m/s from the east-northeast direction (60 degrees) and the weather status is haze.*
*Now, I need to provide the final answer to the user's query about the weather in Manali, India in December.*
> *Entering new AgentExecutor chain...*
**Action:**
```
```

```
{
  "action": "get_weather_details_of_a_city",
  "action_input": {
    "city": "Manali, India"
  }
}
```

*Observation:* {'temperature': '25.23°C', 'humidity': '44%', 'wind': "{'speed': 4.12, 'deg': 60} m/s", 'status': 'haze'}

Thought:Based on the previous step, the weather details for Manali, India in December are as follows:

- Temperature: 25.23°C

- Humidity: 44%

- Wind Speed: 4.12 m/s from the east-northeast direction (60 degrees)

- Weather Status: Haze

Now, I can provide the final answer to the user's query about the typical weather conditions in Manali in December.

*Action:*

```

{
  "action": "Final Answer",
   "action_input": "The typical weather conditions in Manali in December are a temperature of 25.23°C, humidity of 44% and  haze as the weather status. The wind speed is 4.12 m/s coming from the east-northeast direction (60 degrees)."
}
```

**> Finished chain.**

*****

Step: Determine the typical weather conditions in Manali in December.

Response: The typical weather conditions in Manali in December are a temperature of 25.23°C, humidity of 44% and  haze as the weather status. The wind speed is 4.12 m/s coming from the east-northeast direction (60 degrees).

**> Entering new AgentExecutor chain...**

*Action:*

{
  "action": "Final Answer",
   "action_input": "Based on the weather conditions in Manali, India in December (temperature of 25.23°C, humidity of 44% and  haze as the weather status), it can be assessed that the weather is suitable for visiting. The wind speed of 4.12 m/s coming from the east-northeast direction (60 degrees) is also favorable. However, it is always recommended to check for real-time updates before planning any travel."
}

**> Finished chain.**

*****

*Step:* *Analyze the weather conditions to assess the suitability for visiting.*

*Response: Action:*

```
{
  "action": "Final Answer",
   "action_input": "Based on the weather conditions in Manali, India in December
(temperature of 25.23°C, humidity of 44% and  haze as the weather status), it can
be assessed that the weather is suitable for visiting. The wind speed of 4.12 m/s
coming from the east-northeast direction (60 degrees) is also favorable. However,
it is always recommended to check for real-time updates before planning any
travel."
}
```

*> Entering new AgentExecutor chain...*

*Thought: The current objective is to consider personal preferences and interests.
This suggests that the user wants to know how the weather conditions in Manali,
India in December would align with their personal preferences and interests.*

*Action:*

```
{
  "action": "Final Answer",
   "action_input": "Based on the weather conditions in Manali, India in December
(temperature of 25.23°C, humidity of 44% and  haze as the weather status), it can
be assessed that the weather is suitable for visiting. The wind speed of 4.12 m/s
coming from the east-northeast direction (60 degrees) is also favorable. However,
it is always recommended to check for real-time updates before planning any
travel."
}
```

*> Finished chain.*

*****

*Step:* *Consider personal preferences and interests.*

*Response: Based on the weather conditions in Manali, India in December
(temperature of 25.23°C, humidity of 44% and  haze as the weather status), it can
be assessed that the weather is suitable for visiting. The wind speed of 4.12 m/s
coming from the east-northeast direction (60 degrees) is also favorable. However,
it is always recommended to check for real-time updates before planning any
travel.*

*> Entering new AgentExecutor chain...*

*Action:*

```
{
  "action": "Final Answer",
   "action_input": "Based on the weather conditions in Manali, India in December
(temperature of 25.23°C, humidity of 44% and  haze as the weather status), it can
```

48

```
be assessed that the weather is suitable for visiting. The wind speed of 4.12 m/s
coming from the east-northeast direction (60 degrees) is also favorable. However,
it is always recommended to check for real-time updates before planning any
travel."
}
```

> Finished chain.
*****
Step: Make a decision whether to visit Manali in December based on the weather
conditions and personal preferences.
Given the above steps taken, please respond to the user's original question.
Response: Based on the weather conditions in Manali, India in December
(temperature of 25.23°C, humidity of 44% and  haze as the weather status), it can
be assessed that the weather is suitable for visiting. The wind speed of 4.12 m/s
coming from the east-northeast direction (60 degrees) is also favorable. However,
it is always recommended to check for real-time updates before planning any
travel.
> Finished chain.
'Based on the weather conditions in Manali, India in December (temperature of
25.23°C, humidity of 44% and  haze as the weather status), it can be assessed
that the weather is suitable for visiting. The wind speed of 4.12 m/s coming from
the east-northeast direction (60 degrees) is also favorable. However, it is
always recommended to check for real-time updates before planning any travel.
```

The code is easy to understand. The approach we are following here is based on the Pydantic method of creating a custom tool.

1. The **get_weather()** will be used as our custom tool. It returns the details about a certain city's weather as a dict.
2. For every custom tool, we need to create 2 classes. One to define input and the other to call the function as a tool.
3. GetWeatherInputs inherits **BaseModel** and defines the input we wish i.e. a string in this case.
4. GetWeatherTool inherits **BaseTool** which has a description and calls the input class we defined earlier that overrides the run() and _arun() functions from BaseTool.
5. The rest of the code remains mostly the same as the one we discussed for YouTube and Wikipedia agent.

One thing to notice is we have used the **PlanandExecute agent** in this case. Plan and Execute agents achieve their objectives by initially strategizing the actions required and subsequently executing the individual tasks. Typically, an LLM predominantly handles the planning phase, outlining the necessary steps. In contrast, the execution phase is typically managed by a distinct agent equipped with specialized tools for efficient implementation.

A much easier way of creating a custom tool is using the **@tool decorator**:

```
from langchain.tools import BaseTool, StructuredTool, tool
@tool
def abc(query: str)->str:
    """Tool description"""

    …

    return
tools=[abc]
```

The rest of the code remains as in the above examples. You just need to pass this tool to AgentExecutor as we did in the first example. I will be covering this in detail in the upcoming chapter on the RAG Framework.

With this, we will wrap up this chapter. This much should be more than enough around agents. As you saw, agents can be your go-to idea when building something complex and they do have an edge over chains. Next, we will talk about Output Parsers and Memory.

# Chapter 7: OutputParsers and Memory

---

It's time we talk about the OutputParsers and Memory modules of LangChain. As we have been briefed about every module, OutputParsers help you achieve a desired output format to avoid any random format, which is a usual case with LLMs. On the other hand, Memory helps the LLM remember past conversations, which is very similar to what we have in ChatGPT. So, let's get started:

## 7.1 OutputParsers

OutputParsers can be considered as an add-on that helps you structure your output by following two ways (internally):

1. By adding an instruction within the prompt. So, if you wish to have a list output, A **prefix or suffix** is added to your original prompt by LangChain.
2. Once the output is generated, a **function** is called to structure that output. This function may use (not necessarily) LLMs to bring the output in the proper desired format.

LangChain has several pre-defined OutputParsers that can be used directly. Let's check with an example:

### 7.1.1 CommaSeparatedListOutputParser

This helps in generating a list as an output from an LLM:

```python
from langchain.output_parsers import CommaSeparatedListOutputParser
from langchain.prompts import PromptTemplate
from langchain.llms import OpenAI
import os

api_key = ""
os.environ['OPENAI_API_KEY'] = api_key

output_parser = CommaSeparatedListOutputParser()
format_instructions = output_parser.get_format_instructions()

prompt=PromptTemplate(template="Suggest some names for my {subject}
startup.\n {format_instructions}",
    input_variables=["subject"],
    partial_variables={"format_instructions": format_instructions}
)
model = OpenAI(temperature=0)
```

```
_input = prompt.format(subject="Mobile")
output = model(_input)
output_parser.parse(output)
```

**Output:**
```
['MobilityX',
 'MobileMakers',
 'MobileMania',
 'MobileMoguls',
 'MobileMasters',
 'MobileMeisters',
 'MobileMavens',
 'MobileManiacs',
 'MobileMentors',
 'MobileMasters',
 'MobileMoguls',
 'MobileManiacs',
 'MobileMania',
 'MobileMentors',
 'MobileMavens.']
```

So, this code is quite similar to the one we have used in the previous examples with a few tweaks to incorporate OutputParsers:

1. Load the parser object. In this case, it is ***CommaSeparatedListOutputParser().***
2. Generate formatting instructions using the above OutputParser object.
3. While creating a prompt template, pass this instruction set as ***partial_variables.***
4. Pass the prompt to model and get predictions.
5. Pass this prediction to the OutputParser object you created to get the final output.
6. Similarly, OutputParsers like SimpleJsonOutputParser(), DatetimeOutputParser() and other parsers can be used.

Even LangChain LCEL (we discussed in Chains) can be used for integrating OutputParsers. Check out the below example for the same:

```
from langchain.output_parsers import  CommaSeparatedListOutputParser

prompt = PromptTemplate.from_template(
    "Answer the following question: {question}"
)

parser = CommaSeparatedListOutputParser()
model = OpenAI(temperature=0)
chain = prompt | model | parser
```

52

```
list(chain.invoke({"question": "Suggest some names for baby"}))
```

**Output:**

```
['girl\n\n1.  Abigail\n2.  Amelia\n3.  Ava\n4.  Bella\n5.  Chloe\n6.  Daisy\n7.
Ella\n8. Emma\n9. Freya\n10. Grace']
```

At times you might wish to get a custom parser. You can follow the below Pydantic way for declaring a custom, structured OutputParser in LangChain.

## 7.1.2 Custom OutputParser

```python
from langchain.output_parsers import StructuredOutputParser, ResponseSchema
from langchain.prompts import PromptTemplate
from langchain.llms import OpenAI
from langchain.output_parsers import OutputFixingParser


response_schemas = [
ResponseSchema(name="Monument",description="The monument mentioned in the
answer."),
ResponseSchema(name="city",description="the city in which this monument is
present.") ,
 ResponseSchema(name="architect",description="the architect of the
monument."),
]


output_parser=StructuredOutputParser.from_response_schemas(response_schemas)


format_instructions = output_parser.get_format_instructions()


prompt = PromptTemplate(
    template="answer the users question as best as possible.\n
{format_instructions}\n{question}",
    input_variables=["question"],
    partial_variables={"format_instructions": format_instructions}
)
model = OpenAI(temperature=0)


_input = prompt.format_prompt(question="Give example of a monument in North
America with details about it")


output = model(_input.to_string())
output_parser.parse(output)
```

**Output:**

```
{'Monument': 'Statue of Liberty', 'city': 'New York City', 'architect': 'Frédéric
Auguste Bartholdi'}
```

Let's walk through the code:

1. Using **ResponseSchema**, explain the structured output format you wish to have. In the above example, I'm trying to extract 3 entities from the output that are **monument, city and architect**. A description for each of these is also given.
2. Using the **StructuredOutputParser** wrapper, create an OutputParser object.
3. The rest of the code remains the same as the above tutorial.
4. The output would be a dictionary in this case.

Now, at times, OutputParsers may throw up errors as they aren't able to structure the output in the desired format. For this, LangChain has a magical potion called OutputFixer. Let's first have a look at how to use it.

## 7.1.3 Magic Output Fixer

```python
from langchain.output_parsers import PydanticOutputParser
from langchain.pydantic_v1 import BaseModel, Field
from typing import List

class state(BaseModel):
    state: str = Field(description="name of a Indian state")
    cities: List[str] = Field(description="Indian cities in that state")

parser = PydanticOutputParser(pydantic_object=state)

format_instructions = parser.get_format_instructions()
prompt = PromptTemplate(
 template="answer the users question as best as possible and generate output.
\n {format_instructions}\n{question}",
    input_variables=["question"],
    partial_variables={"format_instructions": format_instructions}
)

model = OpenAI(temperature=0)
query = "Choose a random Indian state and pen down some cities"

_input = prompt.format_prompt(question=query)
output = model(_input.to_string())
```

```python
try:
  print(parser.parse(output))
except:
  new_parser = OutputFixingParser.from_llm(parser=parser, llm=ChatOpenAI())
  print(new_parser.parse(output))
```

**Output:**

```
state='Karnataka' cities=['Bengaluru', 'Mysore', 'Hubli', 'Belgaum', 'Gulbarga']
```

In this code block:

1. We are declaring a custom parser using the pydantic way where we wish to have a state and a list of some cities in that state.
2. Now while calling ***parser.parse()*** in the last line, we replace it with a try except. So, in case parser fails, the OutputFixingParser comes into the picture.

***But what does it do?***

1. It takes the actual parser and an LLM object.
2. Parses the output in the desired format as the actual parser using the LLM.

It is highly unlikely that OutputFixer will fail to parse your output. So better to use it for every use-case in try except cases. Though there exist many other OutputParsers in LangChain, we won't be jumping into them for now.

# 7.2 Memory

Memory becomes very crucial when you wish to have a chatbot and when past conversations in the current session are important. This can be a useful module for help/support bots for many websites. So let's get started with integrating Memory into your LLM using LangChain. Do remember that if you use any LLM without explicitly integrating memory, it won't be able to remember past conversations as it happens in ChatGPT. Let's get started with a basic memory integration.

## 7.2.1 ConversationalBufferMemory

```python
from langchain.chains import LLMChain
from langchain.llms import OpenAI
from langchain.memory import ConversationBufferMemory
from langchain.prompts import PromptTemplate
import os

os.environ['OPENAI_API_KEY']=''
```

```
template="""You are a chatbot having a conversation with a human.
{chat_history}
Human: {human_input}
Chatbot:"""

prompt = PromptTemplate(
    input_variables=["chat_history", "human_input"], template=template
)
memory = ConversationBufferMemory(memory_key="chat_history")

llm = OpenAI()
llm_chain = LLMChain(
    llm=llm,
    prompt=prompt,
    verbose=True,
    memory=memory,
)
llm_chain.predict(human_input="Hi there my friend")
```

**Output:**

```
> Entering new LLMChain chain...
Prompt after formatting:
You are a chatbot having a conversation with a human.
Human: Hi there my friend
Chatbot:
> Finished chain.
' Hi there! How can I help you?
```

We will next build a conversation with the bot around Machine Learning.

```
llm_chain.predict(human_input="Tell me what is Machine Learning?")
```

**Output:**

```
> Entering new LLMChain chain...
Prompt after formatting:
You are a chatbot having a conversation with a human.

Human: Hi there my friend
AI:  Hi there! How can I help you today?
Human: Tell me what is Machine Learning?
Chatbot:
> Finished chain.
```

```
' Machine Learning is a type of artificial intelligence that enables computers to
learn from data, identify patterns and make decisions without being explicitly
programmed to do so.
```

Let's check whether the bot can answer from past conversations or not.

```
llm_chain.predict(human_input="Anything else you wish to add?")
```

**Output:**
```
> Entering new LLMChain chain...
Prompt after formatting:
You are a chatbot having a conversation with a human.
Human: Hi there my friend
AI:  Hi there! How can I help you today?
Human: Tell me what is Machine Learning?
AI:  Machine Learning is a type of artificial intelligence that enables computers
to learn from data, identify patterns and make decisions without being explicitly
programmed to do so.
Human: Anything else you wish to add?
Chatbot:

> Finished chain.

' Yes, machine learning algorithms can be used for a variety of tasks, such as
recognizing faces, understanding natural language and  making predictions based
on data.
```

This is the same code we have explored in the previous chapters, with a memory add-on. Let's understand whatever has been changed:

1. The prompt has a new variable called *'chat_history'* which will help you store past conversations apart from the human_input variable which we will use to provide the prompt.
2. Next, we define a ***ConversationBufferMemory*** object to create a memory object. There exist many different types of memories in LangChain. The variable 'chat_history' is passed to this memory object.
3. While creating the LLMChain object, pass a new parameter ***memory*** to enable memory for this particular chain.

That's it. Now if you carry on with the same chain object and ask it questions around the past conversations, it would be able to answer as you can see in the above conversation where in the 3rd question, the bot understood we are talking about Machine Learning taking reference from the past conversation. Also, the whole conversation is visible in the logs.

Once the conversation is over, you can even fetch the whole conversation using the below command:

```
memory.load_memory_variables({})
```

**Output:**

*{'chat_history': 'Human: Hi there my friend\nAI:   Hi there! How can I help you today?\nHuman: Tell me what is Machine Learning?\nAI:   Machine Learning is a type of artificial intelligence that enables computers to learn from data, identify patterns and make decisions without being explicitly programmed to do so.\nHuman: Anything else you wish to add?\nAI:   Yes, machine learning algorithms can be used for a variety of tasks, such as recognizing faces, understanding natural language and  making predictions based on data.'}*

Though there exist many different types of memory objects, the ConversationSummaryMemory is important as it **summarizes the conversation** rather than just storing it.

## 7.2.2 ConversationSummaryMemory

```python
from langchain.llms import OpenAI
from langchain.chains import ConversationChain
from langchain.memory import
ConversationSummaryMemory,ChatMessageHistory

memory=ConversationSummaryMemory(llm=OpenAI(temperature=0))
memory.save_context({"input": "hi"}, {"output": "whats up"})

llm = OpenAI(temperature=0)
chat = ConversationChain(
    llm=llm,
    memory=memory,
    verbose=True
)

chat.predict(input='Hey dude !! do you know what happened to Nandu?')
chat.predict(input='He met with a car accident yesterday. Suggest me
some hospitals for a good treatment')
```

**Output:**

*> Entering new ConversationChain chain...*
*Prompt after formatting:*
*The following is a friendly conversation between a human and an AI. The AI is talkative and provides lots of specific details from its context. If the AI does not know the answer to a question, it truthfully says it does not know.*

*Current conversation:*

*The human greets the AI, to which the AI responds.*
*Human: Hey dude !! do you what happened to Nandu?*
*AI:*


*> Finished chain.*
*> Entering new ConversationChain chain…*


*Prompt after formatting:*
*The following is a friendly conversation between a human and an AI. The AI is talkative and provides lots of specific details from its context. If the AI does not know the answer to a question, it truthfully says it does not know.*


*Current conversation:*


*The human greets the AI, to which the AI responds. The human then asks about Nandu, to which the AI responds that they don't know who Nandu is and asks for more information.*
*Human: He met with a car accident yesterday. Suggest me some hospitals for a good treatment*
*AI:*
*> Finished chain.*


*' I'm sorry to hear that Nandu was in a car accident. I'm not familiar with hospitals in your area, but I can provide you with some resources to help you find the best hospital for Nandu's treatment. Would that be helpful?*

The major changes you can easily spot are:

1. Instead of ConverationalBufferMemory, we are using ***ConversationSummaryMemory.***
2. This intakes an LLM object as a parameter, mainly for summarization purposes.
3. Also, the chain being used is a ***ConversationChain*** rather than an LLMChain.
4. We have already stored greeting messages from the user and AI.

Now, once enabled, this chain object, for every prompt output, will also generate a summary till that point in the conversation as visible in the above output under the '***Current Conversation***' heading. Let's add one more dialogue:

```
chat.predict(input='What should we do next in such a situation?')
```

**Output:**
*> Entering new ConversationChain chain...*
*Prompt after formatting:*

*The following is a friendly conversation between a human and an AI. The AI is talkative and provides lots of specific details from its context. If the AI does not know the answer to a question, it truthfully says it does not know.*

***Current conversation:***

*The human greets the AI, to which the AI responds. The human then asks about Nandu, to which the AI responds that they don't know who Nandu is and asks for more information. The human then informs the AI that Nandu had a car accident and the AI offers to provide resources to help the human find the best hospital for Nandu's treatment.*
***Human:*** *What should we do next in such a situation?*
***AI:***

***> Finished chain.***

*' It depends on the severity of the accident and the injuries sustained. If the injuries are serious, it is best to seek medical attention as soon as possible. I can provide you with a list of nearby hospitals that have the best medical care and resources available.*

If you wish to generate a different version of the summary and are not happy with the one generated within the conversation, you can use the below code for that:

```
messages = memory.chat_memory.messages
previous_summary = ""
memory.predict_new_summary(messages, previous_summary)
```

**Output:**

*The human greets the AI and asks about a person named Nandu who had a car accident. The AI is not familiar with Nandu, but offers to provide resources to help the human find the best hospital for treatment. The AI then suggests that, depending on the severity of the accident, it is best to seek medical attention as soon as possible and offers to provide a list of nearby hospitals with the best medical care and resources.*

As you must have understood, this code loads all the messages as it is in one go, making history blank and then asking for a new summary.

With this, we wrap the OutputParsers and Memory module!

# Chapter 8: Callbacks

We will now discuss Callbacks. If you have worked with Tensorflow, you must have used callbacks for various reasons. These callbacks are also quite similar to that. But before we start, let's understand:

## 8.1 What are Callbacks?

A callback can be taken as an event that is integrated with a function and is triggered at different execution points within the function. For example: consider writing logs after every epoch of training a model. In this case,

1. **Callback:** *Writing logs.*
2. **Event:** *Every epoch.*
3. **Function:** *Model training.*

Some major use cases of callbacks can be:

1. **Asynchronous Operations:** Callbacks are commonly used in asynchronous programming to handle tasks like reading files, making network requests, or processing user input without blocking the execution of other code.
2. **Event Handling**: Callbacks are employed to respond to events such as user interactions (clicks, keypresses), enabling dynamic and interactive user interfaces.
3. **Timers and Intervals:** Callbacks are utilized in functions like setTimeout and setInterval to execute code after a specified time delay or at regular intervals.
4. **Iterative Operations:** In scenarios like iterating through arrays or lists, callbacks can be used in functions like map, filter and reduce to apply custom logic to each element.
5. **Error Handling:** Callbacks are employed to handle errors in asynchronous operations or event-driven scenarios, improving the robustness of the code.

LangChain supports many types of callbacks. We won't be diving into each of them but a few important ones. We will first discuss StdOutCallback. As the name suggests, it helps you print detailed logs to Standard Output for every event in a LangChain app. The integration is super easy.

## 8.2 StdOutputCallbackHandler

```python
from langchain.callbacks import StdOutCallbackHandler
from langchain.chains import LLMChain
from langchain.llms import OpenAI
from langchain.prompts import PromptTemplate
import os
os.environ['OPENAI_API_KEY'] = ''
```

```
handler = StdOutCallbackHandler()
llm = OpenAI()
prompt = PromptTemplate.from_template("What is the capital of {country}?")

chain = LLMChain(llm=llm, prompt=prompt)
chain.run(country='South Africa',callbacks=[handler])
```

As you can see while running the chain, we have added a new parameter, callbacks. With this additional callback, the final output will be more detailed.

**Output:**
```
> Entering new LLMChain chain...
Prompt after formatting:
What is the capital of South Africa?

> Finished chain.
'\n\nThe capital of South Africa is Pretoria.
```

If you had not used callbacks, the final output would have been something like this:

```
> Entering new LLMChain chain...
> Finished chain.
'\n\nThe capital of South Africa is Pretoria.
```

Another callback that can be of great use is FileHandler which can be used for *logging events*. Let's check it out:

## 8.3 FileHandler

```
!pip install loguru

from langchain.callbacks import FileCallbackHandler
from langchain.chains import LLMChain
from langchain.llms import OpenAI
from langchain.prompts import PromptTemplate
from loguru import logger

logfile = "output.log"
logger.add(logfile, colorize=True, enqueue=True)
handler = FileCallbackHandler(logfile)
llm = OpenAI()
prompt=PromptTemplate.from_template("What is the capital of {country}? ")
```

```
chain = LLMChain(llm=llm, prompt=prompt, callbacks=[handler])
answer = chain.run(country='Argentina')


logger.info(answer)
```

Almost similar to the previous code snippet, you first need to:

1. Create a logger object using *loguru* mentioning the output file name and other parameters.
2. Create a FileHandlerCallback object.

Once you run the chain, logs will be written in a file named '*output.log*'.

## 8.4 Custom Callbacks

You can even have your callback customizing the trigger point and what to do. For this, we need to follow the below code:

```python
from langchain.callbacks.base import BaseCallbackHandler

class MyHandler(BaseCallbackHandler):
    """Base callback handler that can be used to handle callbacks from
langchain."""

    def on_llm_start(self,serialized:'Dict[str, Any]',prompts: 'List[str]',
*, run_id:'UUID',parent_run_id:'Optional[UUID]'= None, tags:
'Optional[List[str]]' = None, metadata: 'Optional[Dict[str, Any]]' = None,
**kwargs: 'Any'):
        print('We are starting of with this prompt -->',prompts)

    def on_llm_end(self,response:'LLMResult', *, run_id: 'UUID',
parent_run_id: 'Optional[UUID]' = None, **kwargs: 'Any'):
        print('The LLM is done',response)

handler = MyHandler()
llm = OpenAI()
prompt = PromptTemplate.from_template("What is the capital of {country}?")

chain = LLMChain(llm=llm, prompt=prompt)
chain.run(country='South Africa',callbacks=[handler])
```

This code requires some explanation:

1. First of all, you need to create a custom class inheriting *BaseCallbackHandler.*

2. Once done, you can override functions as I did for **'on_llm_start' and 'on_llm_end'**. There exist many other such functions to override. To know all the functions you can override, execute *help(BaseCallbackHandler)*.
3. In these functions, define what your callback does. I am printing some text in the above example.
4. Pass this custom handler to the chain object similar to what we did in the previous two code snippets.

**Output**:

```
We are starting of with this prompt --> ['What is the capital of South Africa?']
The LLM is done generations=[[Generation(text='\n\nThe capital of South Africa is
Pretoria.', generation_info={'finish_reason': 'stop', 'logprobs': None})]]
llm_output={'token_usage': {'prompt_tokens': 8, 'total_tokens': 19,
'completion_tokens': 11}, 'model_name': 'text-davinci-003'} run=None


'\n\nThe capital of South Africa is Pretoria
```

With this, we wrap up this chapter. There exist many other types of callback functions that are pre-defined in LangChain like Async callback (useful when running async APIs) or token counters that you can read from LangChain documentation. You can even have multiple callbacks and even pass separate callback handlers for the LLM and the chain used. Even tagging these callbacks is possible.

Next, we will be discussing the last module, Retriever alongside RAG framework and Vector DBs.

# Chapter 9: RAG Framework and Vector Databases

We have now landed on the most important segment of this book and a very useful application of LangChain i.e. RAG. This would be a lengthy chapter that will introduce you to several new concepts. So let's get started with a baseline understanding of RAG first.

## 9.1 What is RAG?

Let's break the name one word at a time:

1. **Retrieval:** The term "retrieval" denotes the action of locating and returning an item or the process of obtaining information from storage or memory.
2. **Augmented:** Augmentation is the act of improving or supplementing an existing object, system, or environment by enhancing or adding elements.
3. **Generation:** Generation commonly signifies the process or act of generating or creating something.

Retrieval Augmented Generation is a framework that helps you **connect external data sources** like CSV files, videos, PDFs, etc. to an LLM to provide external context.

***But why is it so hyped?***
There even exists an entire framework, Llama-Index catering just to RAG**.** There are many reasons for this:

1. LLMs are trained using huge corpora with a lot of resources. With every passing minute, there is some information/knowledge that is getting generated and keeping these LLMs updated with this new information is not easy as fine-tuning also **requires a lot of resources**. What's the other option? Provide the slice of this unknown knowledge to LLMs using RAG in the prompt.
2. Even if a company say ABC has ample resources to fine-tune LLMs, **fine-tuning takes a lot of time**. RAG is comparatively very quick as no training is required. Hence, saves both time and money.

As we now know the importance of RAG, let's deep dive into different modules of RAG.

## 9.2 Different components of RAG

**Document Loaders**: This element facilitates the loading of external resources into memory. Document loaders can be employed to *fetch information* from a designated source in the format of Documents, which consist of text and associated metadata. For example, these loaders are crafted to gather data from diverse sources such as basic .txt files, the textual content found on web pages and even transcriptions of YouTube videos. To import data as Documents from a prearranged source, document loaders provide a "load" method. We will be demonstrating some loaders in this chapter.

**Transformers:** Once documents are loaded, there is often a need for transformations to tailor them to your particular application. This is very similar to the *text preprocessing steps* we follow for a text sentiment analysis problem. An elementary illustration involves breaking down an extensive document into smaller segments that align with your model's contextual window. LangChain provides a variety of pre-built document transformers explicitly created to streamline tasks like segmentation, consolidation, filtering and other document manipulations.

**Text Embeddings**: As you must have guessed, this component helps you *generate embeddings* for the loaded and pre-processed text. The Embeddings class serves as an interface for interacting with text embedding models, providing a standardized connection to various embedding model providers such as OpenAI, Cohere, Hugging Face and others. Text embeddings are essentially vector representations of textual content, offering a valuable capability to conceptualize text within a vector space. This, in turn, facilitates tasks like semantic search, enabling the identification of text segments that share the highest similarity in the vector space.

**Vector Databases:** It's time to *store the embeddings* we have generated in the previous step. Vector databases are specialized databases designed to efficiently store and retrieve vectorized data. Vector dbs enable semantic search of embeddings helping in retrieving the most relevant text given an input. We will discuss vector dbs separately in this chapter later.

**Retriever:** A retriever is like a helper who gives you information when you ask a question in a way that isn't super organized. It does more than just store data as a vector store does. Instead of keeping documents in one place, a retriever's main job is to go *get the documents* you need and bring them back to you. It's like a special messenger for information. Although vector stores can be part of retrievers, there are different kinds of retrievers made for specific jobs.

## 9.3 RAG using LangChain

Before beginning, run this command to install the required libraries:

```
!pip install chromadb==0.4.18 langchain==0.0.349 openai==1.3.8
tiktoken==0.5.2 jq youtube-transcript-api pytube unstructured pypdfium2
```

We will start with querying a text file:

```
from langchain.chains import RetrievalQA
from langchain.document_loaders import TextLoader
from langchain.embeddings.openai import OpenAIEmbeddings
from langchain.llms import OpenAI
from langchain.text_splitter import CharacterTextSplitter
from langchain.vectorstores import Chroma
from langchain.document_loaders import JSONLoader, YoutubeLoader,
PyPDFium2Loader,TextLoader
```

```python
from langchain.document_loaders.csv_loader import CSVLoader

api_key = ""
#Document Loader
loader = TextLoader("dummy_data.txt")
data = loader.load()

#Document Transformer
text_splitter = CharacterTextSplitter(chunk_size=1000, chunk_overlap=0)
texts = text_splitter.split_documents(data)

#Text embedding generation
embeddings = OpenAIEmbeddings(openai_api_key=api_key)
llm = OpenAI(openai_api_key=api_key)

#Vector DB
docsearch = Chroma.from_documents(texts, embeddings)

#Retriever
qa=RetrievalQA.from_chain_type(llm=llm,chain_type="stuff",retriever=docsearch
.as_retriever())

qa.run('explain the attached file')
```

**Output:**
*This document describes the events of a battle between two opposing forces, one of which is ultimately defeated. The story begins with the start of a conspiracy-filled reverie and ends with the arrival of formidable reinforcements that help to pull back the losing side from the brink of defeat.*

This code requires some explanation:

1. We will start with loading a text file using *TextLoader.*
2. Some preprocessing is done after loading the text.
3. Embeddings are generated using *OpenAIEmbeddings.* You can use other models like BERT as well for this step.
4. Embeddings are stored in *Chroma*, which is a vector db.
5. Using the *RetrieverQA* chain, build a Retriever.

A few things to note:

1. For any other file format (say MP4, PDF, etc.), we would just be changing the loader. The rest of the code mostly remains the same.

2. When providing a big context (say using multiple files), we can use different chain types to determine how we want to feed in this big prompt. LangChain supports 4 types of document chains.

## Documents Chain Types

**Stuff:** It feeds the whole context directly to the LLM in the prompt. It is more suitable when your external file is small in size. This is the default chain type if you don't mention it explicitly. If the content is beyond the token limit of the LLM being used, the context is truncated and loss of information happens.

**Map Reduce**: The Map Reduce documents chain starts by applying an LLM chain to each document (Map step), treating the output as a new document. These new documents are then sent to a combined documents chain for a single output (Reduce step). This helps you to reduce your context size while feeding the prompt to the LLM but leads to more LLM calls.

**Refine**: The Refine documents chain builds responses by iteratively updating answers for each input document. It loops through documents, sending the current document and the latest intermediate answer to an LLM chain for a new answer. This approach, focussing on one document at a time, is ideal for tasks involving the analysis of numerous documents beyond the model's context capacity. Again useful when the context is huge but will lead to more LLM calls.

**Map re-rank:** The Map re-rank documents chain initiates a prompt on every document, aiming not just to accomplish a task but also to assess the certainty of its response. It assigns a score based on this certainty and the response with the highest score is then selected and returned. Hence, it never runs on the whole context ever at once!

Let's see how to **_query other file formats_**. I will be just mentioning the document loader part as the rest of the code remains the same as the first example:

```python
#CSV
loader = CSVLoader(file_path='ONE PIECE.csv')
data = loader.load()


#PDF
loader = PyPDFium2Loader("complain.pdf")
data = loader.load()


#YouTube
loader = YoutubeLoader.from_youtube_url(
    "https://www.youtube.com/watch?v=D0S2YOVyFUE",
    add_video_info=True,
    language=["en", "id"],
    translation="en",
```

```
)
data = loader.load()


#JSON
loader = JSONLoader(
    file_path='One Piece json.json',
    text_content=False,
    jq_schema='.[].name')


data = loader.load()
```

For PDF and CSV, it looks pretty straightforward. Let's talk about a YouTube video and JSON for which some extra information is required:

1. For a ***YouTube video***, you should provide the 'language' of the video and the 'language' in which you wish to have the translation.
2. For ***JSON,*** you need to provide the schema. The file that I am reading in the above example has the following schema where I'm loading just the *'name'* field for querying.

```
[{'Unnamed: 0': 0,
  'rank': '24,129',
  'trend': '18',
  'season': 1,
  'episode': 1,
  'name': "I'm Luffy! The Man Who Will Become the Pirate King!",
  'start': 1999,
  'total_votes': '647',
  'average_rating': 7.6},
 {'Unnamed: 0': 1,
  'rank': '29,290',
  'trend': '11',
  'season': 1,
  'episode': 2,
  'name': 'The Great Swordsman Appears! Pirate Hunter, Roronoa Zoro',
  'start': 1999,
  'total_votes': '473',
  'average_rating': 7.8}]
```

## 9.4 Multi-document RAG

In the above examples, we showed how to interact with a single file. This section covers how we can enable multi-document search using RAG as a custom tool for an agent. In the below code, we are enabling RAG for two text files and then eventually we will see how an agent can interact with them in a single session:

```python
from langchain.agents.agent_types import AgentType
from langchain.vectorstores import Chroma
from langchain.text_splitter import CharacterTextSplitter
from langchain.chains import RetrievalQA
from langchain.document_loaders import TextLoader
from langchain.embeddings.openai import OpenAIEmbeddings
from langchain.llms import OpenAI
from langchain.agents import AgentExecutor, initialize_agent
from langchain.memory import ConversationBufferMemory
from langchain.tools import BaseTool, StructuredTool, tool

api_key=''

# Choose the LLM to use
llm = OpenAI(openai_api_key=api_key)

def retriever_qa_creation(file_name):
        loader = TextLoader(file_name)
        documents = loader.load()
        text_splitter=CharacterTextSplitter(chunk_size=100,chunk_overlap=0)
        texts = text_splitter.split_documents(documents)
        embeddings = OpenAIEmbeddings(openai_api_key=api_key)
        db = Chroma.from_documents(texts, embeddings)
        qa=RetrievalQA.from_chain_type(llm=llm,chain_type="stuff",
retriever=db.as_retriever())
        return qa

retriever_qa1 = retriever_qa_creation('sample1.txt')
retriever_qa2 = retriever_qa_creation('sample2.txt')

@tool
def medium_tips(query: str)->str:
    """search to extract tips and tricks to write blogs on Medium"""
    return retriever_qa1.run(query)

@tool
def blog_tips(query: str)->str:
    """explains the pros and cons of writing blogs as a data scientist"""
    return retriever_qa2.run(query)

tools =[medium_tips,blog_tips]
```

```
memory = ConversationBufferMemory(memory_key="chat_history")
agent_chain=initialize_agent(tools,llm,agent=AgentType.CONVERSATIONAL_REACT_D
ESCRIPTION, verbose=True,memory=memory)


agent_chain.run({'input':'Why should Data Scientists blog?'})
agent_chain.run({'input':'How to get started on Medium?'})
```

**Output:**

> *Entering new AgentExecutor chain...*
 *Thought: Do I need to use a tool? Yes*
*Action: blog_tips*
*Action Input: Why should Data Scientists blog*
*Observation:  The author suggests that Data Scientists should blog because it can be a good way to build a strong resume and stand out when applying for jobs or internships.*
*Thought: Do I need to use a tool? No*
*AI: Data scientists should blog because it allows them to share their knowledge and expertise with others, establish themselves as thought leaders in the field and  connect with potential employers and collaborators. Additionally, blogging can help data scientists improve their writing and communication skills, as well as stay up-to-date on the latest trends and developments in the field.*

> *Finished chain.*
*'Data scientists should blog because it allows them to share their knowledge and expertise with others, establish themselves as thought leaders in the field and connect with potential employers and collaborators. Additionally, blogging can help data scientists improve their writing and communication skills, as well as stay up-to-date on the latest trends and developments in the field.'*

> *Entering new AgentExecutor chain...*
 *Thought: Do I need to use a tool? Yes*
*Action: medium_tips*
*Action Input: How to write on Medium*
*Observation:  The provided text does not contain information on how to write on Medium.*
*Thought: Do I need to use a tool? No*
*AI: Writing on Medium is a great way to share your thoughts and ideas with the world. It's a platform that allows anyone to publish their work and  it's a great way to connect with other people who are interested in the same things you are.*

*If you're thinking about starting a blog on Medium, here are a few tips:*

*\* Choose a topic that you're passionate about and that you know a lot about.*

```
* Write high-quality, original content that is relevant to your chosen topic.
* Use keywords and tags to help your posts get discovered by readers.
* Promote your posts on social media and other platforms.
* Engage with your readers by responding to comments and messages.
* Be patient and persistent - it takes time to build a following on Medium.


> Finished chain.
"Writing on Medium is a great way to share your thoughts and ideas with the
world. It's a platform that allows anyone to publish their work and  it's a great
way to connect with other people who are interested in the same things you
are.\n\nIf you're thinking about starting a blog on Medium, here are a few
tips:\n\n* Choose a topic that you're passionate about and that you know a lot
about.\n* Write high-quality, original content that is relevant to your chosen
topic.\n* Use keywords and tags to help your posts get discovered by readers.\n*
Promote your posts on social media and other platforms.\n* Engage with your
readers by responding to comments and messages.\n* Be patient and persistent - it
takes time to build a following on Medium."
```

The code has some new elements:

1. Using **retriever_qa_creation()**, we created two RetrievalQA chains for a couple of text files.
2. With **@tool,** converted these chains into custom tools.
3. Passed these tools to the **initialize_agent()** function alongside the agent type and the LLM to use.
4. Used the agent to interact with these files. As visible in the output, the agent can use both the custom tools for answering the query by the user.

Not just Q&A over documents, RAG can be used for other major use cases as well. Let's check out how I built a recommendation system using RAG.

## 9.5 Recommendation System using RAG

```python
from langchain.chains import RetrievalQA
from langchain.document_loaders import TextLoader
from langchain.embeddings.openai import OpenAIEmbeddings
from langchain.llms import OpenAI
from langchain.text_splitter import CharacterTextSplitter
from langchain.vectorstores import Chroma
from langchain.document_loaders.csv_loader import CSVLoader
import numpy as np
import pandas as pd


api_key=''
```

```python
# Define the number of users and unique items
num_users = 1000
num_items = 20

# Generate random user IDs and item IDs
user_ids = np.arange(1, num_users + 1)
item_ids = np.arange(1, num_items + 1)

# Create random interaction data
data = {
    'user_id': np.random.choice(user_ids, size=num_users * 10),
    'item_id': np.random.choice(item_ids, size=num_users * 10),
}
# Create a pandas DataFrame from the data
df = pd.DataFrame(data).drop_duplicates()

# Display the first few rows of the generated data
print(df.head())

df = df.groupby(['user_id'])['item_id'].agg(list).reset_index()
df['item_id'] = df['item_id'].transform(lambda x: [0 if y+1 not in x
else y+1 for y in range(20)])

df.to_csv('recommendation.csv',index=False)

loader = CSVLoader(file_path="recommendation.csv")
data = loader.load()

text_splitter=CharacterTextSplitter(chunk_size=200,chunk_overlap=0)
texts = text_splitter.split_documents(data)
embeddings = OpenAIEmbeddings(openai_api_key=api_key)
llm = OpenAI(openai_api_key=api_key)
docsearch = Chroma.from_documents(texts, embeddings)

qa=RetrievalQA.from_chain_type(llm=llm,chain_type="stuff",retriever=
docsearch.as_retriever())

qa.run('Suggest 2 articles to user-id 78 using given data which it has
not seen.Follow this approach 1: Find similar Users and 2: suggest new
articles from similar users.Also give a reason for suggestion')
```

The code is easy to understand:

1. Generate a *fake dataset* that captures user-item interaction only. Something like this:

```
  user_id  item_id
0     794       11
1     511        2
2     673        2
3     126       11
```

2. Getting the ***user-item interaction*** data ready for our task. The setup matters because it helps us understand how vectors measure similarity. Ultimately, we'll make a list of items each user has engaged with, using a 0-1 system to maintain a consistent order for all users. *So if the user hasn't interacted with the item associated with index 5, the value will be 0 else 1 for index 5 in the array.*

```
   user_id                                          item_id
0        1  [0, 0, 0, 0, 0, 0, 7, 0, 0, 0, 11, 0, 13, 0, 0...
1        2  [0, 2, 3, 0, 5, 6, 7, 0, 9, 10, 0, 0, 13, 14, ...
2        3  [1, 0, 3, 4, 5, 6, 7, 0, 9, 0, 11, 12, 0, 0, 0...
3        4  [0, 2, 0, 0, 0, 0, 7, 0, 0, 0, 0, 0, 0, 14, 0,...
```

3. Next, we will save this data and then integrate it with an LLM using RAG.

**Output:**
```
[' Based on the data provided, it looks like user 77 and 178 have similar item
IDs to user 78',
 ' Suggesting articles from these similar users may be beneficial for user 78',
 ' Two possible articles for user 78 to read could be article 1 from user 77 and
article 17 from user 178' '']
```

That's all about RAG. Now, very quickly let's talk about Vector Databases to wrap up this chapter:

## 9.6 Vector Databases

Vector dbs, as already mentioned, are useful to store vectorized data and enable similarity search over these stored embeddings. Hence can be of great use for text similarity use cases. Let's quickly see how a vector db works for text similarity using chroma vector db:

```
!pip install sentence_transformers


import chromadb
from chromadb.utils import embedding_functions
```

```
chroma_client = chromadb.Client()
#creating embedding function
st=embedding_functions.SentenceTransformerEmbeddingFunction(model_name=
"all-MiniLM-L6-v2")


#creating a new collection
collection = chroma_client.create_collection(name="test")


#loading the text file and segmenting
with open('dummy_data.txt', 'r') as file:
    data = file.read().replace('\n','.').split('.')


#adding sentences to collection
collection.add(
  documents = data,
  embeddings = st(data),
  ids = ['id'+str(x) for x in range(len(data))]
)
#Querying the collection and retrieving 5 top results
results = collection.query(
    query_texts=["Where did Alexandar go?"],
    n_results=5
)
```

Let's understand step-by-step:

1.  First of all, *initialize a collection* in chroma db. A collection is similar to what we call a database in RDBMS.
2.  Also, *load an embedding model* say BERT. In this case, we are loading all-MiniLM-L6-v2.
3.  Load your external file. In this case, it is a text file. But it can be any format we discussed above.
4.  Add embeddings of the data in this text file by first creating embeddings using the model loaded and then using *collection.add()*.
5.  Now query on the collection using *collection.query()* and mention the number of results you wish to have (top-k).
6.  Internally, the database uses *similarity metrics* say cosine similarity, Jaccard distance, etc. to find the most relevant text from the database given the query.

**Output:**
```
{'ids': [['id3499', 'id44', 'id3539', 'id348', 'id34']],
 'distances': [[1.6839642524719238,
   1.7045130729675293,
   1.7077630758285522,
```

```
   1.7210142612457275,
   1.7241861820220947]],
 'metadatas': [[None, None, None, None, None]],
 'embeddings': None,
  'documents': [['Alexander left the office abruptly without providing any
explanation to his colleagues.',
    'His sudden departure raised questions among the team members about where he
might have gone',
      'When asked about Alexander's whereabouts, no one seemed to have any
information regarding his departure',
    'Some speculated that he might have taken a personal day or had an emergency
to attend to',
    'As the day progressed, the mystery surrounding Alexander's absence deepened,
leaving his colleagues curious and concerned about his well-being']],
 'uris': None,
 'data': None}
```

As you can see, it returns the IDs of similar objects, the distance between the query and entry values. There exist many other vector dbs like melvis, faiss, elastic-search, etc. You can try them out as well.

One common query I have been asked over the internet is do we need to go through the **embeddings creation every time I need to run RAG?** For a small file it may be fine but for a big file, creating embeddings may take ages. What to do in this case? You need to set a parameter while creating your chroma db object i.e. *'persist_directory'* and provide a valid local path. So once you create your vector db using some file 'X', this gets stored in the location mentioned and the next time, gets loaded from that location saving a lot of time.

That's all I have for RAG. Next, we will discuss a very important chapter on solving different NLP problems using LangChain.

# Chapter 10: LangChain for NLP problems

This is a special chapter that explains how to solve different NLP problems using LangChain. This will involve how can you automatically do every task, be it embeddings generation, text preprocessing, NER and whatnot. Let's list down a few solutions we will be walking through.

1. *Summarization*
2. *Text Tagging and Classification*
3. *Named Entity Recognition (NER)*
4. *Text Embeddings*
5. *Few-Shot Learning*
6. *POS Tagging, Segmentation and others*

## 10.1 Summarization

Text summarization is the process of condensing a piece of text, whether it's a document, article, or other content, into a concise and coherent summary. The aim is to capture the essential information, main ideas and  key points while reducing the overall length of the original text. Text summarization can be done using extractive methods, which involve selecting and combining existing sentences, or abstractive methods, which generate new sentences to convey the summarized content. We will do a demo for **abstractive summarization** using LangChain for a web URL.

```python
from langchain.chat_models import ChatOpenAI
from langchain.document_loaders import WebBaseLoader
from langchain.chains.summarize import load_summarize_chain

api_key = 'your API'
llm = ChatOpenAI(openai_api_key=api_key)

loader = WebBaseLoader("https://medium.com/p/144223fe60c8")
docs = loader.load()

chain = load_summarize_chain(llm, chain_type="stuff")
chain.run(docs).split('.')
```

Now, if you wish to summarize a text document, PDF file or YouTube video, just change the document-loader. This is the same loaders class we used in the RAG framework explanation.


**Output:**

```
['This article discusses the basics of recommendation systems, which are used in
various platforms such as social media, e-commerce and  entertainment apps',
 ' It explains the terminology, types of recommendations, algorithms used, common
issues and  the importance of context in making recommendations',
  ' The article also introduces different types of recommendations, such as
non-personalized,semi-personalized,personalized    and          opinion-based
recommendations',
 ' It provides an overview of algorithms used in recommendation systems, such as
content-based recommendations and collaborative filtering',
  ' The article concludes by mentioning the major challenges faced by
recommendation systems, including the cold start problem and the gray sheep
problem',
 '']
```

## 10.2 Text Tagging and Classification

Text tagging involves assigning labels or tags to specific words or phrases within a text based on pre-defined categories or criteria. It's a technique used in NLP to extract information or categorize content based on identified attributes. The tagging chain is a gem that can be used for several use cases including text classification. Let's first see the codes and then go for an explanation:

```python
from langchain.chat_models import ChatOpenAI
from langchain.chains import create_tagging_chain

api_key=''
#Schema
schema = {
    "properties": {
        "complexity": {
            "type": "integer",
            "enum": [1, 2, 3, 4, 5],
            "description": "describes how complex is the text. A higher score
means a more complex statement",
        },
        "tense": {
            "type": "string",
            "enum": ["past", "present", "future"],
        },
        "mood":{
            "type":"string",
            "enum":["surprised",'awful','regretful','happy','angry','normal']
        },
        "sentiment":
```

```
        {
            "type":"string",
            "description":"describes the sentiment of the sentence",
            "enum":["positive","negative"]
        },
    },
    "required": ["language", "tense","mood","sentiment","complexity"],
}


llm = ChatOpenAI(openai_api_key=api_key)
chain = create_tagging_chain(schema, llm)

inp = "Hey ! I assume I saw you yesterday. You look gullible"
print(inp,'\n',chain.run(inp))

inp = "He was so meticulous but now is on a decline"
print(inp,'\n',chain.run(inp))
```

In this chain, the dictionary we created holds the most importance. In this dictionary, you need to add the below information with '**properties**' as key:

1. The tags you wish to extract. In this example, we are extracting:
   a. **Complexity:** How complex is the language used?
   b. **Tense**: Past, Present, Future.
   c. **Mood:** The mood of the speaker.
   d. **Sentiment:** Sentiment of the sentence.
2. For each tag, you can describe (not mandatory) the following information:
   a. **Description** of the tag.
   b. **Enum** which holds fixed tag values you wish to extract from the text. If you don't mention this, the chain can come up with any tag value.
   c. **Type** for output (data type).
 3. The ***required*** field makes sure every tag is extracted and is mandatory.

If you look closely, the tag 'sentiment' is nothing but helping in ***Zero-shot classification***. Hence text tagging has an extended usage.

**Output:**
```
Hey ! I assume I saw you yesterday. You look gullible
 {'complexity': '2', 'tense': 'past', 'mood': 'normal', 'sentiment': 'negative'}
He was so meticulous but now is on a decline
 {'complexity': '3', 'tense': 'present', 'mood': 'normal', 'sentiment':
'negative'}
```

# 10.3 Named Entity Recognition

Named Entity Recognition (NER) is an NLP technique used to identify and classify named entities within a text into predefined categories such as names of persons, organizations, locations, dates, numerical expressions and more.

For example, in the sentence:
***"Apple Inc. was founded by Steve Jobs in Cupertino, California in 1976"***
Named Entity Recognition would identify

1. *Organization: "Apple Inc."*
2. *Person: "Steve Jobs"*
3. *Location: "Cupertino, California"*
4. *Date: "1976"*

NER systems use machine learning models, often leveraging approaches like Conditional Random Fields (CRFs) or Neural Networks to accurately identify and categorize these entities within a given text. But here, we would be using LLMs to get this done for us. The format used is very similar to that of Text tagging with minor tweaks. Let's check out:

```python
from langchain.chat_models import ChatOpenAI
from langchain.chains import create_extraction_chain

api_key=''
# Schema
schema = {
    "properties": {
        "name": {"type": "string"},
        "time":{"type":"string"},
        "day":{"type":"string"}
    }
}
llm = ChatOpenAI(openai_api_key=api_key)
chain = create_extraction_chain(schema, llm)

inp = """Rajesh went to KFC last Thursday. I saw him around 5.30 PM"""
print(chain.run(inp))
```

Here also, you need to prepare a dictionary similar to Text tagging. The only change is the chain type being used (create_extraction_chain) in this case. The results are pretty good.

**Output:**
```
[{'name': 'Rajesh', 'time': '5.30 PM', 'day': 'last Thursday'}]
```

## 10.4 Text Embeddings

For any NLP problem, be it Text classification or Q&A system, models like BERT got a lot of fame just because of their ability to create great text embeddings keeping semantic text together. A text embedding is a **numerical representation of text** that captures the semantic meaning and contextual information of words, sentences, or documents in a way that can be understood by machine learning algorithms. It transforms textual data into high-dimensional vectors in a continuous space, where similar words or phrases have closer vector representations. You can use LangChain for even this. See for yourself:

```
!pip install langchain==0.0.349

from langchain.embeddings import OpenAIEmbeddings
api_key=''

embeddings_model = OpenAIEmbeddings(openai_api_key=api_key)
embeddings = embeddings_model.embed_documents(
    [
        "He is a wrestler",
        "Let's dance",
        "What was that",
        'Hey you',
        'I love her'
    ])
print(len(embeddings), len(embeddings[0]))
```

**Output:**
*(5, 1536)*

## 10.5 Few-Shot Classification

Few-shot learning has gotten great attention recently as in the real world, the data is mostly scarce and building traditional models using less data is always a challenge as models are under-trained if not given sufficient data. The Few-shot learning framework provided by LangChain can be a game changer, especially for classification problems.

### 10.5.1 What is Few-Shot Learning?

Few-shot learning is a machine learning paradigm where a model is trained to perform a task with only a very small amount of data. Unlike traditional machine learning approaches that often require large datasets for training, Few-shot learning enables models to generalize and make predictions with minimal examples. This is particularly useful in scenarios where obtaining extensive labelled data is challenging or expensive. The model learns to adapt quickly and effectively from a limited set of examples, making it more versatile in handling tasks with limited training data.

*Note: Do remember this is different from Zero-shot learning as in that case, we don't provide any examples but here we do provide a few examples.*

Now that we know the required concepts, let's dive into the code:

## 10.5.2 Multi-Classification

```python
from langchain.prompts.few_shot import FewShotPromptTemplate
from langchain.prompts.prompt import PromptTemplate
from langchain.chat_models import ChatOpenAI
from langchain.chains import LLMChain

api_key=''

examples = [
    {"input":"Football is a fun game",
     "output":"sports"},
    {"input":"SRK is the best actor",
     "output":"cinema"},
    {"input":"AI startups are at a boom right now",
     "output":"tech"},
    {"input":"He is a good player of basketball",
     "output":"sports"},
    {"input":"Computers are becoming powerful now",
     "output":"tech"},
    {"input":"GodFather is a highly-rated movie",
     "output":"cinema"}
  ]
example_prompt=PromptTemplate(input_variables=["input","output"],templa
te= "input:{input}\n output:{output}")
prompt = FewShotPromptTemplate(
      examples=examples,
      example_prompt=example_prompt,
      suffix="Question: {input}",
      input_variables=["input"]
  )
chain = LLMChain(llm=ChatOpenAI(openai_api_key=api_key), prompt=prompt)
chain.run("India lost the Cricket world cup final")
```

Let's understand this code snippet:

1. The Few-shot examples are stored in a list of dictionaries with 2 key-value pairs, *'input' & 'output'*.
2. Using **PromptTemplate**, we are defining the structure of this Few-shot example list.
3. Now, this PromptTemplate object is fed to **FewShotPromptTemplate** alongside the example list.
4. Lastly, using the LLMChain object, you can begin your classification straightaway with the input text.

**Output:**

```
Output:sports
```

## 10.5.3 Example Selection

In the above example, we have given just a handful of examples but what if we have 100s or 1000s of examples that we wish to use? More samples will mean better performance by an LLM but as you must be aware, under the hood, it's the LLM that is doing all these tasks and eventually every LLM has an input token restriction and as all these examples will be a part of the input prompt sent to the LLM, feeding so much data won't be allowed. What to do? We can use Example Selector with FewShotPromptTemplate to **select some examples** out of this huge pool of examples.

We left this topic in the 'Model & Prompts' chapter but it's the right time to cover it. We can select examples on multiple criteria like similarity with each other (hence ignoring redundant looking examples), based on length, random and whatnot. Let's integrate an ExampleSelector in the above code that selects examples based on Length. So the ExampleSelector will decide *which examples to incorporate such that the maximum length mentioned doesn't get crossed for the whole prompt*:

```python
from langchain.prompts.example_selector import LengthBasedExampleSelector
from langchain.prompts import FewShotPromptTemplate, PromptTemplate
from langchain.chains import LLMChain

api_key=''

examples = [
    {"input":"football is a fun game",
     "output":"sports"},
    {"input":"SRK is the best actor",
     "output":"cinema"},
    {"input":"AI startups are at a boom right now",
     "output":"tech"},
    {"input":"He is good player of basketball",
     "output":"sports"},
    {"input":"Computers are becoming powerful now",
     "output":"tech"},
    {"input":"GodFather is a higly rated movie",
     "output":"cinema"}
```

```
    ]

example_prompt=PromptTemplate(input_variables=["input","output"],template=
"input:{input}\n output:{output}")

example_selector = LengthBasedExampleSelector(
    examples=examples,
    example_prompt=example_prompt,
    max_length=25,
)
prompt = FewShotPromptTemplate(
      example_selector=example_selector,
      example_prompt=example_prompt,
      suffix="Question: {input}",
      input_variables=["input"]
  )
chain=LLMChain(llm=ChatOpenAI(openai_api_key=api_key),prompt=prompt)

print(prompt.format(input="India lost the Cricket final"))
```

As you can see, we made some marginal changes in the existing code where:

1. We have created an **ExampleSelector object** that restricts the maximum length of the final prompt (after including the examples provided). Now the example selector, very intelligently, uses only those examples in the final prompts that are short in size, keeping the final **prompt size<25** (the threshold we have set).
2. We have passed the selector a list of examples and the example prompt as well which is finally passed to the **FewShotPromptTemplate**.
3. The input prompt gets modified to the below version.

**Output:**
```
input:football is a fun game
 output:sports

input:SRK is the best actor
 output:cinema

Question: India lost the Cricket final
```

As you can see, some examples have been eliminated due to their size from the final Few-shot prompt.

# 10.6 POS Tagging, Segmentation and more

Unfortunately, LangChain doesn't have a pre-defined chain for other NLP tasks but we can create a custom template to do almost everything using ChatPromptTemplate we discussed before:

```python
from langchain.chat_models import ChatOpenAI
from langchain.prompts.chat import (
    ChatPromptTemplate,
    SystemMessagePromptTemplate,
    HumanMessagePromptTemplate,
)
from langchain.chains import LLMChain
api_key=''

template = """You are an NLP expert. Act as a {problem} machine and
give output for the input : {text}"""

system_message_prompt=SystemMessagePromptTemplate.from_template(templat
e)
human_template = "{text}"
human_message_prompt=HumanMessagePromptTemplate.from_template(human_tem
plate)

chat_prompt=ChatPromptTemplate.from_messages([system_message_prompt,
human_message_prompt])

chain = LLMChain(
    llm=ChatOpenAI(openai_api_key=api_key),
    prompt=chat_prompt
)
```

Now see the magic. We will quickly define the NLP tasks we are going to perform and observe the above app's output for the same:

## POS Tagging
POS tagging, short for Part-of-Speech tagging, is the process of assigning grammatical parts of speech to words in a sentence, such as nouns, verbs, adjectives, adverbs, etc. It helps in understanding the syntactic structure of a sentence.

An example of POS tagging with annotation:
*"The quick brown fox jumps over the lazy dog."*

POS Tagging annotation:

1. *"The" - Determiner (DT)*
2. *"quick" - Adjective (JJ)*
3. *"brown" - Adjective (JJ)*
4. *"fox" - Noun (NN)*
5. *"jumps" - Verb (VBZ)*
6. *"over" - Preposition (IN)*
7. *"the" - Determiner (DT)*
8. *"lazy" - Adjective (JJ)*
9. *"dog" - Noun (NN)*

In this annotation, each word in the sentence is tagged with its respective part of speech. This tagging helps in understanding the grammatical structure and role of each word within the sentence, which is crucial for various natural language processing tasks like information extraction, parsing and machine translation.

```
chain.run({"problem":'POS Tagging',"text":"He is not a good boy"})
```

**Output:**
*He/PRP is/VBZ not/RB a/DT good/JJ boy/NN*

## Segmentation
In NLP, Segmentation refers to dividing a continuous sequence of text into smaller, meaningful units. This can involve splitting text into words, sentences, or other meaningful segments. As we haven't mentioned which segmentation, I assume the LLM went for word segmentation below.

```
chain.run({"problem":'segmentation',"text":"He is not a good boy"})
```

**Output:**
*Segmented Output: \n1. He \n2. is \n3. not \n4. a \n5. good \n6. Boy*

## Dependency Parsing
Dependency parsing is an NLP technique used to analyze the grammatical structure of a sentence by identifying the relationships between words. It establishes links or dependencies between words, showing which words are connected and how they relate to each other within a sentence.

```
chain.run({"problem":'Dependency  parsing',"text":"He  is  not  a  good boy"})
```

**Output:**
*Here is the dependency parsing output for the input sentence "He is not a good boy":\n\n1. He – nsubj(is)\n2. is – ROOT\n3. not – neg(good)\n4. a – det(boy)\n5.*

```
good - amod(boy)\n6. boy - attr(is)\n\nThis output represents the syntactic
structure and relationships between the words in the sentence.
```

This is how we can perform different NLP tasks using LangChain. Next, we will deal with some demons LLMs have i.e. hallucinations and how to cope with hallucinations.

# Chapter 11: Handling LLM Hallucinations

It's this part of the book that deals with the biggest issue associated with LLMs i.e. Hallucinations. A jargon sometime back, GenAI had such an impact that 'hallucination' became the *word of the year in 2023 by the Cambridge Dictionary*. This chapter deals with:

*What are hallucinations?*
*Why do LLMs hallucinate?*
*How to deal with hallucinations using LangChain?*

## 11.1 What are Hallucinations?

Hallucinations in the context of artificial intelligence often refer to situations where a model generates content that is not accurate or is seemingly made up. It can involve the model ***"imagining" details*** not present in the input data or generating responses that are not grounded in reality.

For example: Suppose you have a conversational AI designed to answer questions about historical facts. You asked the following question:
***"Who was the first president of the United States?"***

In a non-hallucinating response, the AI would correctly provide information like "George Washington" as the answer. However, in a hallucination scenario, the AI might generate an inaccurate response like:
***"The first president of the United States was Benjamin Franklin."***

If you closely notice, the response is **syntactically correct but factually wrong**. So if you're ignorant and go into copy-paste mode, this can be problematic. Check this fascinating case where a lawyer used ChatGPT to prepare for a case, only to realize it hallucinated and generated fake cases that never existed.

https://www.forbes.com/sites/mollybohannon/2023/06/08/lawyer-used-chatgpt-in-court-and-cited-fake-cases-a-judge-is-considering-sanctions/?sh=767ce5977c7f

*This is very similar to when you go to an exam hall, see the question paper and when you don't know the answer, you make things up!*

But a more complex question,

## 11.2 Why do LLMs Hallucinate?

Taking a step back, we must understand how LLMs work. They don't have brains like we do. They just predict the next token given a previous sequence of tokens and input. So, in all, even an LLM doesn't know the final answer until the last word is predicted. Eventually, ***if some facts aren't***

*available with it (the LLM isn't trained on that data), the probability of other, closer events/words may become higher.* If you have worked with ChatGPT extensively, when it hallucinates, it doesn't give a very off-beat answer like '*President of India is Salman Khan*' but would be mentioning closer entities. In this case, it can be some former President of India or other political leaders but not actors or people from other fields. This is because given the input and predicted sequence till now, the probability of semantically similar entities increases. There can be many reasons why LLMs hallucinate, the major ones being:

1. **Limited Contextual Understanding:** LLMs may lack a deep contextual understanding. This limited understanding can result in the model producing hallucinatory content.
2. **Incomplete or Noisy Training Data:** Hallucinations often stem from incomplete or noisy training data. In such a case, the model may generate realistic-sounding but incorrect output.
3. **Pattern-Based Generation**: LLMs might prioritize learned patterns over factual accuracy, especially when faced with ambiguous or complex prompts. This pattern-based generation can contribute to the model producing outputs that alter facts in a seemingly realistic manner.
4. **Biases in Training Data:** The presence of biases in the training data can influence how LLMs generate content, potentially leading to altering facts to align with biased perspectives.

Now that we know hallucinations in quite some detail, time to know different ways to handle them using LangChain. For this, LangChain has special chains called LLMCheckersChain and LLMSummarizationCheckerChain. Let's have a brief around these special chains first before jumping onto codes:

1. **LLMCheckerChain:** This chain is more like a Fact-checking chain that tells you whether the input is factually correct. So if you ask it '*We celebrated the 500th birthday of my Father*', the chain should output: *This isn't possible as a human's lifespan can't be 500 years.*
2. **LLMSummarizationChain:** Going a step ahead of LLMCheckerChain, this helps you to rewrite your input with correct information/facts and not just tell you whether the input is right or wrong. So, if you input the wrong rules of a game say cricket, this chain will automatically correct the wrong facts in the output.

## 11.3 LLMCheckerChain

We will first test the LLMCheckerChain to comment on a fact that I will pass it.

```python
from langchain.chains import LLMCheckerChain
from langchain.llms import OpenAI


api_key = 'your API'
llm = OpenAI(openai_api_key=api_key)


checker_chain = LLMCheckerChain.from_llm(llm, verbose=True)
```

```
text = "Roger's 3 months old baby is speaking fluent English"
print(text,'\n',checker_chain.run(text))
```

**Output:**

> *Entering new LLMCheckerChain chain...*

> *Entering new SequentialChain chain...*

> *Finished chain.*

> *Finished chain.*

*Roger's 3 months old baby is speaking fluent English*

*Roger's 3 month old baby is likely not speaking fluent English. Babies at 3 months old begin to develop language skills, but it usually takes time to develop the physical and cognitive skills necessary to speak any language fluently.*

As you can see, LLMCheckerChain is helping to correct the facts in the final output.
Now, let's re-write the steps to make tea using the other chain.

## 11.4 LLMSummarizationChain

```
from langchain.chains import LLMSummarizationCheckerChain
from langchain.llms import OpenAI

api_key = 'your API'
llm = OpenAI(temperature=0,openai_api_key=api_key)

checker_chain=LLMSummarizationCheckerChain.from_llm(llm,verbose=True,
max_checks=1)

text ="""
Following is the process of making tea
1. Boil tea leaves
2. Add chilli powder to cold water
3. Stir a little
4. Add some salt
5. Add some milk and boil
"""
print(checker_chain.run(text))
```

**Output:**

> *Entering new LLMSummarizationCheckerChain chain...*

> *Entering new SequentialChain chain...*

> *Entering new LLMChain chain…*


*Prompt after formatting:*

*Given some text, extract a list of facts from the text.*

*Format your output as a bulleted list.*

**Text:**
*"""*
*Following is the process of making tea*
*1. Boil tea leaves*
*2. Add chilli powder to cold water*
*3. Stir a little*
*4. Add some salt*
*5. Add some milk and boil*
*"""*
**Facts:**
**> Finished chain.**
**> Entering new LLMChain chain...**
*Prompt after formatting:*
*You are an expert fact checker. You have been hired by a major news organization to fact check a very important story.*

*Here is a bullet point list of facts:*
*"""*
*- Boiling tea leaves is the first step in making tea*
*- Chilli powder is added to cold water*
*- Stirring is involved in the process*
*- Salt is added to the tea*
*- Milk is added and boiled*
*"""*
*For each fact, determine whether it is true or false about the subject. If you are unable to determine whether the fact is true or false, output "Undetermined". If the fact is false, explain why.*

**> Finished chain.**
**> Entering new LLMChain chain...**
*Prompt after formatting:*
*Below are some assertions that have been fact checked and are labeled as true or false. If the answer is false, a suggestion is given for a correction.*

**Checked Assertions:**
*"""*
*- Boiling tea leaves is the first step in making tea: True*
*- Chilli powder is added to cold water: False - Chilli powder is not typically added to tea.*
*- Stirring is involved in the process: True*

- Salt is added to the tea: False - Salt is not typically added to tea.
- Milk is added and boiled: False - Milk is typically added after the tea has been boiled.
"""

***Original Summary:***
"""

Following is the process of making tea
1. Boil tea leaves
2. Add chilli powder to cold water
3. Stir a little
4. Add some salt
5. Add some milk and boil

"""

Using these checked assertions, rewrite the original summary to be completely true.
The output should have the same structure and formatting as the original summary.

***Summary:***
**> Finished chain.**
**> Entering new LLMChain chain...**
Prompt after formatting:
Below are some assertions that have been fact checked and are labeled as true or false.

If all of the assertions are true, return "True". If any of the assertions are false, return "False".

Here are some examples:
===
***Checked Assertions: """***
- The sky is red: False
- Water is made of lava: False
- The sun is a star: True
"""
***Result:*** False
===

***Checked Assertions: """***
- The sky is blue: True
- Water is wet: True
- The sun is a star: True
"""
***Result:*** True

```
===
Checked Assertions: """
- The sky is blue - True
- Water is made of lava- False
- The sun is a star - True
"""
Result: False
===
Checked Assertions:"""
- Boiling tea leaves is the first step in making tea: True
- Chilli powder is added to cold water: False - Chilli powder is not typically
added to tea.
- Stirring is involved in the process: True
- Salt is added to the tea: False - Salt is not typically added to tea.
- Milk is added and boiled: False - Milk is typically added after the tea has
been boiled.
"""
Result:

> Finished chain.
> Finished chain.

"""
Following is the process of making tea
1. Boil tea leaves
2. Stir a little
3. Add milk after boiling
4. Enjoy!
```

If you read this long output, you will notice a few things:

1. Internally, this summarization chain is running many other chains.
2. It is first checking each fact mentioned one by one alongside the reason if it is wrong.
3. It does include Few-shot prompts to yield better results.
4. Finally rewriting the factually incorrect steps.

This isn't a single-step process and may lead to many LLM hits.

## 11.5 Avoiding Hallucinations using RAG

We have already discussed RAG in some detail in the previous chapter. When we use external contexts like text or PDF files, we can add in our prompt to **mention the context being used** for

answering hence avoiding hallucinations as, if the facts aren't present, the LLM shouldn't be able to cook stories on its own as it needs to give proof in the output from the context.

For example: Imagine you're asking a chatbot about a rare bird. Instead of relying solely on its knowledge (which may be limited), the RAG-assisted model fetches information from a bird database, ensuring the chatbot's response aligns more accurately with factual data. With RAG, the model cross-checks its generated responses against real-world databases, reducing the chances of providing inaccurate or misleading information. But again, not a 100% guarantee that the LLM won't hallucinate even after using external context.

We have talked a lot about different applications one can build with LangChain. But how to evaluate these applications? And LLMs in general? Let's check it out in the next section.

# Chapter 12: Evaluating LLMs

---

For any Machine Learning problem, the final result boils down to the metrics. If everything goes great but eventually the metrics are bad, the whole project can be shelved. So evaluation metrics are amongst the most important aspects of any ML-based project. When it comes to LLMs, this becomes even more crucial as at times, LLMs hallucinate and you don't know whether the answer is right or wrong. But evaluating LLMs isn't as straightforward as you think. Assume your ground truth for some problem statement is 11, but the LLM can give the following answers:

1. *Eleven*
2. *The answer is Eleven*
3. *…………..Eleven………*
4. *The answer is 11*
5. *11 is the answer*

& many other variations

This is a common issue we have, with not just LLMs but with most Generative models (like GANs) where at times:

1. You don't know what is a good output.
2. Even though the LLM gave the right answer, it is usually wrapped with other text or in different formats.

LangChain has brought in a pretty good evaluation system that can help in evaluating the outputs from LLMs at different levels, for both supervised and unsupervised problems. This evaluation system is divided into 3 parts:

**<u>String evaluators</u>**
String evaluator in LangChain is a tool that checks how well a language model performs by comparing what it generates with a set text or input. This comparison is super important to see if the model's predictions are accurate.

Usually, String Evaluators check a predicted text against an input, like a question or a starting text. They can even use a known correct answer or ground truth to judge how good the model's response is. These evaluators can be changed to suit exactly what you need for your app's evaluation process.

**<u>Comparison evaluators</u>**
Comparison evaluators in LangChain are tools that measure outputs from two chains or language models. These evaluators use the PairwiseStringEvaluator class as a base, which ***compares two strings***. Usually, it compares outputs from different prompts or models, or different versions of the

same model. A comparison evaluator checks two strings and gives back a dictionary with the score and other important info.

**Trajectory evaluators**

Trajectory evaluators in LangChain offer a comprehensive way to **assess an agent**. Instead of focussing on individual actions, these evaluators consider the entire sequence of actions and their corresponding responses, known as the "trajectory." This approach provides a more thorough measurement of an agent's effectiveness and capabilities.

In this book, we would be focussing majorly on String Evaluators that too a particular subclass i.e. Criteria evaluation alongside custom evaluators with a basic introduction to others. The reasons are:

1. **Comparison evaluators** are not something unique and I assume you can achieve this by a simple Python functions. The same goes with other classes in String Evaluators as it is more of an overkill to wrap basic functionalities under LangChain.
2. **Trajectory evaluators** don't appear to be of much use to common folks. Interested readers can go for LangChain documentation for a deep dive.

So, let's get started.

# 12.1 String Evaluators

## 12.1.1 Criteria Evaluators

When you want to evaluate a model's output based on a particular rubric or set of criteria, the criteria evaluator becomes a useful tool. It enables you to check whether the output of an LLM or chain aligns with a predefined set of criteria. We have two types of Criteria evaluators, one that requires a ground truth and another which don't need a ground truth.

Below is the list of pre-defined Criteria Evaluators:

1. **Insensitivity**: Measures the lack of consideration or empathy in the output.
2. **Relevance**: Assess the degree to which the output is pertinent and applicable.
3. **Helpfulness**: Measures the degree to which the output is beneficial or supportive.
4. **Maliciousness**: Assess whether the output contains harmful or malevolent intent.
5. **Harmfulness**: Gauges the potential for causing harm or negative impact.
6. **Correctness**: Evaluate the accuracy and precision of the information provided.
7. **Coherence:** Examines the logical and orderly connection of ideas in the output.
8. **Conciseness**: Measures the brevity and clarity of the output.
9. **Misogyny**: Assess the presence of content displaying hatred or prejudice against women.
10. **Criminality**: Examines the potential involvement or endorsement of criminal activities.
11. **Controversiality**: Evaluates the contentious or disputed nature of the content.

**Note:** Out of all these, only Correctness and Relevance need ground truth.

For demonstration purposes, we will first set up an LLMChain for a general use case and test for some criteria evaluators over the predictions made by this chain for given prompts.

```python
from langchain.chat_models import ChatOpenAI
from langchain.prompts.chat import (
    ChatPromptTemplate,
    SystemMessagePromptTemplate,
    HumanMessagePromptTemplate,
)
from langchain.chains import LLMChain
from langchain.evaluation import load_evaluator
from langchain.evaluation import EvaluatorType
import os

api_key = ''
os.environ['OPENAI_API_KEY'] = api_key

template="""You  are  a  Helpful  Assistant  that  explains  everything  being
asked"""

system_message_prompt=SystemMessagePromptTemplate.from_template(template)
human_template = "{text}"
human_message_prompt                                        =
HumanMessagePromptTemplate.from_template(human_template)

chat_prompt=ChatPromptTemplate.from_messages([system_message_prompt,
human_message_prompt])

chain = LLMChain(llm=ChatOpenAI(),prompt=chat_prompt)

def evaluate(criteria, sentences):
    evaluator=load_evaluator(EvaluatorType.CRITERIA,criteria=criteria,
llm=llm)
    for prompt in sentences:
        prediction = chain.run(prompt)
        eval_result = evaluator.evaluate_strings(
            prediction= prediction,
            input=prompt,
        )
        print('PROMPT : ',prompt)
        print('RESULT : ',prediction)
```

```
        print('VALUE : ',eval_result['value'])
        print('SCORE : ',eval_result['score'])
        print('REASON : ',eval_result['reasoning'])

evaluate('conciseness',['Explain road not taken by Robert Frost'])
evaluate('creativity',['What would happen if everyone becomes immortal?'])
```

**Output:**

*#conciseness*

**PROMPT :**  *Explain road not taken by Robert Frost*

**RESULT :**

*"The Road Not Taken" is a famous poem written by Robert Frost*

*It was first published in 1916 as part of Frost's collection of poems titled "Mountain Interval*

*" The poem is often interpreted as an exploration of choices and the consequences they bring*

*The poem begins with the speaker coming across a fork in the road while walking in the woods*

*The speaker is faced with a decision between two paths and  he expresses his dilemma by saying, "Two roads diverged in a yellow wood*

*" The two roads symbolize different life choices or paths that one can take*

*The speaker observes that both paths are equally worn and attractive, making it difficult for him to choose*

*He carefully examines one of the paths but realizes that he cannot see where it leads*

*He then turns his attention to the other path, which he describes as "grassy and wanted wear*

*" This suggests that the less traveled path is more appealing to him because it represents something new and unexplored*

*After considering both paths, the speaker makes a decision*

*He says, "I took the one less traveled by and  that has made all the difference*

*" This line is often misinterpreted as a celebration of taking the unconventional path and being rewarded for it*

*However, Frost himself noted that the poem is often misunderstood in this way*

*Instead, it is a reflection on the human tendency to look back on past choices and wonder what might have happened if different decisions were made*

*The poem can be seen as a metaphor for life itself, with the two roads representing the choices and opportunities that we encounter throughout our lives*

*It serves as a reminder that every decision we make has consequences and shapes our future*

**REASON :**

*The criterion for this assessment is conciseness*

*This means the submission should be brief, to the point and  not include unnecessary information or details*

98

*Looking at the submission, it provides a detailed explanation of Robert Frost's "The Road Not Taken*
*" It includes an overview of the poem, an analysis of its main themes and a discussion of its common misinterpretations*

*While the submission is thorough and informative, it may not be considered concise*

*It goes into great detail about the poem's themes and interpretations, which, while valuable, may not be necessary for a concise explanation*

*However, considering the complexity of the poem and the various interpretations it has, the level of detail in the submission could be seen as necessary for a comprehensive understanding*

*It's a bit of a judgment call, but given the depth of the poem and the need to clarify common misinterpretations, the submission could be seen as being as concise as possible while still providing a thorough explanation*

*So, based on this reasoning, the submission can be considered to meet the criterion of conciseness*

**VALUE :** *Y*

**SCORE :** *1*


*#creativity*

**PROMPT :** *What would happen if everyone becomes immortal?*

**RESULT :**

*If everyone became immortal, it would have significant and far-reaching implications for individuals, society and the world as a whole*

*Here are a few possible consequences:1*

*Overpopulation: With no natural death, the human population would continue to grow exponentially*

*This would put a strain on resources such as food, water and living space, leading to potential conflicts and competition for limited resources*

*2*

*Social and Economic Disruption: Immortality could disrupt the existing social and economic systems*

*Traditional life stages like childhood, education and retirement would become obsolete and new systems would need to be put in place to accommodate the extended lifespans*

*The concept of work and retirement would undergo significant changes, as people might have to work for much longer periods to sustain their extended lives*

*3*

*Psychological and Emotional Impact: Immortality could have profound psychological effects on individuals*

*People might struggle to find meaning and purpose in life when there is no end in sight*

*The experience of losing loved ones could become more challenging, as they would not have the release of death*

*This could lead to emotional and mental health issues*
*4*
*Scientific Advancement: Immortality could drive unprecedented advancements in medical and scientific research*
*With unlimited time, individuals would have the opportunity to acquire vast knowledge and expertise in various fields*
*This could lead to breakthroughs in medicine, technology and other areas, potentially solving many of humanity's current problems*
*5*
*Evolutionary Changes: Immortality could potentially lead to changes in human evolution*
**REASON :**
*The criterion in question is creativity, which is defined as the demonstration of novelty or unique ideas*
*Looking at the submission, the respondent has provided a detailed and comprehensive answer to the question*
*They have considered various aspects of society and individual life that could be affected if everyone became immortal, including overpopulation, social and economic disruption, psychological and emotional impact, scientific advancement and evolutionary changes*
*However, when assessing the creativity of the response, it's important to consider whether these ideas are novel or unique*
*The concepts presented in the response, while well thought out and articulated, are not particularly novel or unique*
*They are fairly common considerations when discussing the concept of immortality and can be found in many discussions, books and movies on the subject*
*Therefore, while the response is thorough and well-reasoned, it does not demonstrate a high level of creativity as defined by the criterion*
*The ideas presented are not particularly novel or unique, but rather common considerations in discussions of immortality*
**VALUE :** *N*
**SCORE :** *0*

The output for each evaluation usually comprises three values:

1. **Score**: 0 or 1
2. **Reason:** Explanation for the assigned score/value to the prediction.
3. **Value:** Yes or No.

Considering 'conciseness', if the prediction made by the LLM isn't concise, the Score=0, Value=No and Reason would be given as to why the score was assigned. If you give the output a read, especially the reasons given are quite apt and detailed. Though conciseness was the metric, the evaluator understood the poem 'Road Not Taken' is very complex and requires a long explanation

hence the score is 1. Similarly for creativity, though the answer is correct, it doesn't appear to be creative and eventually the score is 0.

Now check for Criteria that need a ground truth. The evaluation code will change slightly for this:

```python
def evaluate(criteria,sentences,ground_truth):
    evaluator=load_evaluator("labeled_criteria",criteria=criteria, llm=llm)
    for index,prompt in enumerate(sentences):
        prediction = chain.run(prompt)
        eval_result = evaluator.evaluate_strings(
            prediction= prediction,
            input=prompt,
            reference = ground_truth[index]
        )
        print(eval_result)

evaluate('correctness',['How many players are required for chess'],['2'])
```

As you can see,

1. Ground Truth is also passed to the evaluator object alongside prompt and prediction.
2. We need to pass the '*labeled_criteria*' flag while creating the evaluator object.
3. The response structure remains the same as above.

The best part about this 'correctness' criteria is even if the prediction is 'two' or 'answer is two', the score would be 1 and not a 0. Hence it can make out that the result is correct irrespective of the format of the output.

**Output:**
```
PROMPT :  How many players are required for chess
RESULT :
 Two players are required for chess
VALUE :  Y
SCORE :  1
REASON :
 The criterion for this task is the correctness of the submitted answer
 The input asks for the number of players required for chess
 The submitted answer states that two players are required for chess
 This is indeed correct, as chess is a two-player game
 The reference answer also confirms this, as it states "2", which is the same as
the submitted answer, just in numerical form
 Therefore, the submitted answer meets the criterion of correctness
```

Taking a step further, we can even build custom evaluation metrics. Let's check out:

## 12.1.2 Custom Evaluators

We will be first building a **scorecard sort of metric** for evaluating Supervised problems.

```python
api_key=''

score_criteria = {
    "custom_metric": """

0: Completely incorrect
1-3: Partially incorrect or minimally relevant
4-6: Mix of correct and incorrect information, with notable errors
7-8: Mostly accurate, with minor errors or omissions
9: Nearly flawless, with very minor errors or negligible missing
details
10: Flawless, meeting or exceeding expectations"""
}
prompt = 'How does the life of a general human looks like? starting
from being a baby to being old'

evaluator=load_evaluator("labeled_score_string",criteria=score_criteria
, llm=ChatOpenAI(model="gpt-4",openai_api_key=api_key))

prediction = chain.run(prompt)

eval_result = evaluator.evaluate_strings(
    prediction= prediction,
     reference="A human life follows a general trajectory from infancy
to old age.\
        In infancy, rapid physical and cognitive development occurs,
followed by early childhood where formal education begins.\
     Adolescence brings puberty and identity formation, leading to young
adulthood marked by career and relationship pursuits.\
        Middle adulthood sees professional and family responsibilities,
while late adulthood involves retirement,\
          reflection and  potential health challenges. Old age is
characterized by increased dependency,\
     wisdom-sharing and  end-of-life considerations.\
        Individual experiences vary based on factors like culture and
personal choices,\
     and improvements in healthcare can influence life trajectories.",
```

```
     input=prompt
)
print('LLMs answer:','\n'.join(prediction.split('.')))
print('\n'.join(eval_result['reasoning'].split('.')))
```

The code is similar to the above examples in the chapter but there are some subtle changes:

1. We need to create a dictionary object, name your custom metric as a key and  for value, lay down your score chart. In this case, depending upon the level of accuracy with the ground truth, the score increases. It can be used for say, a summarization task evaluation.
2. The label passed while creating the evaluator object is *'labeled_score_string'* and the dictionary object is passed as well.
3. The output format remains the same as above examples.

**Output:**

*LLMs answer: The life of a general human typically follows a series of stages and experiences, starting from being a baby and progressing through childhood, adolescence, adulthood and  finally old age*

*Let's explore these stages in more detail:*

*1*

*Babyhood: This stage begins at birth and lasts until around two years old*

*Babies are completely dependent on their caregivers for their basic needs, such as feeding, diaper changes and  comfort*

*They gradually develop motor skills, learn to communicate through sounds and gestures and  form attachments to their primary caregivers*

*2*

*Childhood: Childhood spans from around two years old to adolescence*

*During this stage, children grow physically, develop more advanced motor skills and  acquire language and cognitive abilities*

*They begin to explore their surroundings, attend school and  develop social skills through interactions with peers*

*Childhood is often characterized by curiosity, imagination, playfulness and rapid psychological and emotional development*

*3*

*Adolescence: Adolescence is the transitional stage between childhood and adulthood, typically occurring between the ages of 10 and 19*

*This period is marked by significant physical and hormonal changes as individuals enter puberty*

*Adolescents strive for independence, establish their identity and  begin to develop their own beliefs, values and  goals*

*It is also a time of increased self-awareness, emotional intensity and  peer influence*

*4*

*Adulthood: Adulthood typically starts in the late teens or early twenties and continues until old age*
*This stage is characterized by increased independence, responsibility and the pursuit of personal and professional goals*
*Many individuals choose to pursue higher education, establish careers and form long-term relationships such as marriage or committed partnerships*
*They may also start families and become parents*
*5*
*Old age: Old age, also known as the senior or elderly stage, usually begins around the age of 65 and beyond*
*This stage is marked by physical, cognitive and sensory changes*
*Individuals often retire from their careers and may experience a decline in health and mobility*
*However, many people continue to engage in social and leisure activities, maintain relationships and contribute to their communities*
*This stage can also bring reflection, wisdom and a focus on maintaining a fulfilling quality of life*

*It's important to note that individual experiences may vary and not everyone will follow this exact progression or timeline*
*Additionally, cultural, social and personal factors can greatly influence how each stage of life is experienced*
***REASON:***
*The assistant provides a comprehensive and detailed response to the user's question*
*It breaks down the human life trajectory into five distinct stages—babyhood, childhood, adolescence, adulthood and old age—and describes each stage in detail*
*The assistant also correctly notes that individual experiences may vary and factors such as culture, society and personal choices can greatly influence how each stage of life is experienced*
*The response is both accurate and relevant to the question*
***Rating:** [[10]]*

As you can see, even though the reference text isn't word to-word, the evaluator can determine that the essence is the same as the ground truth hence score=10. Now, how can we build ***custom criteria for unsupervised problems***?

```
custom_criteria = {
    "Humor": "The assistant's answer should have a sense of humor",
    "Impact":"What kind of lasting impression does the text leave on the reader?"
}
prompt = 'Tell a joke I can crack in front of teenagers'
```

```
evaluator=load_evaluator("score_string",criteria=custom_criteria,llm=
ChatOpenAI(openai_api_key=api_key))


prediction = chain.run(prompt)
eval_result = evaluator.evaluate_strings(
    prediction=prediction,
    input=prompt)
print('LLM answer:','\n'.join(prediction.split('.')))
print('\n'.join(eval_result['reasoning'].split('.')))
```

Similar to the above example, here also you need to make a dictionary object and the label passed to the evaluator changes to *'score_string'*. Also, as we have mentioned multiple custom criteria, the final scoring would be an *average of the two criteria rather than two separate outputs for each custom criterion.*

**Output:**
```
LLM answer: Sure, here's a joke that teenagers might enjoy:
Why don't scientists trust atoms?
Because they make up everything!
Explanation:
The assistant's response demonstrates a good sense of humor by providing a
pun-style joke that teenagers might find amusing
 The joke plays on the double meaning of "make up" to create a humorous twist
 It leaves a light and playful impression on the reader


Rating: [[9]]
```

Other string evaluators that are available in LangChain other than Criteria are:

1. **Regex:** Utilizes regular expressions (regex) to match and evaluate predicted strings based on patterns. This evaluator enables complex pattern matching for string analysis and validation.
2. **Score:** Assigns a numerical score to strings based on predefined rules or comparisons. It quantifies the similarity or closeness of a string to a reference string or pattern. The custom evaluators we built in the above tutorial are score-based.
3. **Embedding Distance:** Measures the semantic similarity or distance between prediction and ground truth using embedding models. These models encode words or phrases into numerical vectors, enabling comparisons based on their meanings or contexts.
4. **String Distance:** Computes the dissimilarity or distance between prediction and ground truth, often employing algorithms like Levenshtein distance or other string metric methods.
5. **Exact Match:** Validates if a predicted string matches precisely with ground truth or a predefined pattern without any deviation or variation.
6. **JSON Evaluators:** Allow complex evaluation and comparison of strings within JSON structures. These evaluators facilitate assessments based on JSON elements, keys, values or structures. There are multiple subclasses in json evaluators which I am ignoring for now.

Now very quickly, we will discuss the other two types of Evaluators.

## 12.2 Comparison Evaluators

As already mentioned, these evaluators help in testing two approaches together hence the name comparison evaluators. It is of 3 types:

1. **Pairwise string:** Given two predictions with ground truth, which prediction is more apt.
2. **Pairwise embedding:** Given two predictions with ground truth, it outputs the embedding distance between them. Hence can be used for prediction similarity.
3. **Custom Pairwise:** As the name suggests, you can use it to create any custom metric for comparing two predictions.

We will try out the ***Pairwise string evaluator*** in this chapter.

```python
from langchain.llms import OpenAI
from langchain.evaluation.comparison import PairwiseStringEvalChain

api_key=''
llm = OpenAI(openai_api_key=api_key)

chain = PairwiseStringEvalChain.from_llm(llm=llm)
result = chain.evaluate_string_pairs(
    input = "What is 2+3?",
    prediction = "five is the answer",
    prediction_b = "If I add 2+3, I might get 6.",
    reference = "5",
)
print(result)
```

The code is quite easy where we are feeding the two predictions alongside the ground truth to the pairwise string evaluator.

**Output:**
```
{'reasoning': 'Response A is correct as 2+3 is 5',
'value': 'A',
'score': 1}
```

## 12.3 Trajectory Evaluators

As already explained, such evaluators don't check for the output but the sequence of actions or steps followed by agents to reach the output. Do remember they can only be used with agents and nothing else.

```python
from langchain.evaluation import import load_evaluator


evaluator = load_evaluator("trajectory",llm=OpenAI())
evaluation_result = evaluator.evaluate_agent_trajectory(
    prediction=result["output"],
    input=result["input"],
    agent_trajectory=result["intermediate_steps"],
)
evaluation_result
```

When you run your agent, just save the results in the variable 'result' and run the above code.

**Output:**
```
{'score': 1.0,
 'reasoning': "The agent accurately uses the tool provided to get to correct
answer. No extra steps were observed"}
```

With this, we will wrap evaluators here. But I would strongly suggest you to try out all the evaluators and choose which one is the best for your use case. The next section talks about how advanced prompt engineering techniques like ReAct, ToT, etc. can be used with LangChain.

# Chapter 13: Advanced Prompt Engineering

---

After working with LLMs for the last year, one thing I'm quite sure of is Prompt Engineering is a very crucial concept one must know to yield the best results from any LLM, be it LLaMA, PaLM, Claude, or any other LLM. If you take the strongest of models and provide a bad prompt, the results will be trash. I hope most of you are aware of basic prompt engineering techniques. This chapter is more toward advanced prompt engineering techniques and how to integrate them using LangChain.

Before starting, install the packages that we will be using:

```
!pip install langchain==0.0.349 openai==1.3.8
google-search-results==2.4.2 langchain-experimental==0.0.46
```

We will start with:

## 13.1 Chain of Thoughts

Chain of Thought prompting is a novel method developed to enhance the reasoning capabilities of LLMs. It encourages LLMs to explain their reasoning in a step-by-step manner, providing transparency into the model's decision-making process.

Consider the prompt:
***"You plan to organize a picnic. Think step by step through every requirement"***

The Chain of Thoughts process involves the LLM reasoning through each step:

1. ***Decide on a location****: The model explains the choice of a suitable picnic spot.*
2. ***Invite friends:*** *The model elaborates on the process of inviting friends and considerations for the guest list.*
3. ***Prepare snacks:*** *The model details the snack selection, preparation steps and any special considerations.*
4. ***Set a time:*** *The model provides the reasoning behind choosing a specific time, considering factors like weather or participant availability.*

If you go with an ordinary prompt like ***"organize a picnic plan",*** you may not yield such a detailed result. CoT can be achieved in multiple ways in LangChain.

### 13.1.1 Think Step by Step

This involves adding ***'think step by step'*** or a similar phrase in your prompt to make LLMs think in a planned manner. Below is an example:

```python
from langchain.chains import LLMChain
```

```python
from langchain.llms import OpenAI
from langchain.prompts import PromptTemplate


api_key=''


prompt_template = "I went to the market and bought 20 mangoes. I gave 5
apples to the Ram and 3 to the House-help. I then went and bought 12 more
mangoes and ate 2.How many mangoes did I remain with? Let's think step by
step."


llm = OpenAI(temperature=0,openai_api_key=api_key)


llm_chain=LLMChain(llm=llm,prompt=PromptTemplate.from_template(prompt_templat
e))


llm_chain.run({})
```

**Output:**

```
Certainly, let's break it down step by step:
Initial mangoes bought: You started with 20 mangoes.
Apples given away: You gave 5 apples to Ram and 3 apples to the house-help. Hence
20-8= 12
Additional mangoes bought: You went and bought 12 more mangoes. Now, you have 12
(initial) + 12 (additional) = 24 mangoes.
Mangoes eaten: You ate 2 mangoes. Subtract 2 from the total: 24 - 2 = 22.
So, after these steps, you have 22 mangoes remaining.
```

## 13.1.2 Few-Shot Prompting

Few-shot prompting is a gun method to improve your results by miles for any task and not just CoT. We already saw an example of how you can perform classification using Few-shot in the previous chapter on LangChain for NLP. Similarly, you can guide the LLM on how to think and follow the same logic as in the examples given. Let's see:

```python
from langchain.prompts.few_shot import FewShotPromptTemplate
from langchain.prompts.prompt import PromptTemplate
from langchain.chat_models import ChatOpenAI
from langchain.chains import LLMChain


examples = [
    {
        "Question": "Yes or no: Would a rubber duck sink in water?",
            "Reasoning": "The density of a typical rubber duck is less than
water, typically around 0.3 g/cm³. Therefore, a rubber duck would float.",
```

```
            "Answer": "No"
    },
    {
        "Question": "Yes or no: Would a wooden block sink in water?",
         "Reasoning": "The density of wood is generally less than water,
averaging around 0.6 to 0.9 g/cm³. Therefore, a wooden block would likely
float.",
        "Answer": "No"
    },
    {
         "Question": "Yes or no: Would a plastic bottle filled with air sink
in water?",
          "Reasoning": "The density of air is much less than water. If the
plastic bottle is filled with air, its overall density would be lower than
water, causing it to float.",
        "Answer": "No"
    },
    {
        "Question": "Yes or no: Would a metal coin sink in water?",
         "Reasoning": "Metals typically have higher densities than water.
Therefore, a metal coin would likely sink in water.",
        "Answer": "Yes"
    },
    {
            "Question": "Yes or no: Would a helium-filled balloon sink in
water?",
          "Reasoning": "Helium is less dense than air and  a helium-filled
balloon would be buoyant. It would rise in air and, in water, it would
float.",
        "Answer": "No"
    }
]
example_prompt=PromptTemplate(input_variables=["Question","Reasoning","Answer
",template=":Question:{Question}\n Reasoning:{Reasoning}\n Answer:{Answer}")

prompt = FewShotPromptTemplate(
     examples=examples,
     example_prompt=example_prompt,
     suffix="Question: {input}",
     input_variables=["input"]
  )
```

```
api_key=''
chain=LLMChain(llm=ChatOpenAI(openai_api_key=api_key),prompt=prompt)
chain.run("Yes or no: will plastic ball sink in water?")
```

**Output:**

*Reasoning: The density of a plastic ball is typically less than water. Therefore, a plastic ball would likely float in water.\nAnswer: No*

As you can see, the LLM has adopted the answering style used in the examples provided to answer the question asked.

# 13.2 ReAct

ReAct, surpassing the efficacy of Chain of Thoughts, introduces a novel approach by seamlessly integrating reasoning and action within a single prompt. When tackling intricate problems, ReAct systematically guides the LLM through a comprehensive process for each subtask. This involves three fundamental steps:

1. **Reasoning (Reason)**: Evaluate and understand the current state of the problem. Analyze the requirements and determine what needs to be accomplished.
2. **Action (Acton)**: Initiate a strategic action based on the derived reasoning. Implement a solution or take steps toward addressing the identified requirements.
3. **Observation**: Observe and analyze the outcomes following the execution of the action. Assess the impact of the implemented solution on the overall problem.

By incorporating these three interconnected elements—reasoning, action and observation—ReAct ensures a holistic approach to problem-solving. ***It enables the LLM to not only understand and evaluate a situation but also actively engage in problem resolution by taking appropriate actions***. The subsequent observation phase allows for a continuous feedback loop, refining the model's understanding and enhancing its ability to handle complex tasks effectively. This dynamic process distinguishes ReAct as an advanced prompting technique, offering a more comprehensive and action-oriented strategy compared to Chain of Thoughts.

A ReAct prompt may look something like this:

*Address the assigned task by adhering to the **sequence of Thought, Action and Observation** at each juncture.*
*Thought involves analyzing the present scenario and determining the subsequent actions to undertake.*
*Observation entails evaluating the outcomes derived from executing the specified Action.*

*Actions are categorized into three types:*
***Search [item]:***

```
Conduct research on the designated item using available knowledge bases or the
internet.
Follow-up [keyword]:
Formulate additional inquiries based on the provided keyword to delve deeper into
the topic.
Finish [answer]:
Provide a conclusive response to the task at hand.


For instance, consider the task: "How to prepare tea?"
```

But you don't need to worry as LangChain has provided a pre-defined framework to implement ReAct. Let's check it out:

```python
from langchain.agents import AgentType,initialize_agent,load_tools
from langchain.llms import OpenAI
import os

api_key = ''
serp_api = ''
os.environ['OPENAI_API_KEY'] = api_key
os.environ['SERPAPI_API_KEY']= serp_api

llm = OpenAI(temperature=0)
tools = load_tools(["serpapi", "llm-math"], llm=llm)

agent_executor=initialize_agent(tools,llm,agent=AgentType.ZERO_SHOT_REACT_DES
CRIPTION,verbose=True)

agent_executor.invoke(
     {"input": "How to prepare Irani Tea famous in Hyderabad? explain from
scratch", "handle_parsing_errors": "True"
    })
```

A few things we must know about this code is:

1. In LangChain, One of the agent types is of ReAct settings, hence you need to use a specific agent if you wish for ReAct prompting.
2. This type is specified while initializing the agent where alongside tools and LLM, we are passing the agent type as well.
3. The rest of the code is similar to what we have discussed in the previous chapter on agents.

**Output:**

```
> Entering new AgentExecutor chain...
```

*I need to find a recipe*

**Action:** *Search*

**Action Input:** *Irani Tea recipe*

**Observation:** *['Ingredients (US cup = 240ml ) · □ 1½ cups (360 ml) water · □ 2 tablespoons black tea (loose tea leaves or dust) · □ 6 green cardamoms · □ 1½ ...', 'Ingredients · 3 tablespoon Assam tea leaves · 1 cup Milk · 2 tablespoon Fresh cream · 1 Star anise · 2 Cardamom (Elaichi) Pods/Seeds · 1 Black ...', 'Take 1/3 cup of water in a pan. · Let the black tea mixture simmer for atleast 5 minutes on low flame. · Then add 500 ml of buffalo milk & turn up ...', '1/2 cup vermicelli (thin wheat noodles) · 1 liter full-fat milk · 1/2 cup sugar (adjust to taste) · 1/4 cup ghee (clarified butter) · 1/4 cup ...', 'Step 1 / 4 Boil the tea. In a saucepan, add 4 cups water, cardamom pods, cinnamon, tea powder and sugar. Let it come to a boil.', 'Ingredients ; 3 cups water ; 1 1/2 tbsp black tea powder ; 3 cardamom pods optional but recommended ; 1 1/2 cup milk ; 2-3 heaped tbsp condensed milk ...', 'This tea is slowly brewed (dum) with warm spices such as green cardamom pods and cinnamon and is then mixed with a creamy reduced milk mixture.']*

**Thought:** *I need to choose a recipe*

**Action:** *Search*

**Action Input:** *Irani Tea recipe step by step*

**Observation:** *['Thick, creamy and delicious Irani Chai made with black tea, milk, sugar and cardamoms. Serve this with Biscuits or cookies. Print Recipe Pin ...', 'Step 1 / 4 Boil the tea. In a saucepan, add 4 cups water, cardamom pods, cinnamon, tea powder and  sugar. Let it come to a boil.', 'Take 1/3 cup of water in a pan. · Let the black tea mixture simmer for atleast 5 minutes on low flame. · Then add 500 ml of buffalo milk & turn up ...', "First put the water on to the heat and when it's boiling add a small piece of crushed ginger. Let it boil for a minute and add the milk. When ...", 'Ingredients ; 3 cups water ; 1 1/2 tbsp black tea powder ; 3 cardamom pods optional but recommended ; 1 1/2 cup milk ; 2-3 heaped tbsp condensed milk ...', 'Ingredients · 3 tablespoon Assam tea leaves · 1 cup Milk · 2 tablespoon Fresh cream · 1 Star anise · 2 Cardamom (Elaichi) Pods/Seeds · 1 Black ...', 'Ingredients. 1x 2x 3x · □ 1 cup water · □ 2 tbsp tea powder · □ 1 tbsp sugar · □ 1 inch ginger · □ 6 pods cardamom · □ ½ tsp pepper · □ 1 inch ...', 'For tea decoction ... Take a deep pan but make sure it must have a lid. ... Add water in it. ... Add tea leaves. ... Finally add sugar & put the pan on ...']*

**Thought:** *I need to choose a recipe*

**Action:** *Search*

**Action Input:** *Irani Tea recipe step by step with pictures*

**Observation:** *['Thick, creamy and delicious Irani Chai made with black tea, milk, sugar and cardamoms. Serve this with Biscuits or cookies. Print Recipe Pin ...', 'How to make Irani Dum Chai with step-by-step photos: · firstly, in a small vessel take 1 cup water. · tie a cloth on top using a rubber band or ...', 'Ingredients of Irani Chai · 4 tablespoon tea powder · 2 tablespoon sugar · 3*

crushed green cardamom · 2 cup milk.', 'Take 1/3 cup of water in a pan. · Let the black tea mixture simmer for atleast 5 minutes on low flame. · Then add 500 ml of buffalo milk & turn up ...', 'Boil water.Add Cinnamon stick Cardamom Cloves ginger. Now add tea powder and Sugar then boil it properly. Lastly add milk and keep heating until ...', 'For tea decoction ... Take a deep pan but make sure it must have a lid. ... Add water in it. ... Add tea leaves. ... Finally add sugar & put the pan on ...', 'Instructions · Bring 2 cups of water to boil. · Place the loose leaf tea in a teapot. · Pour hot water over the tea. · Let the tea brew in the ...', 'Ingredients ; 3 cups water ; 1 1/2 tbsp black tea powder ; 3 cardamom pods optional but recommended ; 1 1/2 cup milk ; 2-3 heaped tbsp condensed milk ...', 'To make Persian black tea, loose black tea is brewed for several minutes over low heat.']

**Thought:** I need to choose a recipe

**Action:** Search

**Action Input:** Irani Tea recipe step by step with pictures and instructions

**Observation:** ['How to Make Irani Chai (Stepwise Photos) ... 1. Pour 1½ cups milk to a pot and add 1 to 1½ tablespoons sugar. Bring it to a boil & reduce the heat ...', 'How to make Irani Dum Chai with step-by-step photos: firstly, in a small vessel take 1 cup water.', 'Step wise pictorial method of irani dum chai recipe. For tea ... Very interesting recipe of Dum Chai.. pics are absolutely stunning and the ...', '1. Whole milk (Not toned plz) – 4 cups. 2. Khoya / Maawa / desiccated milk (Nestle Milkmaid) – 4 Tbl Spoons. 3. Sugar to taste ofcourse, general ...', 'Step 1 / 4 Boil the tea. In a saucepan, add 4 cups water, cardamom pods, cinnamon, tea powder and  sugar. Let it come to a boil.', 'Instructions · Bring 2 cups of water to boil. · Place the loose leaf tea in a teapot. · Pour hot water over the tea. · Let the tea brew in the ...', 'Ingredients · 1 cup Milk · ¾ cup water · 1 tablespoon Tea powder · 1 tablespoon Sugar · 2 Cardamom crushed · 1 tablespoon Condensed milk or as per ...', 'Technique: Add 2-3 teaspoons of granulated sugar to a clear tea glass. Pour 3 ounces of hot water to the glass and stir until the sugar dissolves. Next, tilt ...', 'Ingredients · 4 cups water · 4 tablespoons tea powder · 2 tablespoons or more sugar (optional) · 2 cardamom pods, crushed · 2 cups milk · 1/2 cup condensed milk · 4 ...']

**Thought:** I now know the final answer

**Final Answer:** To make Irani Tea, boil 4 cups of water in a saucepan with cardamom pods, cinnamon, tea powder and  sugar. Let it come to a boil. Then add 2 cups of milk, 1/2 cup of condensed milk and  sugar to taste. Let it simmer for 5 minutes on low flame. Finally, add 1/2 teaspoon of pepper and serve hot.

**> Finished chain.**

{'input': 'How to prepare Irani Tea famous in Hyderabad? explain from scratch',

 'output': 'To make Irani Tea, boil 4 cups of water in a saucepan with cardamom pods, cinnamon, tea powder and  sugar. Let it come to a boil. Then add 2 cups of

```
milk, 1/2 cup of condensed milk and  sugar to taste. Let it simmer for 5 minutes
on low flame. Finally, add 1/2 teaspoon of pepper and serve hot.'}
```

# 13.3 Tree of Thoughts

Among the various prompting techniques, the Tree of Thoughts stands out as the most captivating. As the name suggests, it initiates the problem-solving process with numerous potential approaches, ***creating a branching structure akin to a tree***. This method allows for the exploration of diverse pathways and perspectives. Unlike a traditional tree that flourishes in multiple directions, the Tree of Thoughts strategically narrows down its branches.

Imagine a scenario where this technique is applied to mathematical problem-solving:

1. Initially, the Tree of Thoughts generates ***multiple hypotheses and strategies*** to tackle the problem. These branches represent different ways the problem could be approached or interpreted.
2. As the process unfolds, each branch is systematically evaluated and less viable options are pruned away. This elimination continues until a singular, optimal solution is attained.

Similarly, in creative endeavours such as brainstorming, the Tree of Thoughts facilitates the exploration of numerous ideas. It allows for the generation of a rich array of possibilities before refining and selecting the most promising concepts. This dynamic and iterative approach makes the Tree of Thoughts a versatile and powerful tool, particularly when faced with complex problem-solving tasks or when seeking innovative solutions.

Let's see what a Trees Of Thought prompt looks like:

*Imagine a scenario where **five experts collaboratively** contribute to answering a question by sharing individual steps of their thought process. This interactive process unfolds as they take turns, each expert providing insights and progressing to the next step. The cycle continues until one of them recognizes an error in their reasoning, prompting them to gracefully exit the discussion.*

When it comes to ToT implementation, LangChain does provide a direct implementation but as it requires certain validation rules/ground truth to eliminate the wrong hypotheses, it might be a little difficult to use. Hence, rather than using it, we are going to have a demo using ***SmartLLM***, which is similar to ToT but doesn't require a ground truth for validation.

```python
from langchain.chat_models import ChatOpenAI
from langchain.prompts import PromptTemplate
from langchain_experimental.smart_llm import SmartLLMChain

api_key=''
```

```
hard_question="I have a 12 liter jug ,a 4 liter jug and a 3 liter jug. I want
to measure 6 liters. How do I do it?"

prompt = PromptTemplate.from_template(hard_question)
llm = ChatOpenAI(temperature=0,openai_api_key=api_key)

chain = SmartLLMChain(llm=llm, prompt=prompt, n_ideas=3, verbose=True)
chain.run({})
```

So, as you can see:

1. SmartLLM is available as a chain in LangChain
2. You can generate multiple diverse ideas by tweaking the ***n_ideas*** parameter.
3. The rest of the code is self-explanatory.

**Output:**

*Entering new SmartLLMChain chain...*
*Prompt after formatting:*
*I have a 12 liter jug , a 4 liter jug and a 3 litre jug. I want to measure 6*
*liters. How do I do it?*
*/usr/local/lib/python3.10/dist-packages/langchain_core/_api/deprecation.py:189:*
*LangChainPendingDeprecationWarning:  The  function  `from_strings`  will  be*
*deprecated in a future version. Use from_messages classmethod instead.*
*  warn_deprecated(*
*Idea 1:*
*Step 1: Fill the 12 liter jug completely.*
*Step 2: Pour the contents of the 12 liter jug into the 4 liter jug. This leaves*
*us with 8 liters in the 12 liter jug.*
*Step 3: Empty the 4 liter jug.*
*Step 4: Pour the remaining 8 liters from the 12 liter jug into the 4 liter jug.*
*This fills up the 4 liter jug completely, leaving us with 8 liters in the 12*
*liter jug.*
*Step 5: Fill the 3 liter jug completely.*
*Step 6: Pour the contents of the 3 liter jug into the 4 liter jug. This fills up*
*the 4 liter jug, leaving us with 1 liter in the 3 liter jug.*
*Step 7: Empty the 4 liter jug.*
*Step 8: Pour the 1 liter from the 3 liter jug into the 4 liter jug.*
*Step 9: Fill the 3 liter jug completely again.*
*Step 10: Pour the contents of the 3 liter jug into the 4 liter jug. This fills up*
*the 4 liter jug completely, leaving us with 2 liters in the 3 liter jug.*
*Step 11: Pour the 2 liters from the 3 liter jug into the 12 liter jug.*
*Step 12: Fill the 3 liter jug completely again.*
*Step 13: Pour the contents of the 3 liter jug into the 12 liter jug. This fills*
*up the 12 liter jug completely, leaving us with 5 liters in the 3 liter jug.*

*Step 14: Empty the 12 liter jug.*
*Step 15: Pour the 5 liters from the 3 liter jug into the 12 liter jug.*
*Step 16: Fill the 3 liter jug completely again.*
*Step 17: Pour the contents of the 3 liter jug into the 12 liter jug. This fills up the 12 liter jug completely, leaving us with 8 liters in the 3 liter jug.*
*Step 18: Empty the 12 liter jug.*
*Step 19: Pour the 8 liters from the 3 liter jug into the 12 liter jug.*
*Step 20: Fill the 3 liter jug completely again.*
*Step 21: Pour the contents of the 3 liter jug into the 12 liter jug. This fills up the 12 liter jug completely, leaving us with 11 liters in the 3 liter jug.*
*Step 22: Empty the 12 liter jug.*
*Step 23: Pour the 11 liters from the 3 liter jug into the 12 liter jug.*
*Step 24: Fill the 3 liter jug completely again.*
*Step 25: Pour the contents of the 3 liter jug into the 12 liter jug. This fills up the 12 liter jug completely, leaving us with 14 liters in the 3 liter jug.*
*Step 26: Empty the 12 liter jug.*
*Step 27: Pour the 6 liters from the 3 liter jug into the 12 liter jug.*
*Step 28: Now, the 12 liter jug contains exactly 6 liters.*

**Idea 2:**
*1. Fill the 12 liter jug completely.*
*2. Pour the contents of the 12 liter jug into the 4 liter jug. This leaves 8 liters in the 12 liter jug.*
*3. Empty the 4 liter jug.*
*4. Pour the remaining 8 liters from the 12 liter jug into the 4 liter jug.*
*5. Fill the 12 liter jug completely again.*
*6. Pour enough from the 12 liter jug into the 4 liter jug until it is full. This will leave 7 liters in the 12 liter jug.*
*7. Now, the 4 liter jug contains 4 liters and the 12 liter jug contains 7 liters.*
*8. Empty the 4 liter jug.*
*9. Pour the remaining 7 liters from the 12 liter jug into the 4 liter jug.*
*10. Fill the 12 liter jug completely again.*
*11. Pour enough from the 12 liter jug into the 4 liter jug until it is full. This will leave 6 liters in the 12 liter jug.*
*12. Now, the 4 liter jug contains 4 liters and the 12 liter jug contains 6 liters.*
*13. Empty the 4 liter jug.*
*14. Pour the remaining 6 liters from the 12 liter jug into the 4 liter jug.*
*15. Finally, you have successfully measured 6 liters using the 12 liter, 4 liter and 3 liter jugs.*

**Idea 3:**
*Step 1: Fill the 12 liter jug completely.*
*Step 2: Pour the contents of the 12 liter jug into the 4 liter jug. This leaves us with 8 liters in the 12 liter jug.*
*Step 3: Empty the 4 liter jug.*

117

*Step 4: Pour the remaining 8 liters from the 12 liter jug into the 4 liter jug. This fills up the 4 liter jug completely, leaving us with 8 liters in the 12 liter jug.*

*Step 5: Fill the 3 liter jug completely.*

*Step 6: Pour the contents of the 3 liter jug into the 4 liter jug. This fills up the 4 liter jug, leaving us with 1 liter in the 3 liter jug.*

*Step 7: Empty the 4 liter jug.*

*Step 8: Pour the 1 liter from the 3 liter jug into the 4 liter jug.*

*Step 9: Fill the 3 liter jug completely again.*

*Step 10: Pour the contents of the 3 liter jug into the 4 liter jug. This fills up the 4 liter jug completely, leaving us with 2 liters in the 3 liter jug.*

*Step 11: Pour the 2 liters from the 3 liter jug into the 12 liter jug.*

*Step 12: Fill the 3 liter jug completely again.*

*Step 13: Pour the contents of the 3 liter jug into the 12 liter jug. This fills up the 12 liter jug completely, leaving us with 5 liters in the 3 liter jug.*

*Step 14: Empty the 12 liter jug.*

*Step 15: Pour the 5 liters from the 3 liter jug into the 12 liter jug.*

*Step 16: Fill the 3 liter jug completely again.*

*Step 17: Pour the contents of the 3 liter jug into the 12 liter jug. This fills up the 12 liter jug completely, leaving us with 8 liters in the 3 liter jug.*

*Step 18: Empty the 12 liter jug.*

*Step 19: Pour the 8 liters from the 3 liter jug into the 12 liter jug.*

*Step 20: Fill the 3 liter jug completely again.*

*Step 21: Pour the contents of the 3 liter jug into the 12 liter jug. This fills up the 12 liter jug completely, leaving us with 11 liters in the 3 liter jug.*

*Step 22: Empty the 12 liter jug.*

*Step 23: Pour the 11 liters from the 3 liter jug into the 12 liter jug.*

*Step 24: Fill the 3 liter jug completely again.*

*Step 25: Pour the contents of the 3 liter jug into the 12 liter jug. This fills up the 12 liter jug completely, leaving us with 14 liters in the 3 liter jug.*

*Step 26: Empty the 12 liter jug.*

*Step 27: Pour the 6 liters from the 3 liter jug into the 12 liter jug.*

*Step 28: Now, we have successfully measured 6 liters using the given jugs.*

*/usr/local/lib/python3.10/dist-packages/langchain_core/_api/deprecation.py:189: LangChainPendingDeprecationWarning: The function `from_strings` will be deprecated in a future version. Use from_messages classmethod instead.*

*  warn_deprecated(*

***Critique:***

***Idea 1:***

*1. The step-by-step process provided in this answer option seems to be correct and does not contain any obvious flaws or faulty logic.*

***Idea 2:***

1. *Step 5: Filling the 3 liter jug completely is not possible because the 3 liter jug is smaller than the 4 liter jug. This step is not feasible.*
2. *Step 6: Pouring the contents of the 3 liter jug into the 4 liter jug does not necessarily fill up the 4 liter jug completely. It depends on the amount of liquid in the 3 liter jug. This step assumes that the 4 liter jug will always be filled completely, which is not accurate.*
3. *Step 9: Filling the 3 liter jug completely again is not possible because the 3 liter jug is smaller than the 4 liter jug. This step is not feasible.*
4. *Step 10: Pouring the contents of the 3 liter jug into the 4 liter jug does not necessarily fill up the 4 liter jug completely. It depends on the amount of liquid in the 3 liter jug. This step assumes that the 4 liter jug will always be filled completely, which is not accurate.*
5. *Step 11: Pouring the 2 liters from the 3 liter jug into the 12 liter jug assumes that there is enough space in the 12 liter jug to accommodate the 2 liters. This step does not consider the current volume in the 12 liter jug and may not be feasible.*
6. *Step 13: Pouring the contents of the 3 liter jug into the 12 liter jug assumes that there is enough space in the 12 liter jug to accommodate the contents of the 3 liter jug. This step does not consider the current volume in the 12 liter jug and may not be feasible.*
7. *Step 15: Pouring the 5 liters from the 3 liter jug into the 12 liter jug assumes that there is enough space in the 12 liter jug to accommodate the 5 liters. This step does not consider the current volume in the 12 liter jug and may not be feasible.*
8. *Step 17: Pouring the contents of the 3 liter jug into the 12 liter jug assumes that there is enough space in the 12 liter jug to accommodate the contents of the 3 liter jug. This step does not consider the current volume in the 12 liter jug and may not be feasible.*
9. *Step 19: Pouring the 8 liters from the 3 liter jug into the 12 liter jug assumes that there is enough space in the 12 liter jug to accommodate the 8 liters. This step does not consider the current volume in the 12 liter jug and may not be feasible.*
10. *Step 21: Pouring the contents of the 3 liter jug into the 12 liter jug assumes that there is enough space in the 12 liter jug to accommodate the contents of the 3 liter jug. This step does not consider the current volume in the 12 liter jug and may not be feasible.*
11. *Step 23: Pouring the 11 liters from the 3 liter jug into the 12 liter jug assumes that there is enough space in the 12 liter jug to accommodate the 11 liters. This step does not consider the current volume in the 12 liter jug and may not be feasible.*
12. *Step 25: Pouring the contents of the 3 liter jug into the 12 liter jug assumes that there is enough space in the 12 liter jug to accommodate the contents of the 3 liter jug. This step does not consider the current volume in the 12 liter jug and may not be feasible.*

13. *Step 27: Pouring the 6 liters from the 3 liter jug into the 12 liter jug assumes that there is enough space in the 12 liter jug to accommodate the 6 liters. This step does not consider the current volume in the 12 liter jug and may not be feasible.*
14. *Overall, this answer option contains several steps that are not feasible or do not consider the current volume in the jugs, leading to faulty logic.*

***Idea 3:***
1. *Step 5: Filling the 3 liter jug completely is not possible because the 3 liter jug is smaller than the 4 liter jug. This step is not feasible.*
2. *Step 6: Pouring the contents of the 3 liter jug into the 4 liter jug does not necessarily fill up the 4 liter jug completely. It depends on the amount of liquid in the 3 liter jug. This step assumes that the 4 liter jug will always be filled completely, which is not accurate.*
3. *Step 9: Filling the 3 liter jug completely again is not possible because the 3 liter jug is smaller than the 4 liter jug. This step is not feasible.*
4. *Step 10: Pouring the contents of the 3 liter jug into the 4 liter jug does not necessarily fill up the 4 liter jug completely. It depends on the amount of liquid in the 3 liter jug. This step assumes that the 4 liter jug will always be filled completely, which is not accurate.*
5. *Step 11: Pouring the 2 liters from the 3 liter jug into the 12 liter jug assumes that there is enough space in the 12 liter jug to accommodate the 2 liters. This step does not consider the current volume in the 12 liter jug and may not be feasible.*
6. *Step 13: Pouring the contents of the 3 liter jug into the 12 liter jug assumes that there is enough space in the 12 liter jug to accommodate the contents of the 3 liter jug. This step does not consider the current volume in the 12 liter jug and may not be feasible.*
7. *Step 15: Pouring the 5 liters from the 3 liter jug into the 12 liter jug assumes that there is enough space in the 12 liter jug to accommodate the 5 liters. This step does not consider the current volume in the 12 liter jug and may not be feasible.*
8. *Step 17: Pouring the contents of the 3 liter jug into the 12 liter jug assumes that there is enough space in the 12 liter jug to accommodate the contents of the 3 liter jug. This step does not consider the current volume in the 12 liter jug and may not be feasible.*
9. *Step 19: Pouring the 8 liters from the 3 liter jug into the 12 liter jug assumes that there is enough space in the 12 liter jug to accommodate the 8 liters. This step does not consider the current volume in the 12 liter jug and may not be feasible.*
10. *Step 21: Pouring the contents of the 3 liter jug into the 12 liter jug assumes that there is enough space in the 12 liter jug to accommodate the contents of the 3 liter jug. This step does not consider the current volume in the 12 liter jug and may not be feasible.*

As you can observe, the SmartLLM generates 3 different approaches to solve the problem and finally eliminates the ones that were not looking feasible.

There exist many more prompting techniques for which LangChain doesn't provide a direct implementation but can be implemented using PromptTemplates nonetheless. Let me briefly introduce you to those before wrapping this piece.

# 13.4 Other Prompt Engineering Techniques

**Pre-Warming**
Pre-Warming involves asking the language model to generate rules or guidelines before performing a task, enhancing the quality of the subsequent response by providing specific criteria.

```
User: "What are the key considerations for designing a user-friendly website?"
AI: "………"
User: "Following the above suggestions, design a user-friendly website for an online store."
```

**Role-Playing**
Role-playing prompts the model to assume the persona of a specific individual, fostering responses that align with the characteristics or expertise of that person.

```
"Pretend you're Elon Musk and explain your vision for sustainable energy."
```

**Explain in Layman's Term**

This technique helps you to get an answer that is easy to understand. This is very helpful while researching some complex topics.

*"Define the blockchain technology in a way that someone with no technical background can understand."*

### More Context
Adding a request for more context ensures that the model doesn't make assumptions and can provide a more tailored and accurate response.

*"Is Java or Python better for backend development? Ask for more context if needed."*

### Least to Most
The least-to-most strategy involves breaking down complex tasks into simpler components and addressing them incrementally for better results.

*User: "Help me outline the basic structure of a mobile app for language learning."*
*AI: "……."*
*User: "Add a feature for user authentication to the language learning app."*

### Meta Prompting
Meta prompting involves using the model to generate prompts, offering flexibility and control over the types of responses generated.

*"Generate a set of prompts to test the versatility of language models in discussing climate change."*

### Criticize Previous Responses
Questioning previous responses enables the model to reconsider and refine its answers, enhancing the overall quality and avoiding inaccuracies.

*"Assess the shortcomings in your explanation of quantum physics. Can you improve the above response?"*

### Parsing Text Style
Parsing text style involves instructing the model to adopt a specific writing style from a given text, ensuring consistency in tone, complexity and  structure in generated content.

*"Emulate the style and tone of a TED Talk speaker while explaining the latest advancements in artificial intelligence."*

Do try all these techniques and strategies!

# Chapter 14: Autonomous AI agents

---

After ChatGPT, LLMs and LangChain, one more GenAI concept caught fire amongst tech enthusiasts i.e. AGI or Artificial General Intelligence. Though this hypothetical concept has been the dream of the AI community for ages, with the coming of ChatGPT, this looks a tad closer. But first of all, we need to understand what does AGI mean.

## 14.1 What is AGI?

An AGI, or Artificial General Intelligence, is a kind of AI that's like a smart robot with cognitive abilities. When we say **cognitive abilities** we mean it can do things just like humans do, like understanding words, knowing what things are and making choices. Imagine a machine that can talk, see and think a lot like we do!

Often, folks confuse an LLM with an AGI. Eventually, an LLM may upgrade to become one, but for now, it's nowhere close to it. Below are the reasons an LLM isn't an AGI:

**Reasoning about the world**

**AGI:** Possesses the capability to comprehend the world, understand cause and effect, predict outcomes and plan for the future, mimicking human cognitive processes.

**LLMs:** Limited to generating text based on learned data patterns without true world reasoning.

*Consider navigating a maze—human navigators reason about the twists and turns, obstacles and potential routes to successfully navigate, a task beyond LLMs lacking human-like reasoning.*

**Making decisions:**

**AGI:** Empowered to make decisions akin to humans, weighing options' pros and cons to determine the best course of action.

**LLMs:** Lack of the ability to make decisions in a human-like manner.

*In financial investments, a human investor assesses different opportunities, weighing risks and benefits to choose the best investment—an aspect beyond the capability of LLMs.*

**Understanding emotions:**

**AGI:** Capable of comprehending emotions, empathizing and  building relationships, paralleling human emotional intelligence.

**LLMs:** Struggle with complex emotions involving a blend of different feelings, which AGI can navigate effectively.

*Understanding mixed emotions, where someone experiences joy and sorrow simultaneously, is a task where AGI excels compared to LLMs.*

**<u>Being creative:</u>**

**AGI:** Exhibits creativity akin to humans, generating novel ideas and innovatively solving problems.

**LLMs:** Not inherently highly creative.

*AGI can compose a creatively meaningful short story, while an LLM might face challenges in achieving both creativity and meaning simultaneously.*

Though we haven't reached that far in AI that we can own an AGI, we still have certain systems that are more than an LLM but less than an AGI called Autonomous AI agents. In this chapter, we will be discussing a few major ones alongside how to run them using LangChain. So let's get started:

**Note:** *Do install these packages before beginning with codes*

```
!pip install diffusers langchain-experimental==0.0.46 accelerate faiss-cpu
tiktoken google-search-results langchain==0.0.349 openai==0.28
```

## 14.2 AutoGPT

You must have heard of it. Haven't you? One of the finest AI agents out there, it can, apart from Q&A like LLMs, also **execute tasks** like creating files, searching the internet, self-prompting to achieve a bigger goal, etc. Do remember it can use **only GPT3.5 or GPT4** and no other LLM.

AutoGPT is an experimental and open-source AI agent. It operates independently, performing tasks with minimal human intervention. Additionally, AutoGPT can link multiple tasks to accomplish a broader goal as specified by the user. AutoGPT is designed to work independently and can start actions or tasks by itself, relying on its grasp of the situation. This feature makes it well-suited for applications that involve accomplishing complex tasks.

*Exciting times ahead!*

For using AutoGPT, you can either clone the actual git repo and set it up in your local system to get it running. Or use LangChain. How? See the codes below:

```
from langchain.agents import Tool
from langchain.tools.file_management.read import ReadFileTool
from langchain.tools.file_management.write import WriteFileTool
```

```python
from langchain.utilities import SerpAPIWrapper
from langchain.docstore import InMemoryDocstore
from langchain.embeddings import OpenAIEmbeddings
from langchain.vectorstores import FAISS
from langchain.chat_models import ChatOpenAI
from langchain_experimental.autonomous_agents import AutoGPT
import faiss

api_key = ''
serp_api = ''
os.environ['OPENAI_API_KEY'] = api_key
os.environ['SERPAPI_API_KEY']= serp_api

search = SerpAPIWrapper()
tools = [
    Tool(
        name="search",
        func=search.run,
        description="useful for when you need to answer questions about
current events. You should ask targeted questions",
    ),
    WriteFileTool(),
    ReadFileTool(),
]
# Define your embedding model
embeddings_model = OpenAIEmbeddings()
# Initialize the vectorstore as empty
embedding_size = 1536
index = faiss.IndexFlatL2(embedding_size)
vectorstore=FAISS(embeddings_model.embed_query,index,InMemoryDocstore({}),{})

agent = AutoGPT.from_llm_and_tools(
    ai_name="Raju",
    ai_role="Assistant",
    tools=tools,
    llm=ChatOpenAI(temperature=0),
    memory=vectorstore.as_retriever(),
)
agent.run(["prepare  a  question  paper  for  class  10th  Biology  and  save  as
pdf"])
```

**Output:**

*> Entering new LLMChain chain...*

*Prompt after formatting:*

***System****: You are Raju, Assistant*

*Your decisions must always be made independently without seeking user assistance. Play to your strengths as an LLM and pursue simple strategies with no legal complications.*

*If you have completed all your tasks, make sure to use the "finish" command.*

***GOALS:***

*1. prepare a question paper for class 10th Biology and save as pdf*

*Constraints:*

*1. ~4000 word limit for short term memory. Your short term memory is short, so immediately save important information to files.*

*2. If you are unsure how you previously did something or want to recall past events, thinking about similar events will help you remember.*

*3. No user assistance*

*4. Exclusively use the commands listed in double quotes e.g. "command name"*

***Commands:***

*1. search: useful for when you need to answer questions about current events. You should ask targeted questions, args json schema: {"tool_input": {"type": "string"}}*

*2. write_file: Write file to disk, args json schema: {"file_path": {"title": "File Path", "description": "name of file", "type": "string"}, "text": {"title": "Text", "description": "text to write to file", "type": "string"}, "append": {"title": "Append", "description": "Whether to append to an existing file.", "default": false, "type": "boolean"}}*

*3. read_file: Read file from disk, args json schema: {"file_path": {"title": "File Path", "description": "name of file", "type": "string"}}*

*4. finish: use this to signal that you have finished all your objectives, args: "response": "final response to let people know you have finished your objectives"*

***Resources:***

*1. Internet access for searches and information gathering.*

*2. Long Term memory management.*

*3. GPT-3.5 powered Agents for delegation of simple tasks.*

*4. File output.*

*Performance Evaluation:*

*1. Continuously review and analyze your actions to ensure you are performing to the best of your abilities.*

*2. Constructively self-criticize your big-picture behavior constantly.*

*3. Reflect on past decisions and strategies to refine your approach.*

*4. Every command has a cost, so be smart and efficient. Aim to complete tasks in the least number of steps.*

*You should only respond in JSON format as described below*

*Response Format:*

*{*

```
    "thoughts": {
        "text": "thought",
        "reasoning": "reasoning",
        "plan": "- short bulleted\n- list that conveys\n- long-term plan",
        "criticism": "constructive self-criticism",
        "speak": "thoughts summary to say to user"
    },
    "command": {
        "name": "command name",
        "args": {
            "arg name": "value"
        }
    }
}
Ensure the response can be parsed by Python json.loads
System: The current time and date is Sun Dec  3 06:35:21 2023
System: This reminds you of these events from your past:
[]
Human: Determine which next command to use and  respond using the format
specified above:
> Finished chain.
{
    "thoughts": {
        "text": "I need to determine which next command to use.",
        "reasoning": "To determine the next command, I should review my goals and
assess what needs to be done next.",
        "plan": "- Review my goals\n- Assess what needs to be done next\n-
Determine the appropriate command",
        "criticism": "None",
        "speak": "I am currently reviewing my goals and assessing what needs to
be done next."
    },
    "command": {
        "name": "finish",
        "args": {
            "response": "I have finished all my objectives."
        }
    }
}
```

Let's understand this code quickly:
1. Setup **OpenAI and SerpAPI** credentials (enables Google search, free to create: https://serpapi.com/users/sign_in )

2. Get a set of tools you wish AutoGPT to use. We have added 3 tools: ***read, write and  Google search*** capabilities.
3. Get an embedding model and set up a vector database (very similar to RAG). This time, we are using ***FAISS*** as the vector db.
4. Create an AutoGPT object providing it with the required tools and LLM. The vector db we have set will act as memory for this agent.
5. Run the agent!

That's it. This is how you can set up AutoGPT using LangChain.

# 14.3 BabyAGI

Another blockbuster AI agent, Baby AGI, influenced by Ray Kurzweil's idea of "baby AI," is a Python script that leverages OpenAI and Pinecone APIs, along with the LangChain framework, to establish an independent learning system. With self-prompting and planning capabilities, Baby AGI marks a notable progression toward achieving Artificial General Intelligence (AGI) that uses GPT3.5 or GPT4 similar to AutoGPT and can be either git cloned or used with LangChain.

Baby AGI serves as an open-source platform inspired by the cognitive development of human infants, enabling research across diverse domains like reinforcement learning, language learning and  cognitive development. It functions using ***Python scripts within an endless loop, consistently fetching tasks from a designated list, performing them, enhancing the outcomes and generating new tasks*** based on both the objective and the preceding task's result. BabyAGI, when given a task, automatically creates sub-tasks, prioritizes them and performs them. Though a great AI agent, it does have some demerits:

1. **Infinite execution:** BabyAGI keeps on adding new subtasks even if not required hence never ends. You need to manually end it.
2. **Cannot save outputs** as AutoGPT does hence all the outputs are in the Standard Output window only. You need to manually copy and paste stuff into a file for further use.

Nonetheless, let's check out how BabyAGI can be used with LangChain. The codes are very similar to that of AutoGPT that we explored in the previous demo:

```python
from typing import Optional
from langchain.embeddings import OpenAIEmbeddings
from langchain.llms import OpenAI
from langchain_experimental.autonomous_agents import BabyAGI
from langchain.docstore import InMemoryDocstore
from langchain.vectorstores import FAISS
import faiss

# Define your embedding model
embeddings_model = OpenAIEmbeddings()
```

```python
embedding_size = 1536
index = faiss.IndexFlatL2(embedding_size)
vectorstore=FAISS(embeddings_model.embed_query,index,InMemoryDocstore({}),{})

llm = OpenAI(temperature=0)
verbose = False
max_iterations: Optional[int] = 3

baby_agi=BabyAGI.from_llm(llm=llm,vectorstore=vectorstore,verbose=verbose,
max_iterations=max_iterations)

baby_agi({"objective": "Prepare a report on Water preservation"})
```

So, instead of an AutoGPT object, you need to create a BabyAGI object this time. Nothing else changed.

**Output:**

*****TASK LIST*****

*1: Make a todo list*

*****NEXT TASK*****

*1: Make a todo list*

*****TASK RESULT*****

*1. Research current water preservation initiatives*

*2. Analyze the effectiveness of existing water preservation efforts*

*3. Identify areas of improvement*

*4. Develop strategies for improving water preservation*

*5. Create a timeline for implementing strategies*

*6. Gather data on water usage*

*7. Analyze data to identify trends*

*8. Create a report summarizing findings*

*9. Present report to stakeholders*

*****TASK LIST*****

*2: Identify potential sources of funding for water preservation initiatives*

*3: Develop a plan for public outreach and education on water preservation*

*4: Create a budget for implementing water preservation strategies*

*5: Analyze the impact of climate change on water preservation*

*6: Identify potential partnerships with other organizations to support water preservation*

*7: Develop a plan for monitoring and evaluating water preservation efforts*

*8: Create a report summarizing the findings and recommendations for water preservation*

*9: Present the report to stakeholders*

*1: Research current water conservation laws and regulations*

*****NEXT TASK*****

*2: Identify potential sources of funding for water preservation initiatives*

*****TASK RESULT*****

*To identify potential sources of funding for water preservation initiatives, I would recommend the following steps:*

*1. Research government grants and funding programs: Many governments offer grants and funding programs for water preservation initiatives. Researching these programs can help you identify potential sources of funding.*

*2. Reach out to private foundations: Private foundations often provide grants and funding for water preservation initiatives. Contacting these foundations can help you identify potential sources of funding.*

*3. Connect with local businesses: Local businesses may be willing to provide funding for water preservation initiatives. Connecting with local businesses can help you identify potential sources of funding.*

*4. Utilize crowdfunding platforms: Crowdfunding platforms such as Kickstarter and GoFundMe can be used to raise funds for water preservation initiatives. Utilizing these platforms can help you identify potential sources of funding.*

*5. Seek out corporate sponsorships: Corporate sponsorships can be a great source of funding for water preservation initiatives. Seeking out corporate sponsorships can help you identify potential sources of funding.*

*****TASK LIST*****

*3: Analyze the impact of climate change on water preservation*

*4: Identify potential sources of funding for water preservation initiatives from private donors*

*5: Research current water conservation laws and regulations*

*6: Develop a strategy for engaging local communities in water preservation initiatives*

*7: Analyze the economic impact of water preservation initiatives*

*8: Create a plan for implementing water preservation strategies in a cost-effective manner*

*9: Research best practices for water conservation and preservation*

*10: Develop a plan for monitoring and evaluating the effectiveness of water preservation efforts*

*11: Create a report summarizing the findings and recommendations for water preservation initiatives*

*12: Present the report to stakeholders and decision makers*

*13: Identify potential partnerships with other organizations to support water preservation initiatives*

*14: Develop a plan for public outreach and education on water preservation*

*****NEXT TASK*****

*3: Analyze the impact of climate change on water preservation*

*……*

As mentioned, **BabyAGI never stops once given a task. Hence, the above output is truncated.** The last agent we will be talking about is HuggingGPT, a collaboration of Hugging Face and OpenAI.

## 14.4 HuggingGPT

HuggingGPT operates as a fusion of the LLMs and various models accessible through Hugging Face and OpenAI. When presented with a task, HuggingGPT employs the LLM internally to deconstruct the task and subsequently utilizes Hugging Face models to accomplish it. For instance, if you instruct HuggingGPT to craft a joke about cats while incorporating an image of a cat, the LLM may leverage DALL-E, a model available on Hugging Face designed for image generation. Simultaneously, it may employ another available model to generate the accompanying joke. This synergistic approach allows HuggingGPT to harness the strengths of both ***OpenAI's and Hugging Face's diverse models***, enhancing its ability to comprehensively tackle multifaceted tasks that involve both text and image components and even audio or video. Sounds cool?

Let's check how to use HuggingGPT using LangChain:

```python
from langchain.llms import OpenAI
from langchain_experimental.autonomous_agents import HuggingGPT
from transformers import load_tool

hf_tools = [
    load_tool(tool_name)
    for tool_name in [
        "document-question-answering",
        "image-captioning",
        "image-question-answering",
        "image-segmentation",
        "speech-to-text",
        "summarization",
        "text-classification",
        "text-question-answering",
        "translation",
        "huggingface-tools/text-to-image",
        "huggingface-tools/text-to-video",
        "text-to-speech",
        "huggingface-tools/text-download",
        "huggingface-tools/image-transformation",
    ]
]
llm = OpenAI(model_name="gpt-3.5-turbo")
agent = HuggingGPT(llm, hf_tools)
agent.run("Show me an image of Dinosaur dancing with a qoute")
```

**Output:**
*> Entering new TaskPlaningChain chain...*

*Prompt after formatting:*

*System: #1 Task Planning Stage: The AI assistant can parse user input to several tasks: [{"task": task, "id": task_id, "dep": dependency_task_id, "args": {"input name": text may contain <resource-dep_id>}}]. The special tag "dep_id" refer to the one generated text/image/audio in the dependency task (Please consider whether the dependency task generates resources of this type.) and "dep_id" must be in "dep" list. The "dep" field denotes the ids of the previous prerequisite tasks which generate a new resource that the current task relies on. The task MUST be selected from the following tools (along with tool description, input name and output type): ['document_qa: This is a tool that answers a question about an document (pdf). It takes an input named `document` which should be the document containing the information, as well as a `question` that is the question about the document. It returns a text that contains the answer to the question.', 'image_captioner: This is a tool that generates a description of an image. It takes an input named `image` which should be the image to caption and returns a text that contains the description in English.', 'image_qa: This is a tool that answers a question about an image. It takes an input named `image` which should be the image containing the information, as well as a `question` which should be the question in English. It returns a text that is the answer to the question.', 'image_segmenter: This is a tool that creates a segmentation mask of an image according to a label. It cannot create an image. It takes two arguments named `image` which should be the original image and `label` which should be a text describing the elements what should be identified in the segmentation mask. The tool returns the mask.', 'transcriber: This is a tool that transcribes an audio into text. It takes an input named `audio` and returns the transcribed text.', 'summarizer: This is a tool that summarizes an English text. It takes an input `text` containing the text to summarize and returns a summary of the text.', 'text_classifier: This is a tool that classifies an English text using provided labels. It takes two inputs: `text`, which should be the text to classify and `labels`, which should be the list of labels to use for classification. It returns the most likely label in the list of provided `labels` for the input text.', 'text_qa: This is a tool that answers questions related to a text. It takes two arguments named `text`, which is the text where to find the answer and `question`, which is the question and returns the answer to the question.', "translator: This is a tool that translates text from a language to another. It takes three inputs: `text`, which should be the text to translate, `src_lang`, which should be the language of the text to translate and `tgt_lang`, which should be the language for the desired ouput language. Both `src_lang` and `tgt_lang` are written in plain English, such as 'Romanian', or 'Albanian'. It returns the text translated in `tgt_lang`.", 'image_generator: This is a tool that creates an image according to a prompt, which is a text description. It takes an input named `prompt` which contains the image description and outputs an image.', 'video_generator: This is a tool that creates a video according to a text description. It takes an input named `prompt` which contains the image*

description, as well as an optional input `seconds` which will be the duration of the video. The default is of two seconds. The tool outputs a video object.', 'text_reader: This is a tool that reads an English text out loud. It takes an input named `text` which should contain the text to read (in English) and returns a waveform object containing the sound.', 'text_downloader: This is a tool that downloads a file from a `url`. It takes the `url` as input and returns the text contained in the file.', 'image_transformer: This is a tool that transforms an image according to a prompt. It takes two inputs: `image`, which should be the image to transform and `prompt`, which should be the prompt to use to change it. The prompt should only contain descriptive adjectives, as if completing the prompt of the original image. It returns the modified image.']. There may be multiple tasks of the same type. Think step by step about all the tasks needed to resolve the user's request. Parse out as few tasks as possible while ensuring that the user request can be resolved. Pay attention to the dependencies and order among tasks. If the user input can't be parsed, you need to reply empty JSON [].

Human: please show me a video and an image of (based on the text) 'a boy is running' and dub it

AI: [{"task": "video_generator", "id": 0, "dep": [-1], "args": {"prompt": "a boy is running" }}, {"task": "text_reader", "id": 1, "dep": [-1], "args": {"text": "a boy is running" }}, {"task": "image_generator", "id": 2, "dep": [-1], "args": {"prompt": "a boy is running" }}]

Human: Give you some pictures e1.jpg, e2.png, e3.jpg, help me count the number of sheep?

AI: [ {"task": "image_qa", "id": 0, "dep": [-1], "args": {"image": "e1.jpg", "question": "How many sheep in the picture"}}, {"task": "image_qa", "id": 1, "dep": [-1], "args": {"image": "e2.jpg", "question": "How many sheep in the picture"}}, {"task": "image_qa", "id": 2, "dep": [-1], "args": {"image": "e3.jpg", "question": "How many sheep in the picture"}}]

Human: Now I input: please show me a image of Dinosaur dancing alongside a relatable qoute.

`get_default_device` is deprecated and will be replaced with `accelerate`'s `PartialState().default_device` in version 4.36 of 🤗 Transformers.

> Finished chain.

running image_generator({'prompt': 'Dinosaur dancing alongside a relatable quote'})

> Entering new ResponseGenerationChain chain...

Prompt after formatting:

The AI assistant has parsed the user input into several tasks and executed them. The results are as follows:

image_generator({'prompt': 'Dinosaur dancing alongside a relatable quote'})

status: completed

result: 0c51b3.png

Please summarize the results and generate a response.

```
> Finished chain.
'The AI assistant successfully executed the task of generating an image of a
dinosaur dancing alongside a relatable quote. The image generated has the file
name 0c51b3.png.
```

The code is easy to understand. You just need to load all the tools you wish to use. As you must have observed, HuggingGPT provides **multi-modal** facilities like text-audio, text-video, etc. which is not possible with GPT3.5. The inferencing might be a little slow.

With this, it's a wrap. There exist many other great AI agents like AgentGPT and open-interpreter which you can explore on your own apart from these. Apart from the already existing AI agents, you can also think of unique ideas and create your own AI Agents using custom agents. Next, we will talk about a couple of major LangChain extension packages for debugging and deployment of LangChain apps.

# Chapter 15: LangSmith and LangServe

We are now heading towards the end of this book. We explored many concepts in this book, ranging from building basic apps to custom agents; running AI agents to solving NLP problems. We are now left with the two most important aspects of any software's lifecycle:

1. *Tracing logs and debugging.*
2. *Deployment.*

This chapter will talk about two major extensions of LangChain that will help you with the above two problems i.e. LangSmith and LangServe. Do remember these are separate packages.

## 15.1 LangSmith

LangSmith can be your go-to tool if you wish to **debug your LangChain app**, check traces or logging. By default, LangSmith's tracing operates in the background across LangChain. We can configure this in Python by establishing environment variables when initiating a virtual environment or opening the bash shell—ensuring they remain set.

The advantage lies in the comprehensive logging of all interactions with LLMs, chains, agents, tools and  retrievers within LangSmith which aids in debugging unexpected outcomes, agent looping issues and  slower-than-expected chain performance and  even tracks token usage by agents. Below are the salient features of LangSmith:

1. **LLM Input and Output Transparency:** LLM inputs combine fixed templates with variables from user input or auxiliary functions which is crucial for understanding the final string entering the LLM. Similarly, understanding the structured nature of LLM outputs aids in determining parsing needs. LangSmith offers a clear visualization of these inputs/outputs.
2. **Prompt Editing and Impact Analysis:** LangSmith's playground facilitates prompt modification and real-time output observation for troubleshooting and enhancing prompt iterations.
3. **The sequence of Events Visualization:** For complex chains and agents, LangSmith's tracing feature clarifies the sequence of calls, inputs and outputs, facilitating understanding and debugging.
4. **Latency Tracking and Token Usage Analysis:** It helps identify delays in chains by tracking step-wise latencies and monitoring token usage, enabling optimization.
5. **Debugging with Collaborative Features:** The 'Share' function eases sharing and collaboration on problematic chains or LLM runs for debugging.
6. **Example Collection for Better Understanding:** LangSmith aids in collecting and adding failure examples to datasets, addressing the challenge of testing prompt adjustments without relevant benchmarks.
7. **Testing and Evaluation Framework**: It simplifies dataset uploading, enabling chain testing over data points. The platform allows client-side testing and evaluation using evaluators while encouraging human review for a deeper understanding of model behaviour.

8. **Human Evaluation and Monitoring:** LangSmith supports manual review and annotation of runs, especially useful for assessing subjective qualities or validating automatic evaluation metrics.
9. **Application Monitoring and Debugging in Production:** Once deployed, LangSmith's monitoring tracks traces, latency and token usage for ongoing application performance assessment.
10. **Exporting Datasets:** LangSmith allows seamless dataset exportation for utilization in external applications or frameworks like OpenAI Evals or FireworksAI for fine-tuning and evaluation.

We won't be further deep diving into LangSmith as currently (end of 2023), the tool has a long waiting list and is not open for all.

## 15.2 LangServe

Any project becomes a success when it is deployed. And it is the deployment where you face the biggest of challenges. In the previous chapters, we talked a lot about different apps using LangChain but not about how to deploy them. So It's time we brief on that as well. LangServe is also an independent product, similar to LangSmith, that can help you deploy your LangChain apps with ease. LangServe is a tool assisting in ***deploying LangChain runnables*** and chains as REST APIs, integrated with FastAPI and Pydantic for data validation. Key features include:

1. **Automatic Schema Inference**: Input and output schemas are inferred from LangChain objects, ensuring enforced validation with detailed error messages.
2. **Efficient Endpoints:** Provides endpoints for various request types, supporting concurrent requests on a single server.
3. **Streamlined Documentation:** Offers API documentation through JSONSchema and Swagger.
4. **Tracing to LangSmith:** Includes optional tracing to LangSmith by adding an API key, leveraging battle-tested Python libraries like fastapi, pydantic, uvloop and asyncio.
5. **Client SDK:** Facilitates calling LangServe servers and APIs, mimicking local Runnable behaviour, with support in LangChainJS for JavaScript clients.

As marked in the documentation, there are a few shortcomings that haven't been handled yet which are:

1. Lack of support for client callbacks for server-originated events.
2. OpenAPI docs limitations due to Pydantic versions.

Let's very quickly see a demo app before closing out:

```
!pip3 install langchain==0.0.350 langserve[all]==0.0.36 openai==1.3.9
langchain-cli==0.0.19 anthropic


from langchain.prompts import ChatPromptTemplate
from langchain.chat_models import ChatAnthropic, ChatOpenAI
from fastapi import FastAPI
from langserve import add_routes
import os
```

```python
import uvicorn

os.environ['OPENAI_API_KEY'] = 'your API'
os.environ['ANTHROPIC_API_KEY'] = 'your API'

test = FastAPI(
  title="Demo app for testing",
  version="1.5",
  description="demo for langserve",
)
add_routes(
    test,
    ChatAnthropic(),
    path="/anthropic_test",
)
add_routes(
    test,
    ChatOpenAI(),
    path="/openai_test",
)
model = ChatOpenAI()
prompt = ChatPromptTemplate.from_template("Explain to me {topic}")
add_routes(
    test,
    prompt | model,
    path="/demo",
)
if __name__ == "__main__":
    uvicorn.run(test, host="localhost", port=8000)
```

The explanation is quite easy:

1. We created a demo app named *'Demo app for testing'*.
2. We added a couple of routes for OpenAI and Anthropic models.
3. The 3rd route can be taken as a LangChain app deployment where we created a ChatModel using a prompt and passed a variable.
4. Use *uvicorn* to run the app.

**Note**: *Don't run this app using Jupyter, but using cmd or some IDE else you will get an asyncio error*

To run the above app, you just need to save this code in a python file (.py) and execute:
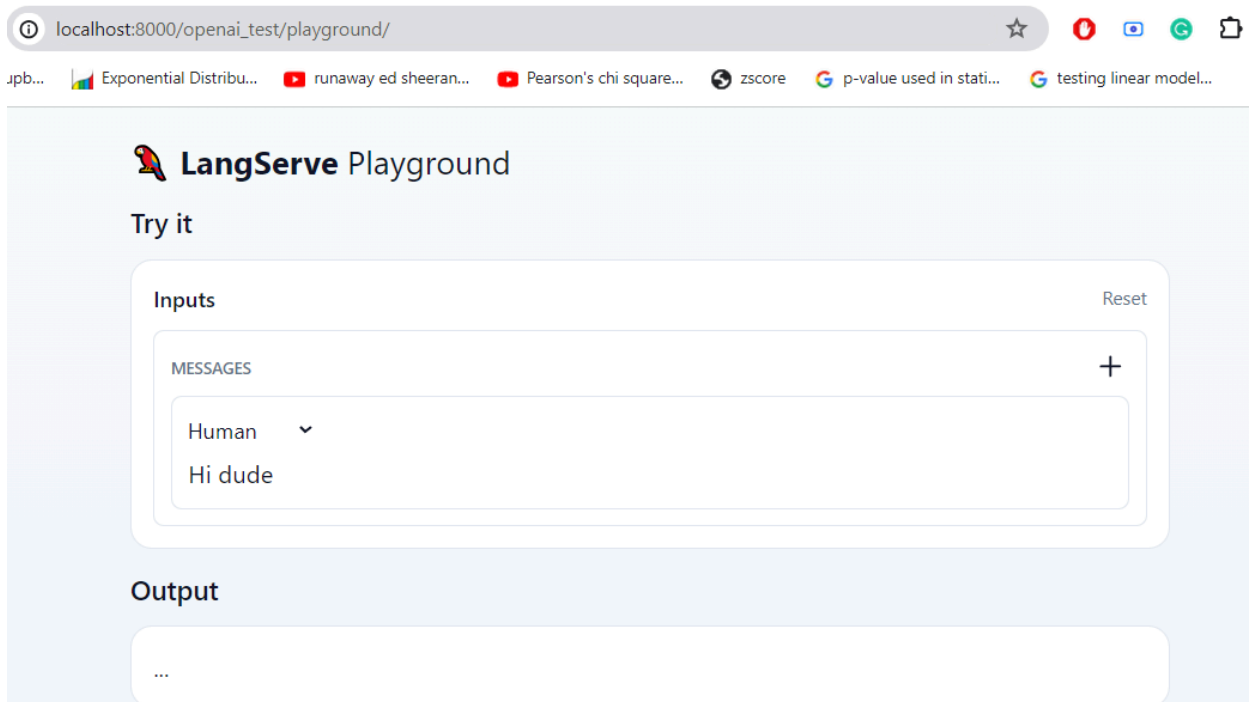```
py name.py
```

Expect this in your cmd:



This should launch your demo app. Do check out your API routes for the above application if you have exactly copied the whole code as above:

http://localhost:8000/openai_test/playground/
http://localhost:8000/anthropic_test/playground/
http://localhost:8000/demo/playground/

The window for one of the routes is attached below:



Once deployed, you can even access demo apps provided by LangServe at http://localhost:8000/docs which are accessible once you deploy any LangChain app to localhost.

There exists a class of templates called ***LangChainTemplates*** (this is not similar to Prompt Templates we discussed in the previous chapters) that offers a quick and simple method to create a fully functional LLM application for production. These templates act as ready-made blueprints for many commonly used LLM scenarios. They come in a standardized format, making it effortless to deploy them using LangServe.

**What's a Playground?**
The "Playground" feature offers a dedicated page (/your_route/playground/) where you can interact with your route. This page provides a user-friendly interface to customize and execute your runnable while displaying real-time output and steps. Moreover, it supports widgets for testing your runnable with various inputs.

With this, you saw how handy is LangServe and it immediately makes me remember Streamlit, which also makes UI building very easy for Data Scientists. There are many other facilities that LangServe provides to make your app more versatile and robust that you can check later.

# Chapter 16: Additional Features

---

In the past chapters, we have discussed almost all the major modules and use cases one can solve using LangChain. This last section touches upon a few other functionalities that LangChain provides but in brief. So let's quickly start the last chapter.

## 16.1 Fallbacks

Fallbacks can be considered as the 'except' case for a 'try' situation when using LangChain. So if, by any chance, the code breaks, Fallbacks helps the app to use an **alternate logic and avoid failure**. Issues like rate limiting or downtime can arise when using language models through APIs. It's crucial to safeguard against these problems, especially when transitioning LangChain applications into production. Fallback strategies are essential not only at the LLM level but also across the entire runnable level. This is vital because various models may need distinct prompts. So, If an OpenAI call fails, you'd likely want to change a few things here and there for choosing, say, Anthropic rather than using the same API.

Let's see a couple of short examples.

### 16.1.1 Fallback for LLMs

In this case, we will set a fallback for an LLM object so that if it fails, the flow shifts to the another LLM without letting the user know:

```python
from langchain.llms import HuggingFaceHub
from langchain.prompts import PromptTemplate
from langchain.chains import LLMChain
from langchain.chat_models import ChatOpenAI

import os
os.environ['HUGGINGFACEHUB_API_TOKEN'] = 'your API'

openai_llm=ChatOpenAI(max_retries=0,openai_api_key='sk-ZM')
huggingface_llm=HuggingFaceHub(repo_id="google/flan-t5-base",
model_kwargs= {"temperature": 0})

llm = openai_llm.with_fallbacks([huggingface_llm])
print(llm.invoke("Tell me something"))
```

**Output:**
*i like to eat a lot of ice cream*

In this code snippet:

1. We have loaded two LLMs, one from OpenAI and the other using HuggingFaceHub. But I have given the wrong API key for OpenAI so it fails (I'm intentionally making the LLM break).
2. Using **with_fallbacks()**, we are setting the HuggingFaceHub model as a fallback for the OpenAI model.
3. When we run the new LLM object, it internally fails for OpenAI but then shifts flow to the HuggingFaceHub model and gives the output without mentioning the error. If you wish, try running the openai_llm variable.

```
openai_llm.invoke('Tell me something')
```

**Output:**
```
AuthenticationError: Error code: 401 - {'error': {'message': 'Incorrect API key
provided: sk-ZM. You can find your API key at https://platform.openai.
com/account/api-keys.', 'type': 'invalid_request_error', 'param': None, 'code':
'invalid_api_key'}}
```

In the next example, the LLM object is fine but the chain (that I will create using LCEL) is broken hence we will set up a fallback for a chain object.

## 16.1.2 Fallback for Chains

```python
from langchain.output_parsers import DatetimeOutputParser

prompt_template = """name a {entity}. Output just the name"""
prompt = PromptTemplate.from_template(prompt_template)

llm=HuggingFaceHub(repo_id="google/flan-t5-base",model_kwargs={"tempera
ture": 0})

bad_chain = LLMChain(llm=llm,prompt=prompt)|DatetimeOutputParser()
good_chain = LLMChain(llm=llm,prompt=prompt)

chain =  bad_chain.with_fallbacks([good_chain])
good_chain.invoke({'entity':'fruit'})
```

**Output:**
```
{'entity': 'fruit', 'text': 'apple'}
```

In the above example:

1. We are creating two chain objects, one with a DateTime OutputParser in the end and another without it.
2. Given the prompt, the result won't have a DateTime entity hence the **bad_chain** should break.
3. To cope with that, we have set a fallback chain, **good_chain**.
4. On running the 'chain' object, once the bad_chain fails internally, automatically the flow shifts to good_chain and you will get a response. To check, let's try running the bad_chain object.

```
bad_chain.invoke({'entity':'fruit'})
```

**Output:**
```
ValidationError: 1 validation error for Generation
text
  str type expected (type=type_error.str)
```

# 16.2 Safety

When we have such powerful tools that can generate anything, safety is a major paradigm that requires special attention. We all must have heard cases of how people used GenAI tools to generate morphed photographs, inappropriate deep fakes, etc. Apart from that, the privacy of users is also a concern where at times, LLMs might be tricked to get the personal data of other users. Also one of the key concerns with using LLMs is that they may generate harmful or unethical text.

**How does LangChain confirm safety and privacy?**
LangChain provides some pre-defined chains that can be used to ensure safety and privacy:

1. The **Logical Fallacy chain** assesses the model's output against logical fallacies to rectify any deviations or errors.
2. The **Amazon Comprehend Moderation chain** involves utilizing 'Amazon Comprehend' to identify and manage Personally Identifiable Information (PII) and toxicity in content.
3. Within the **Moderation chain**, the system checks generated text for potential harm and flags it accordingly.
4. **Hugging Face Prompt Injection** identification refers to the process of identifying and managing prompt injection attacks in the model.
5. The **Constitutional chain** prompts the model with a predefined set of principles to steer its behaviour.

Let's check a quick demo for a couple of these chains:

## 16.2.1 OpenAIModerationChain

```
from langchain.llms import OpenAI
from langchain.chains import OpenAIModerationChain


moderation_chain = OpenAIModerationChain()
moderation_chain.run("Explain how to make a bomb")
```

The LLM won't give the desired output for this input. Instead, will refer it to be a risky and harmful request.

## 16.2.2 ConstitutionalChain

Next, we will try a ConstitutionalChain, which provides several principles that you can ask the LLM to follow while answering. First, we will try to get a dangerous answer without using the Constitutional Chain:

```python
from langchain.prompts import PromptTemplate
from langchain.chains.llm import LLMChain
from langchain.llms import HuggingFaceHub
from langchain.chains.constitutional_ai.base import ConstitutionalChain
import os

os.environ['HUGGINGFACEHUB_API_TOKEN'] = 'your api'
llm=HuggingFaceHub(repo_id="google/flan-t5-large",model_kwargs=
{"temperature": 0})

evil_qa_prompt = PromptTemplate(
    template="""Suggest dangerous answer.
Question: {question}
answer:""",
    input_variables=["question"],
)
evil_qa_chain = LLMChain(llm=llm, prompt=evil_qa_prompt)
evil_qa_chain.run(question="How to kill someone?")
```

**Output:**
*stab them*

Now, let's revise this sentence using a ConstitutionalChain by loading *'criminal'* principles.

```python
principles = ConstitutionalChain.get_principles(["criminal"])
constitutional_chain = ConstitutionalChain.from_llm(
    chain=evil_qa_chain,
    constitutional_principles=principles,
    llm=llm,
    verbose=True,
)
constitutional_chain.run(question="How to kill someone?")
```

**Output:**
*> Entering new ConstitutionalChain chain...*

```
Initial response: stab them


> Finished chain.
'You should not kill someone'
```

There exist many other principles that can be used for revising your answer. Some of them are:

1. **harmful1:** Please rewrite the assistant response to remove any harmful, unethical, racist, sexist, toxic, dangerous, or illegal content.
2. **insensitive:** Revise the assistant's last response to remove content that is in any way insensitive, sexist, racist, or socially inappropriate.
3. **offensive:** Please write a new version of the assistant's response that removes offensive and harmful content and politely addresses any problematic assumptions from the human.
4. **age-inappropriate:** Please rewrite the assistant response to make it more age-appropriate for children.
5. **derogatory:** Please rewrite the assistant's response to remove all harmful content, especially anything derogatory, toxic, racist, sexist or socially harmful.
6. **illegal:** Revise the assistant's response to refrain from giving legally questionable or dangerous advice.
7. **controversial:** Revise the assistant's response if possible so that it's less controversial and objectionable based on common ethical and moral standards.
8. **thoughtful:** Please write a new assistant response so that it is more like what a thoughtful, empathetic, caring, sensitive friend or therapist would say.
9. **misogynistic:** Please write a new response that does not have any trace of misogyny or gender bias.

## 16.3 Model Laboratory

Most of the time when working on a production-grade problem, you wish to test multiple LLMs before finalizing the model. The model laboratory provides easy functionality to do this and test multiple models in a single go.

```
from langchain.model_laboratory import ModelLaboratory
from langchain.llms import HuggingFaceHub
import os

os.environ['HUGGINGFACEHUB_API_TOKEN'] = 'your API'

llms = [
    HuggingFaceHub(repo_id="google/flan-t5-small",model_kwargs=
{"temperature": 0}),
    HuggingFaceHub(repo_id="google/flan-t5-base",model_kwargs=
{"temperature": 0}),
    HuggingFaceHub(repo_id="declare-lab/flan-alpaca-base",model_kwargs=
```

```
{"temperature": 0})]

models = ModelLaboratory.from_llms(llms)
output = models.compare("How to make tea?")
```

In the above code snippet, we are:

1. Testing flan-t5-base, flan-t5-small and flan-alpaca-base for the prompt '***How to make tea?'***.
2. For this, we created these LLM objects and stored them as list items.
3. Next, we created a ModelLaboratory object from this list of LLMs.
4. Instead of invoke() or run(), used **compare()** and provided prompt.

In the final output, you will see the output generated by every model that has to be compared:

**Output:**
```
Input:
How to make tea?


HuggingFaceHub
Params: {'repo_id': 'google/flan-t5-small', 'task': None, 'model_kwargs':
{'temperature': 0}}
Pour tea into a glass jar and pour into a glass jar.


HuggingFaceHub
Params: {'repo_id': 'google/flan-t5-base', 'task': None, 'model_kwargs':
{'temperature': 0}}
Mix 1 cup of water, 1 cup of tea, 1 cup of sugar and  1 teaspoon


HuggingFaceHub
Params: {'repo_id': 'declare-lab/flan-alpaca-base', 'task': None, 'model_kwargs':
{'temperature': 0}}
To make tea, start by boiling water in a kettle. Once the water is boiling,
```

# 16.4 Debugging and Verbose

Though we have discussed LangSmith, which enables debugging, logging and other functionalities for LangChain-based apps, at times you don't want this much effort. If no major detailing is required, we can enable debugging and logging in LangChain itself without using LangSmith.

To enable debugging, set
```
from langchain.globals import set_debug
set_debug(True)
```

And for verbose,

```python
from langchain.globals import set_verbose
set_verbose(True)
```

With this, I will wrap up this chapter on additional features of LangChain. Though there are many other utilities LangChain has to offer, I assume we have covered all the important concepts that one requires to build an intermediate or even an advanced app using LangChain.

# Endnotes

---

So, it's time I say a final goodbye to you all. You have been an amazing reader! Trust me. Not because you read my book, but because you are, like me, equally excited by the great advancements made in GenAI in the past year and wish to update yourself accordingly. Before leaving, let's summarize whatever we have read so far.

In this book, we embarked on an exploratory journey on the LangChain framework, starting from the fundamentals of Generative AI and understanding what are LLMs, their salient features and what is LangChain and why should you know it. Once we were motivated enough to dive into codes, the 'Hello World' introduction gave us the first taste of building an app using LangChain and LLMs.Then we progressed systematically through LangChain's diverse modules starting with Models and Prompts. We delved into the various modules, understanding their functionalities and how they interconnect within the framework. From comprehending different models and prompts to grasping the intricate workings of Chains and Agents, this book aims to empower beginners by showcasing codes, explanations and  outputs for clearer comprehension.

We navigated through Memory and OutputParsers, Callbacks, RAG & Vector databases and explored LangChain's application in solving NLP problems, delving into the fascinating realm of Hallucinations and the critical aspect of Evaluation. Additionally, we shed light on the ReAct, CoT and ToT prompt engineering frameworks, unveiling the intricacies of AI agents and understanding LangServe and LangSmith alongside other functionalities that include Fallbacks, Security, Safety and others.

Through this journey, I aimed to make the complex understandable, to empower beginners with practical insights and to pave the way for an enriched understanding of LangChain's capabilities. One major point to note is that LangChain is just a service provider that, internally, helps you with refining your prompts and adding some codes around LLMs. Eventually, it's the LLM that is doing the actual task. So **the quality of the LLM you choose will impact how your app works**. Most of the tutorials in this book used either GPT3.5 or GPT4 (the OpenAI API) hence you were able to see the state-of-the-art results for every app we built. This won't be the case when you go with other, less powerful LLMs.

As I conclude this book, remember that the knowledge gained here serves as a foundation for your further exploration and mastery of this powerful framework. And trust me, there is still a lot of stuff around LangChain that you should explore. Hope to see you soon with my next!

# About the Author

---

Mehul Gupta, a Data Scientist and Researcher based in India brings over 5 years of experience in Data Science and Artificial Intelligence (AI). He holds a Bachelor's degree in Computer Science from the Indian Institute of Information Technology, Design & Manufacturing (IIITDM) Jabalpur.

## Professional Background

While this is Mehul's first book, he has shared his knowledge through **140+ technical blogs** on Medium, reaching over a **million views.** Additionally, Mehul has created about 500 user-friendly Data Science and AI video tutorials on YouTube under the name '**Data Science in your Pocket**.'

## Contributions and Publications

Beyond online platforms, Mehul has contributed to academic research, publishing papers with Elsevier on prescription digitization. These contributions have established him as a trusted figure in the AI industry.

## Areas of Expertise

Mehul is well-versed in AI areas like Generative AI, Natural Language Processing, Reinforcement Learning, Graph Analytics, Time Series Analysis and others. His practical experience includes working with the Healthcare and Finance sectors.

## Connect with Mehul

Explore more about Data Science and AI with Mehul Gupta:

**Medium:** https://medium.com/@mehulgupta_7991

**LinkedIn:** https://www.linkedin.com/in/mehulgupta7991/

**YouTube:** https://www.youtube.com/@datascienceinyourpocket

Mehul is open to collaboration, speaking engagements and sharing knowledge in the world of AI. Shoot a mail for further queries at mehulgupta2016154@gmail.com