

Preface

I have been working in the field of artificial intelligence since 1982 and without a doubt Large Language Models (LLMs) like GPT-3 and infrastructure projects like LangChain are the largest technology breakthroughs that I have lived through. This book will use the [LangChain](#) and [GPT Index \(LlamaIndex\)](#) projects along with the OpenAI GPT-3 and ChatGPT APIs to solve a series of interesting problems.

Harrison Chase started the LangChain project in October 2022 and as I write this book in February 2023 the GitHub repository for LangChain <https://github.com/hwchase17/langchain> has 171 contributors. Jerry Liu started the GPT Index project (recently renamed to LlamaIndex) at the end of 2022 and the GitHub Repository for LlamaIndex https://github.com/jerryliu/gpt_index currently has 54 contributors.

The GitHub repository for examples in this book is <https://github.com/mark-watson/langchain-book-examples.git>. Please note that I usually update the code in the examples repository fairly frequently for library version updates, etc.

While the documentation and examples online for LangChain and LlamaIndex are excellent, I am still motivated to write this book to solve interesting problems that I like to work on involving information retrieval, natural language processing (NLP), dialog agents, and the semantic web/linked data fields. I hope that you, dear reader, will be delighted with these examples and that at least some of them will inspire your future projects.

About the Author

I have written over 20 books, I have over 50 US patents, and I have worked at interesting companies like Google, Capital One, SAIC, Mind AI, and others. You can find links for reading most of my recent books free on my web site <https://markwatson.com>. If I had to summarize my career the short take would be that I have had a lot of fun and enjoyed my work. I hope that what you learn here will be both enjoyable and help you in your work.

If you would like to support my work please consider purchasing my books on [Leanpub](#) and star my git repositories that you find useful on [GitHub](#). You can also interact with me on social media on [Mastodon](#) and [Twitter](#). I am also available as a consultant: <https://markwatson.com>.

Book Cover

I live in Sedona, Arizona. I took the book cover photo in January 2023 from the street that I live on.

Acknowledgements

This picture shows me and my wife Carol who helps me with book production and editing.

 Mark and Carol Watson

Mark and Carol Watson

I would also like to thank the following readers who reported errors or typos in this book: Armando Flores.

Requirements for Running and Modifying Book Examples

I show full source code and a fair amount of example output for each book example so if you don't want to get access to some of the following APIs then you can still read along in the book.

To use OpenAI's GPT-3 and ChatGPT models you will need to sign up for an API key (free tier is OK) at <https://openai.com/api/> and set the environment variable `OPENAI_API_KEY` to your key value.

You will need to get an API key for examples using **Google's Knowledge Graph APIs**.

Reference: [Google Knowledge Graph APIs](#).

The example programs using Google's Knowledge Graph APIs assume that you have the file `~/.google_api_key` in your home directory that contains your key from <https://console.cloud.google.com/apis>.

You will need to install **SerpApi** for examples integrating web search. Please reference: [PyPi project page](#).

You can sign up for a free non-commercial 100 searches/month account with an email address and phone number at <https://serpapi.com/users/welcome>.

You will also need [Zapier](#) account for the GMail and Google Calendar examples.

After reading though this book, you can review the website [LangChainHub](#) which contains prompts, chains and agents that are useful for building LLM applications.

Large Language Model Overview

[Large language models](#) are a subset of artificial intelligence that use deep learning and neural networks to process natural language. [Transformers](#) are a type of neural network architecture that can learn context in sequential data using self-attention mechanisms. They

were introduced in 2017 by a team at Google Brain and have become popular for LLM research. Some examples of [transformer-based](#) LLMs are [BERT](#), [GPT-3](#), [T5](#) and [Megatron-LM](#).

The main points we will discuss in this book are:

- LLMs are deep learning algorithms that can understand and generate natural language based on massive datasets.
- LLMs use techniques such as self-attention, masking, and fine-tuning to learn complex patterns and relationships in language. LLMs can understand and generate natural language because they use transformer models, which are a type of neural network that can process sequential data such as text using attention mechanisms. Attention mechanisms allow the model to focus on relevant parts of the input and output sequences while ignoring irrelevant ones.
- LLMs can perform various natural language processing (NLP) and natural language generation (NLG) tasks, such as summarization, translation, prediction, classification, and question answering.
- Even though LLMs were initially developed for NLP applications, LLMs have also shown potential in other domains such as computer vision and computational biology by leveraging their generalizable knowledge and transfer learning abilities.

[BERT models](#) are one of the first types of transformer models that were widely used. BERT was developed by Google AI Language in 2018. BERT models are a family of masked language models that use transformer architecture to learn bidirectional representations of natural language. BERT models can understand the meaning of ambiguous words by using the surrounding text as context. The “magic trick” here is that training data comes almost free because in masking models, you programmatically chose random words, replace them with a missing word token, and the model is trained to predict the missing words. This process is repeated with massive amounts of training data from the web, books, etc.

Here are some “papers with code” links for BERT (links are for code, paper links in the code repositories):

- <https://github.com/allenai/scibert>
- <https://github.com/google-research/bert>

Big Tech Businesses vs. Small Startups Using Large Language Models

Both Microsoft and Google play both sides of this business game: they want to sell cloud LLM services to developers and small startup companies and they would also like to achieve lock-in for their consumer services like Office 365, Google Docs and Sheets, etc.

Microsoft has been integrating AI technology into workplace emails, slideshows, and spreadsheets as part of its ongoing partnership with OpenAI, the company behind ChatGPT. Microsoft’s Azure OpenAI service offers a powerful tool to enable these outcomes when leveraged

with their data lake of more than two billion metadata and transactional elements.

As I write this, Google has just opened a public wait list for their Bard AI/chat search service. I have used various Google APIs for years in code I write. I have no favorites in the battle between tech giants, rather I am mostly interested in what they build that I can use in my own projects.

Hugging Face, which makes AI products and hosts those developed by other companies, is working on open-source rivals to ChatGPT and will [use AWS](#) for that as well. Cohere AI, Anthropic, Hugging Face, and Stability AI are some of the startups that are using OpenAI and Hugging Face APIs for their products. As I write this chapter, I view Hugging Face as a great source of specialized models. I love that Hugging Face models can be run via their APIs and also self-hosted on our own servers and sometimes even on our laptops. Hugging Face is a fantastic resource and even though I use their models much less frequently in this book than OpenAI APIs, you should embrace the hosting and open source flexibility of Hugging Face. Here I use OpenAI APIs because I want you to get set up and creating your own projects quickly.

Dear reader, I didn't write this book for developers working at established AI companies (although I hope such people find the material here useful!). I wrote this book for small developers who want to scratch their own itch by writing tools that save them time. I also wrote this book hoping that it would help developers build capabilities into the programs they design and write that rival what the big tech companies are doing.

Getting Started With LangChain

Harrison Chase started the LangChain project in October 2022 and as I write this book in February 2023 the GitHub repository for LangChain <https://github.com/hwchase17/langchain> has 171 contributors.

[LangChain](#) is a framework for building applications with large language models (LLMs) through chaining different components together. Some of the applications of LangChain are chatbots, generative question-answering, summarization, data-augmented generation and more. LangChain can save time in building chatbots and other systems by providing a standard interface for chains, agents and memory, as well as integrations with other tools and end-to-end examples. We refer to “chains” as sequences of calls (to an LLMs and a different program utilities, cloud services, etc.) that go beyond just one LLM API call. LangChain provides a standard interface for chains, many integrations with other tools, and end-to-end chains for common applications. Often you will find existing chains already written that meet the requirements for your applications.

For example, one can create a chain that takes user input, formats it using a `PromptTemplate`, and then passes the formatted response to a Large Language Model (LLM) for processing.

While LLMs are very general in nature which means that while they can perform many tasks effectively, they often can not directly provide specific answers to questions or tasks that require deep domain knowledge or expertise. LangChain provides a standard interface for agents, a library of agents to choose from, and examples of end-to-end agents.

LangChain Memory is the concept of persisting state between calls of a chain or agent. LangChain provides a standard interface for memory, a collection of memory implementations, and examples of chains/agents that use memory². LangChain provides a large collection of common utils to use in your application. Chains go beyond just a single LLM call, and are sequences of calls (whether to an LLM or a different utility). LangChain provides a standard interface for chains, lots of integrations with other tools, and end-to-end chains for common applications.

LangChain can be integrated with one or more model providers, data stores, APIs, etc. LangChain can be used for in-depth question-and-answer chat sessions, API interaction, or action-taking. LangChain can be integrated with Zapier's platform through a natural language API interface (we have an entire chapter dedicated to Zapier integrations).

Installing Necessary Packages

For the purposes of examples in this book, you might want to create a new Anaconda or other Python environment and install:

```
1 pip install langchain llama_index openai
2 pip install kor pydrive pandas rdflib
3 pip install google-search-results SPARQLWrapper
```

For the rest of this chapter we will use the subdirectory `langchain_getting_started` and in the next chapter use `llama_index_case_study` in the GitHub repository for this book.

Creating a New LangChain Project

Simple LangChain projects are often just a very short Python script file. As you read this book, when any example looks interesting or useful, I suggest copying the `requirements.txt` and Python source files to a new directory and making your own GitHub private repository to work in. Please make the examples in this book “your code,” that is, freely reuse any code or ideas you find here.

Basic Usage and Examples

While I try to make the material in this book independent, something you can enjoy with no external references, you should also take advantage of the high quality [documentation](Langchain Quickstart Guide) and the individual detailed guides for prompts, chat, document loading, indexes, etc.

As we work through some examples please keep in mind what it is like to use the ChatGPT web application: you enter text and get responses. The way you prompt ChatGPT is obviously important if you want to get useful responses. In code examples we automate and formalize this manual process.

You need to choose a LLM to use. We will usually choose the GPT-3.5 API from OpenAI because it is general purpose and much less expensive than OpenAI's previous model APIs. You will need to [sign up](#) for an API key and set it as an environment variable:

```
1 export OPENAI_API_KEY="YOUR KEY GOES HERE"
```

Both the libraries **openai** and **langchain** will look for this environment variable and use it. We will look at a few simple examples in a Python REPL. We will start by just using OpenAI's text prediction API:

```
1 $ python
2 >>> from langchain.llms import OpenAI
3 >>> llm = OpenAI(temperature=0.8)
4 >>> s = llm("John got into his new sports car, and he dro\
5 ve it")
6 >>> s
7 ' to work\n\nJohn started up his new sports car and drove\
8 it to work. He had a huge smile on his face as he drove,
9 excited to show off his new car to his colleagues. The w
10 ind blowing through his hair, and the sun on his skin, he
11 felt a sense of freedom and joy as he cruised along the
12 road. He arrived at work in no time, feeling refreshed an
13 d energized.'
14 >>> s = llm("John got into his new sports car, and he dro\
15 ve it")
16 >>> s
17 " around town\n\nJohn drove his new sports car around tow\
18 n, enjoying the feeling of the wind in his hair. He stopp
19 ed to admire the view from a scenic lookout, and then spe
20 d off to the mall to do some shopping. On the way home, h
21 e took a detour down a winding country road, admiring the
22 scenery and enjoying the feel of the car's powerful engi
```



```
23 ne. By the time he arrived back home, he had a huge smile
24 on his face."
```

Notice how when we ran the same input text prompt twice that we see different results. Setting the temperature in line 3 to a higher value increases the randomness.

Our next example is in the source file `directions_template.py` and uses the `PromptTemplate` class. A prompt template is a reproducible way to generate a prompt. It contains a text string (“the template”), that can take in a set of parameters from the end user and generate a prompt. The prompt template may contain language model instructions, few-shot examples to improve the model’s response, or specific questions for the model to answer.

```
1 from langchain.prompts import PromptTemplate
2 from langchain.llms import OpenAI
3 llm = OpenAI(temperature=0.9)
4
5 def get_directions(thing_to_do):
6     prompt = PromptTemplate(
7         input_variables=["thing_to_do"],
8         template="How do I {thing_to_do}?",
9     )
10    prompt_text = prompt.format(thing_to_do=thing_to_do)
11    print(f"\n{prompt_text}:")
12    return llm(prompt_text)
13
14 print(get_directions("get to the store"))
15 print(get_directions("hang a picture on the wall"))
```

You could just write Python string manipulation code to create a prompt but using the utility class `PromptTemplate` is more legible and works with any number of prompt input variables.

The output is:

```
1 $ python directions_template.py
2
3 How do I get to the store?:
4
5 To get to the store, you will need to use a mode of trans\
6 portation such as walking, driving, biking, or taking pub
7 lic transportation. Depending on the location of the stor
8 e, you may need to look up directions or maps to determin
9 e the best route to take.
```

```

10
11 How do I hang a picture on the wall?:
12
13 1. Find a stud in the wall, or use two or three wall anch\
14 ors for heavier pictures.
15 2. Measure and mark the wall where the picture hanger wil\
16 l go.
17 3. Pre-drill the holes and place wall anchors if needed.
18 4. Hammer the picture hanger into the holes.
19 5. Hang the picture on the picture hanger.

```

The next example in the file **country_information.py** is derived from an example in the LangChain documentation. In this example we use **PromptTemplate** that contains the pattern we would like the LLM to use when returning a response.

```

1 from langchain.prompts import PromptTemplate
2 from langchain.llms import OpenAI
3 llm = OpenAI(temperature=0.9)
4
5 def get_country_information(country_name):
6     print(f"\nProcessing {country_name}:")
7     global prompt
8     if "prompt" not in globals():
9         print("Creating prompt...")
10        prompt = PromptTemplate(
11            input_variables=["country_name"],
12            template = ""
13 Predict the capital and population of a country.
14
15 Country: {country_name}
16 Capital:
17 Population: "",
18        )
19        prompt_text = prompt.format(country_name=country_name)
20        print(prompt_text)
21        return llm(prompt_text)
22
23 print(get_country_information("Canada"))
24 print(get_country_information("Germany"))

```

You can use the ChatGPT web interface to experiment with prompts and when you find a pattern that works well then write a Python script like the last example, but changing the data you supply in the **PromptTemplate** instance.

The output of the last example is:

```
1 $ python country_information.py
2
3 Processing Canada:
4 Creating prompt...
5
6 Predict the capital and population of a country.
7
8 Country: Canada
9 Capital:
10 Population:
11
12
13 Capital: Ottawa
14 Population: 37,058,856 (as of July 2020)
15
16 Processing Germany:
17
18 Predict the capital and population of a country.
19
20 Country: Germany
21 Capital:
22 Population:
23
24
25 Capital: Berlin
26 Population: 83,02 million (est. 2019)
```

Creating Embeddings

We will reference the [LangChain embeddings documentation](#). We can use a Python REPL to see what text to vector space embeddings might look like:

```
1 $ python
2 Python 3.10.8 (main, Nov 24 2022, 08:08:27) [Clang 14.0.6\
3 ] on darwin
4 Type "help", "copyright", "credits" or "license" for more\
5 information.
6 >>> from langchain.embeddings import OpenAIEmbeddings
7 >>> embeddings = OpenAIEmbeddings()
8 >>> text = "Mary has blond hair and John has brown hair. \
9 Mary lives in town and John lives in the country."
10 >>> doc_embeddings = embeddings.embed_documents([text])
11 >>> doc_embeddings
12 [[0.007727328687906265, 0.0009025644976645708, -0.0033224\
```

```
13 383369088173, -0.01794492080807686, -0.017969949170947075
14 , 0.028506645932793617, -0.013414892368018627, 0.00466768
15 16418766975, -0.0024965214543044567, -0.02662956342101097
16 ,
17 ...]]
18 >>> query_embedding = embeddings.embed_query("Does John live in the city?")
19
20 >>> query_embedding
21 [0.028048301115632057, 0.011499025858938694, -0.009440079\
22 33139801, -0.020809611305594444, -0.023904507979750633, 0
23 .018750663846731186, -0.01626438833773136, 0.018129095435
24 142517,
25 ...]
26 >>>
```

Notice that the **doc_embeddings** is a list where each list element is the embeddings for one input text document. The **query_embedding** is a single embedding. Please read the above linked embedding documentation.

We will use vector stores to store calculated embeddings for future use. In the next chapter we will see a document database search example using LangChain and Llama-Index.

Using LangChain Vector Stores to Query Documents

We will reference the [LangChain Vector Stores documentation](#). You need to install a few libraries:

```
1 pip install chroma
2 pip install chromadb
3 pip install unstructured
```

The example script is **doc_search.py**:

```
1 from langchain.text_splitter import CharacterTextSplitter
2 from langchain.vectorstores import Chroma
3 from langchain.embeddings import OpenAIEmbeddings
4 from langchain.document_loaders import DirectoryLoader
5 from langchain import OpenAI, VectorDBQA
6
7 embeddings = OpenAIEmbeddings()
8
9 loader = DirectoryLoader('../data/', glob="**/*.txt")
```

```

10 documents = loader.load()
11 text_splitter = CharacterTextSplitter(chunk_size=2500, ch\
12 unk_overlap=0)
13
14 texts = text_splitter.split_documents(documents)
15
16 docsearch = Chroma.from_documents(texts, embeddings)
17
18 qa = VectorDBQA.from_chain_type(llm=OpenAI(),
19                                chain_type="stuff",
20                                vectorstore=docsearch)
21
22 def query(q):
23     print(f"Query: {q}")
24     print(f"Answer: {qa.run(q)}")
25
26 query("What kinds of equipment are in a chemistry laborat\
27 ory?")
28 query("What is Austrian School of Economics?")
29 query("Why do people engage in sports?")
30 query("What is the effect of body chemistry on exercise?")

```

The **DirectoryLoader** class is useful for loading a directory full of input documents. In this example we specified that we only want to process text files, but the file matching pattern could have also specified PDF files, etc.

The output is:

```

1 $ python doc_search.py
2 Created a chunk of size 1055, which is longer than the sp\
3 ecified 1000
4 Running Chroma using direct local API.
5 Using DuckDB in-memory for database. Data will be transie\
6 nt.
7 Query: What kinds of equipment are in a chemistry laborat\
8 ory?
9 Answer: A chemistry lab would typically include glasswar\
10 e, such as beakers, flasks, and test tubes, as well as ot
11 her equipment such as scales, Bunsen burners, and thermom
12 eters.
13 Query: What is Austrian School of Economics?
14 Answer: The Austrian School is a school of economic thou\
15 ght that emphasizes the spontaneous organizing power of t
16 he price mechanism. Austrians hold that the complexity of
17 subjective human choices makes mathematical modelling of
18 the evolving market extremely difficult and advocate a "

```

19 laissez faire" approach to the economy. Austrian School e
 20 conomists advocate the strict enforcement of voluntary co
 21 ntractual agreements between economic agents, and hold th
 22 at commercial transactions should be subject to the small
 23 est possible imposition of forces they consider to be (in
 24 particular the smallest possible amount of government in
 25 tervention). The Austrian School derives its name from it
 26 s predominantly Austrian founders and early supporters, i
 27 ncluding Carl Menger, Eugen von Bohm-Bawerk and Ludwig vo
 28 n Mises.

29 Query: Why do people engage in sports?
 30 Answer: People engage in sports for leisure and entertai\
 31 nment, as well as for physical exercise and athleticism.

32 Query: What is the effect of body chemistry on exercise?
 33 Answer: Body chemistry can affect the body's response to\
 34 exercise, as certain hormones and enzymes produced by th
 35 e body can affect the energy levels and muscle performanc
 36 e. Chemicals in the body, such as adenosine triphosphate
 37 (ATP) and urea, can affect the body's energy production a
 38 nd muscle metabolism during exercise. Additionally, the b
 39 ody's levels of electrolytes, vitamins, and minerals can
 40 affect exercise performance.

41 Exiting: Cleaning up .chroma directory

LangChain Overview Wrap Up

We will continue using LangChain for the rest of this book as well as the LlamaIndex library that we introduce in the next chapter.

I cover just the subset of LangChain that I use in my own projects in this book. I urge you to read the LangChain documentation and to explore public LangChain chains that users have written on [Langchain-hub](#).

Overview of LlamaIndex

The popular LlamaIndex project used to be called GPT-Index but has been generalized to work with more models than just GPT-3, for example [using Hugging Face embeddings](#).

[LlamaIndex](#) is a project that provides a central interface to connect your language models with external data. It was created by Jerry Liu and his team in the fall of 2022. It consists of a set of data structures [designed to make it easier to use large external knowledge bases](#)

[with language models](#). Some of its uses are:

- Querying structured data such as tables or databases using natural language
- Retrieving relevant facts or information from large text corpora
- Enhancing language models with domain-specific knowledge

LlamaIndex supports a variety of document types, including:

- Text documents are the most common type of document. They can be stored in a variety of formats, such as .txt, .doc, and .pdf.
- XML documents are a type of text document that is used to store data in a structured format.
- JSON documents are a type of text document that is used to store data in a lightweight format.
- HTML documents are a type of text document that is used to create web pages.
- PDF documents are a type of text document that is used to store documents in a fixed format.

LlamaIndex can also index data that is stored in a variety of databases, including:

- SQL databases such as MySQL, PostgreSQL, and Oracle. NoSQL databases such as MongoDB, Cassandra, and CouchDB.
- Solr is a popular open-source search engine that provides high performance and scalability.
- Elasticsearch is another popular open-source search engine that offers a variety of features, including full-text search, geospatial search, and machine learning.
- Apache Cassandra is a NoSQL database that can be used to store large amounts of data.
- MongoDB is another NoSQL database that is easy to use and scale.
- PostgreSQL is a relational database that is widely used in enterprise applications.

LlamaIndex is a flexible framework that can be used to index a variety of document types and data sources.

We will look first at a short example derived from the LlamaIndex documentation and later look at the parts of the LlamaIndex source code that uses LangChain.

Using LlamaIndex to Search Local Documents Using GPT-3

The following example is similar to the last example in the overview chapter on LangChain. In line 8 we use a utility data loader function provided by LlamaIndex to read documents in an input directory. As a demonstration we save the index (consisting of document embeddings) to disk and reload it. This technique is useful when you have a large number of static documents so the indexing procedure

can take a while and require many OpenAI API calls. As an example, you might have many gigabytes of company documentation that doesn't change often so it makes sense to only occasionally recreate the index.

```
1 # Derived from a documentation example at:
2 # https://github.com/jerryjliu/gpt_index
3
4 # make sure you set the following environment variable
5 # is set: OPENAI_API_KEY
6
7 from llama_index import GPTSimpleVectorIndex,
8                         SimpleDirectoryReader
9 documents = SimpleDirectoryReader('../data').load_data()
10 index = GPTSimpleVectorIndex(documents)
11
12 # save to disk
13 index.save_to_disk('index.json')
14 # load from disk
15 index = GPTSimpleVectorIndex.load_from_disk('index.json')
16
17 # search for a document
18 r = index.query("effect of body chemistry on exercise?")
19 print(r)
```

You may have noticed that the query defined on line 17 is the same query that we used last chapter.

```
1 INFO:root:> [build_index_from_documents] Total LLM token \
2 usage: 0 tokens
3 INFO:root:> [build_index_from_documents] Total embedding \
4 token usage: 2272 tokens
5 INFO:root:> [query] Total LLM token usage: 1018 tokens
6 INFO:root:> [query] Total embedding token usage: 7 tokens
7
8 The effect of body chemistry on exercise can vary dependi\
9 ng on the individual. Factors such as hydration levels, e
10 lectrolyte balance, and the availability of energy-storin
11 g molecules such as adenosine triphosphate (ATP) can all
12 affect the bodys ability to perform exercise. Additionall
13 y, hormones such as adrenaline and cortisol can influence
14 the bodys response to exercise, as can the presence of c
15 ertain nutrients and minerals.
```

Note that in general, many organic chemicals could be listed in response to the query **effect of body chemistry on exercise?** but that just **adenosine triphosphate (ATP)** was mentioned in the test document in the file **data/chemistry.txt**.

Using LlamaIndex for Question Answering from a List of Web Sites

In this example we use the **trafilatura** and **html2text** libraries to get text from a web page that we will index and search. The class **TrafilaturaWebReader** does the work of creating local documents from a list of web page URLs and the index class **GPTListIndex** builds a local index for use with OpenAI API calls to implement search.

The following listing shows the file **web_page_QA.py**:

```
1 # Derived from the example at:
2 # https://github.com/jerryjliu/gpt_index/blob/main/exampl\
3 es/data_connectors/WebPageDemo.ipynb
4
5 # pip install llama-index, html2text, trafilatura
6
7 from llama_index import GPTListIndex
8 from llama_index import TrafilaturaWebReader
9
10 def query_website(url_list, *questions):
11     documents = TrafilaturaWebReader().load_data(url_list)
12     index = GPTListIndex(documents)
13     for question in questions:
14         print(f"\n== QUESTION: {question}\n")
15         response = index.query(question)
16         print(f"== RESPONSE: {response}")
17
18 if __name__ == "__main__":
19     url_list = ["https://markwatson.com"]
20     query_website(url_list, "What programming languages doe\
21 s Mark use?",
22                       "How many books has Mark writte\
23 n?",
24                       "What musical instruments does \
25 Mark play?")
```

This example is not efficient because we create a new index for each web page we want to search. That said, this example (that was derived from an example in the LlamaIndex documentation) implements a pattern that you can use, for example, to build a reusable index of your company's web site and build an end-user web search app.

The output for these three test questions in the last code example is:

```
1 INFO:root:> [build_index_from_documents] Total LLM token \
2 usage: 0 tokens
3 INFO:root:> [build_index_from_documents] Total embedding \
4 token usage: 0 tokens
5
6 == QUESTION: What programming languages does Mark use?
7
8 INFO:root:> [query] Total LLM token usage: 2210 tokens
9 INFO:root:> [query] Total embedding token usage: 0 tokens
10 == RESPONSE:
11 Mark uses a variety of programming languages, including H\
12 askell, Ruby, Clojure, JavaScript, Java, Common Lisp, Pyt
13 hon, and Smalltalk.
14
15 == QUESTION: How many books has Mark written?
16
17 INFO:root:> [query] Total LLM token usage: 2189 tokens
18 INFO:root:> [query] Total embedding token usage: 0 tokens
19 == RESPONSE:
20 Mark has written 20+ books.
21
22 == QUESTION: What musical instruments does Mark play?
23
24 INFO:root:> [query] Total LLM token usage: 2197 tokens
25 INFO:root:> [query] Total embedding token usage: 0 tokens
26 == RESPONSE:
27 Mark plays guitar, didgeridoo, and American Indian flute.
```

LlamaIndex/GPT-Index Case Study Wrap Up

LlamaIndex is a set of data structures and library code designed to make it easier to use large external knowledge bases such as Wikipedia. LlamaIndex creates a vectorized index from your document data, making it highly efficient to query. It then uses this index to identify the most relevant sections of the document based on the query.

LlamaIndex is useful because it provides a central interface to connect your LLM's with external data and offers data connectors to your existing data sources and data formats (API's, PDF's, docs, SQL, etc.). It provides a simple, flexible interface between your external data and LLMs.

Some projects that use LlamaIndex include building personal assistants with LlamaIndex and GPT-3.5, using LlamaIndex for document retrieval, and combining answers across documents.

Using Google's Knowledge Graph APIs With LangChain

Google's Knowledge Graph (KG) is a knowledge base that Google uses to serve relevant information in an infobox beside its search results. It allows the user to see the answer in a glance, as an instant answer. The data is generated automatically from a variety of sources, covering places, people, businesses, and more. I worked at Google in 2013 on a project that used their KG for an internal project.

Google's public Knowledge Graph Search API lets you find entities in the Google Knowledge Graph. The API uses standard schema.org types and is compliant with the JSON-LD specification. It supports entity search and lookup.

You can use the Knowledge Graph Search API to build applications that make use of Google's Knowledge Graph. For example, you can use the API to build a search engine that returns results based on the entities in the Knowledge Graph.

In the next chapter we also use the public KGs DBPedia and Wikidata. One limitation of Google's KG APIs is that it is designed for entity (people, places, organizations, etc.) lookup. When DBPedia and Wikidata it is possible to find a wider range of information using the SPARQL query language, such as relationships between entities. You can use the Google KG APIs to find some entity relationships, e.g., all the movies directed by a particular director, or all the books written by a particular author. You can also use the API to find information like all the people who have worked on a particular movie, or all the actors who have appeared in a particular TV show.

Setting Up To Access Google Knowledge Graph APIs

To get an API key for Google's Knowledge Graph Search API, you need to go to the Google API Console, enable the Google Knowledge Graph Search API, and create an API key to use in your project. You can then use this API key to make requests to the Knowledge Graph Search API.

To create your application's API key, follow these steps:

- Go to the API Console.
- From the projects list, select a project or create a new one.
- If the APIs & services page isn't already open, open the left side menu and select APIs & services.
- On the left, choose Credentials.
- Click Create credentials and then select API key.

You can then use this API key to make requests to the Knowledge Graph Search APIs.

When I use Google's APIs I set the access key in `~/.google_api_key` and read in the key using:

```
1 api_key=open(str(Path.home())+"/.google_api_key").read()
```

You can also use environment variables to store access keys. Here is a code snippet for making an API call to get information about me:

```
1 import json
2 from urllib.parse import urlencode
3 from urllib.request import urlopen
4 from pathlib import Path
5 from pprint import pprint
6
7 api_key =
8     open(str(Path.home()) + "/.google_api_key").read()
9 query = "Mark Louis Watson"
10 service_url =
11     "https://kgsearch.googleapis.com/v1/entities:search"
12 params = {
13     "query": query,
14     "limit": 10,
15     "indent": True,
16     "key": api_key,
17 }
18 url = service_url + "?" + urlencode(params)
19 response = json.loads(urlopen(url).read())
20 pprint(response)
```

The JSON-LD output would look like:

```
1 { '@context': { '@vocab': 'http://schema.org/',
2                 'EntitySearchResult':
3                 'goog:EntitySearchResult',
4                 'detailedDescription':
5                 'goog:detailedDescription',
6                 'goog': 'http://schema.googleapis.com/',
7                 'kg': 'http://g.co/kg',
8                 'resultScore': 'goog:resultScore'},
9   '@type': 'ItemList',
10  'itemListElement': [{ '@type': 'EntitySearchResult',
11                        'result': { '@id': 'kg:/m/0b6_g82',
12                                  '@type': ['Thing',
```

```

13             'Person'],
14             'description': 'Author',
15             'name':
16             'Mark Louis Watson',
17             'url':
18             'http://markwatson.com'},
19         'resultScore': 43}}}]

```

In order to not repeat the code for getting entity information from the Google KG, I wrote a utility **Google_KG_helper.py** that encapsulates the previous code and generalizes it into a mini-library:

```

1  """Client for calling Knowledge Graph Search API."""
2
3  import json
4  from urllib.parse import urlencode
5  from urllib.request import urlopen
6  from pathlib import Path
7  from pprint import pprint
8
9  api_key =
10     open(str(Path.home()) + "/.google_api_key").read()
11
12  # use Google search API to get information
13  # about a named entity:
14
15  def get_entity_info(entity_name):
16     service_url =
17         "https://kgsearch.googleapis.com/v1/entities:search"
18     params = {
19         "query": entity_name,
20         "limit": 1,
21         "indent": True,
22         "key": api_key,
23     }
24     url = service_url + "?" + urlencode(params)
25     response = json.loads(urlopen(url).read())
26     return response
27
28  def tree_traverse(a_dict):
29     ret = []
30     def recur(dict_2, a_list):
31         if isinstance(dict_2, dict):
32             for key, value in dict_2.items():
33                 if key in ['name', 'description',
34                     'articleBody']:

```

```

35         a_list += [value]
36         recur(value, a_list)
37     if isinstance(dict_2, list):
38         for x in dict_2:
39             recur(x, a_list)
40     recur(a_dict, ret)
41     return ret
42
43
44 def get_context_text(entity_name):
45     json_data = get_entity_info(entity_name)
46     return ' '.join(tree_traverse(json_data))
47
48 if __name__ == "__main__":
49     get_context_text("Bill Clinton")

```

The main test script is in the file **Google_Knowledge_Graph_Search.py**:

```

1 """Example of Python client calling the
2   Knowledge Graph Search API."""
3
4 from llama_index import GPTListIndex, Document
5
6 import Google_KG_helper
7
8 def kg_search(entity_name, *questions):
9     ret = ""
10    context_text =
11        Google_KG_helper.get_context_text(entity_name)
12    print(f"Context text: {context_text}")
13    doc = Document(context_text)
14    index = GPTListIndex([doc])
15    for question in questions:
16        response = index.query(question)
17        ret +=
18            f"QUESTION: {question}\nRESPONSE: {response}\n"
19    return ret
20
21 if __name__ == "__main__":
22     s = kg_search("Bill Clinton",
23                 "When was Bill president?")
24     print(s)

```

The example output is:

```
1 $ python Google_Knowledge_Graph_Search.py
2 Context text: Bill Clinton 42nd U.S. President William Je\
3 fferon Clinton is an American retired politician who ser
4 ved as the 42nd president of the United States from 1993
5 to 2001.
6 INFO:root:> [build_index_from_documents] Total LLM token \
7 usage: 0 tokens
8 INFO:root:> [build_index_from_documents] Total embedding \
9 token usage: 0 tokens
10 INFO:root:> [query] Total LLM token usage: 77 tokens
11 INFO:root:> [query] Total embedding token usage: 0 tokens
12 QUESTION: When was Bill president?
13 RESPONSE:
14 Bill Clinton was president from 1993 to 2001.
```

Accessing Knowledge Graphs from Google, DBPedia, and Wikidata allows you to integrate real world facts and knowledge with your applications. While I mostly work in the field of deep learning I frequently also use Knowledge Graphs in my work and in my personal research. I think that you, dear reader, might find accessing highly structured data in KGs to be more reliable and in many cases simpler than using web scraping.

Using DBPedia and WikiData as Knowledge Sources

Both [DBPedia](#) and [Wikidata](#) are public Knowledge Graphs (KGs) that store data as [Resource Description Framework \(RDF\)](#) and are accessed through the [SPARQL Query Language for RDF](#). The examples for this project are in the [GitHub repository for this book](#) in the directory `kg_search`.

I am not going to spend much time here discussing RDF and SPARQL. Instead I ask you to read online the introductory chapter **Linked Data, the Semantic Web, and Knowledge Graphs** in my book [A Lisp Programmer Living in Python-Land: The Hy Programming Language](#).

As we saw in the last chapter, a Knowledge Graph (that I often abbreviate as KG) is a graph database using a schema to define types (both objects and relationships between objects) and properties that link property values to objects. The term “Knowledge Graph” is both a general term and also sometimes refers to the specific Knowledge Graph used at Google which I worked with while working there in 2013. Here, we use KG to reference the general technology of storing knowledge in graph databases.

DBPedia and Wikidata are similar, with some important differences. Here is a summary of some similarities and differences between DBPedia and Wikidata:

- Both projects aim to provide structured data from Wikipedia in various formats and languages. Wikidata also has data from other sources so it contains more data and more languages.
- Both projects use RDF as a common data model and SPARQL as a query language.
- DBpedia extracts data from the infoboxes in Wikipedia articles, while Wikidata collects data entered through its interfaces by both users and automated bots.
- Wikidata requires sources for its data, while DBpedia does not.
- DBpedia is more popular in the Semantic Web and Linked Open Data communities, while Wikidata is more integrated with Wikimedia projects.

To the last point: I personally prefer DBpedia when experimenting with the semantic web and linked data, mostly because DBpedia URIs are human readable while Wikidata URIs are abstract. The following URIs represent the town I live in, Sedona Arizona:

- DBpedia: https://dbpedia.org/page/Sedona,_Arizona
- Wikidata: <https://www.wikidata.org/wiki/Q80041>

In RDF we enclose URIs in angle brackets like `<https://www.wikidata.org/wiki/Q80041>`.

If you read the chapter on RDF and SPARQL in my book link that I mentioned previously, then you know that RDF data is represented by triples where each part is named:

- subject
- property
- object

We will look at two similar examples in this chapter, one using DBpedia and one using Wikidata. Both services have SPARQL endpoint web applications that you will want to use for exploring both KGs. We will look at the DBpedia web interface later. Here is the Wikidata web interface:



In this SPARQL query the prefix **wd:** stands for Wikidata data while the prefix **wdt:** stands for Wikidata type (or property). The prefix **rdfs:** stands for RDF Schema.

Using DBPedia as a Data Source

DBpedia is a community-driven project that extracts structured content from Wikipedia and makes it available on the web as a Knowledge Graph (KG). The KG is a valuable resource for researchers and developers who need to access structured data from Wikipedia. With the use of SPARQL queries to DBpedia as a data source we can write a variety applications, including natural language processing, machine learning, and data analytics. We demonstrate the effectiveness of DBpedia as a data source by presenting several examples that illustrate its use in real-world applications. In my experience, DBpedia is a valuable resource for researchers and developers who need to access structured data from Wikipedia.

In general you will start projects using DBPedia by exploring available data using the web app <https://dbpedia.org/sparql> that can be seen in this screen shot:



The following listing of file **dbpedia_generate_rdf_as_nt.py** shows Python code for making a SPARQL query to DBPedia and saving the results as RDF triples in NT format in a local text file:

```
1 from SPARQLWrapper import SPARQLWrapper
2 from rdflib import Graph
3
4 sparql = SPARQLWrapper("http://dbpedia.org/sparql")
5 sparql.setQuery("""
6     PREFIX dbpedia-owl: <http://dbpedia.org/ontology/>
7     PREFIX dbpedia: <http://dbpedia.org/resource>
8     PREFIX dbpprop: <http://dbpedia.org/property>
9
10    CONSTRUCT {
11        ?city dbpedia-owl:country ?country .
12        ?city rdfs:label ?citylabel .
13        ?country rdfs:label ?countrylabel .
14        <http://dbpedia.org/ontology/country> rdfs:label \
15 "country"@en .
16    }
17    WHERE {
18        ?city rdf:type dbpedia-owl:City .
19        ?city rdfs:label ?citylabel .
20        ?city dbpedia-owl:country ?country .
21        ?country rdfs:label ?countrylabel .
22        FILTER (lang(?citylabel) = 'en')
23        FILTER (lang(?countrylabel) = 'en')
24    }
25    LIMIT 50
26 """)
27 sparql.setReturnFormat("rdf")
28 results = sparql.query().convert()
29
30 g = Graph()
31 g.parse(data=results.serialize(format="xml"), format="xml\
32 ")
33
34 print("\nresults:\n")
35 results = g.serialize(format="nt").encode("utf-8").decode\
36 ('utf-8')
37 print(results)
38
39 text_file = open("sample.nt", "w")
40 text_file.write(results)
41 text_file.close()
```

Here is the printed output from running this script (most output not shown, and manually edited to fit page width):

```

1 $ python generate_rdf_as_nt.py
2 results:
3
4 <http://dbpedia.org/resource/Ethiopia>
5   <http://www.w3.org/2000/01/rdf-schema#label>
6   "Ethiopia"@en .
7 <http://dbpedia.org/resource/Valentin_Alsina,_Buenos_Aire\
8 s>
9   <http://www.w3.org/2000/01/rdf-schema#label>
10  "Valentin Alsina, Buenos Aires"@en .
11 <http://dbpedia.org/resource/Davyd-Haradok>
12   <http://dbpedia.org/ontology/country>
13   <http://dbpedia.org/resource/Belarus> .
14 <http://dbpedia.org/resource/Davyd-Haradok>
15   <http://www.w3.org/2000/01/rdf-schema#label>
16   "Davyd-Haradok"@en .
17 <http://dbpedia.org/resource/Belarus>
18   <http://www.w3.org/2000/01/rdf-schema#label>
19   "Belarus"@en .
20 ...

```

This output was written to a local file **sample.nt**. I divided this example into two separate Python scripts because I thought it would be easier for you, dear reader, to experiment with fetching RDF data separately from using a LLM to process the RDF data. In production you may want to combine KG queries with semantic analysis.

This code example demonstrates the use of the **GPTSimpleVectorIndex** for querying RDF data and retrieving information about countries. The function **download_loader** loads data importers by string name. While it is not a type safe to load a Python class by name using a string, if you misspell the name of the class to load the call to **download_loader** then a Python **ValueError**(**“Loader class name not found in library”**) error is thrown. The **GPTSimpleVectorIndex** class represents an index data structure that can be used to efficiently search and retrieve information from the RDF data. This is similar to other types of **LlamaIndex** vector index types for different types of data sources.

Here is the script **dbpedia_rdf_query.py**:

```

1 "Example from documentation"
2
3 from llama_index import GPTSimpleVectorIndex, Document
4 from llama_index import download_loader
5
6 RDFReader = download_loader("RDFReader")

```

```
7 doc = RDFReader().load_data("sample.nt")
8 index = GPTSimpleVectorIndex(doc)
9
10 result = index.query("list all countries in a quoted Pyth\
11 on array, then explain why")
12
13 print(result.response)
```

Here is the output:

```
1 $ python rdf_query.py
2 INFO:root:> [build_index_from_documents] Total LLM token \
3 usage: 0 tokens
4 INFO:root:> [build_index_from_documents] Total embedding \
5 token usage: 761 tokens
6 INFO:root:> [query] Total LLM token usage: 921 tokens
7 INFO:root:> [query] Total embedding token usage: 12 tokens
8
9 ['Argentina', 'French Polynesia', 'Democratic Republic of\
10 the Congo', 'Benin', 'Ethiopia', 'Australia', 'Uzbekista
11 n', 'Tanzania', 'Albania', 'Belarus', 'Vanuatu', 'Armenia
12 ', 'Syria', 'Andorra', 'Venezuela', 'France', 'Vietnam',
13 'Azerbaijan']
14
15 This is a list of all the countries mentioned in the cont\
16 ext information. All of the countries are listed in the c
17 ontext information, so this list is complete.
```

Why are there only 18 countries listed? In the script used to perform a SPARQL query on DBPedia, we had a statement **LIMIT 50** at the end of the query so only 50 RDF triples were written to the file **sample.nt** that only contains data for 18 countries.

Using Wikidata as a Data Source

It is slightly more difficult exploring Wikidata compared to DBPedia. Let's revisit getting information about my home town of Sedona Arizona.

In writing this example, I experimented with SPARQL queries using the [Wikidata SPARQL web app](#).

We can start by finding RDF statements with the object value being "Sedona" using the Wikidata web app:

```
1 select * where {
2   ?s ?p "Sedona"@en
3 } LIMIT 30
```

First we write a helper utility to gather prompt text for an entity name (e.g., name of a person, place, etc.) in the file **wikidata_generate_prompt_text.py**:

```
1 from SPARQLWrapper import SPARQLWrapper, JSON
2 from rdflib import Graph
3 import pandas as pd
4
5 def get_possible_eitity_uris_from_wikidata(entity_name):
6     sparql = SPARQLWrapper("https://query.wikidata.org/spa\
7 rql")
8     sparql.setQuery("""
9         SELECT ?entity ?entityLabel WHERE {
10             ?entity rdfs:label "%s"@en .
11         } limit 5
12     """ % entity_name)
13
14     sparql.setReturnFormat(JSON)
15     results = sparql.query().convert()
16
17     results = pd.json_normalize(results['results']['bindin\
18 gs']).values.tolist()
19     results = ["<" + x[1] + ">" for x in results]
20     return [*set(results)] # remove duplicates
21
22 def wikidata_query_to_df(entity_uri):
23     sparql = SPARQLWrapper("https://query.wikidata.org/spa\
24 rql")
25     sparql.setQuery("""
26         SELECT ?description ?is_a_type_of WHERE {
27             %s schema:description ?description FILTER (lang(?\
28 description) = 'en') .
29             %s wdt:P31 ?instanceOf .
30             ?instanceOf rdfs:label ?is_a_type_of FILTER (lang\
31 (?is_a_type_of) = 'en') .
32         } limit 10
33     """ % (entity_uri, entity_uri))
34
35     sparql.setReturnFormat(JSON)
36     results = sparql.query().convert()
37     results2 = pd.json_normalize(results['results']['bindi\
38 ngs'])
```

```

39     prompt_text = ""
40     for index, row in results2.iterrows():
41         prompt_text += row['description.value'] + " is a \
42 type of " + row['is_a_type_of.value'] + "\n"
43     return prompt_text
44
45 def generate_prompt_text(entity_name):
46     entity_uris = get_possible_eitity_uris_from_wikidata(e\
47 ntity_name)
48     prompt_text = ""
49     for entity_uri in entity_uris:
50         p = wikidata_query_to_df(entity_uri)
51         if "disambiguation page" not in p:
52             prompt_text += entity_name + " is " + wikidata\
53 _query_to_df(entity_uri)
54     return prompt_text
55
56 if __name__ == "__main__":
57     print("Sedona:", generate_prompt_text("Sedona"))
58     print("California:",
59         generate_prompt_text("California"))
60     print("Bill Clinton:",
61         generate_prompt_text("Bill Clinton"))
62     print("Donald Trump:",
63         generate_prompt_text("Donald Trump"))

```

This utility does most of the work in getting prompt text for an entity.

The **GPTTreeIndex** class is similar to other LlamaIndex index classes. This class builds a tree-based index of the prompt texts, which can be used to retrieve information based on the input question. In LlamaIndex, a **GPTTreeIndex** is used to select the child node(s) to send the query down to. A **GPTKeywordTableIndex** uses keyword matching, and a **GPTVectorStoreIndex** uses embedding cosine similarity. The choice of which index class to use depends on how much text is being indexed, what the granularity of subject matter in the text is, and if you want summarization.

GPTTreeIndex is also more efficient than **GPTSimpleVectorIndex** because it uses a tree structure to store the data. This allows for faster searching and retrieval of data compared to a linear list index class like **GPTSimpleVectorIndex**.

The LlamaIndex code is relatively easy to implement in the script **wikidata_query.py** (edited to fit page width):

```

1 from llama_index import StringIterableReader, GPTTreeIndex
2 from wikidata_generate_prompt_text import generate_prompt\
3 _text

```



```

4
5 def wd_query(question, *entity_names):
6     prompt_texts = []
7     for entity_name in entity_names:
8         prompt_texts +=
9             [generate_prompt_text(entity_name)]
10    documents =
11        StringIterableReader().load_data(texts=prompt_texts)
12    index = GPTTreeIndex(documents)
13    return index.query(question)
14
15 if __name__ == "__main__":
16     print("Sedona:", wd_query("What is Sedona?", "Sedona"))
17     print("California:",
18         wd_query("What is California?", "California"))
19     print("Bill Clinton:",
20         wd_query("When was Bill Clinton president?",
21             "Bill Clinton"))
22     print("Donald Trump:",
23         wd_query("Who is Donald Trump?",
24             "Donald Trump"))

```

Here is the test output (with some lines removed):

```

1 $ python wikidata_query.py
2 Total LLM token usage: 162 tokens
3 INFO:llama_index.token_counter.token_counter:> [build_ind\
4 ex_from_documents] INFO:llama_index.indices.query.tree.le
5 af_query:> Starting query: What is Sedona?
6 INFO:llama_index.token_counter.token_counter:> [query] To\
7 tal LLM token usage: 154 tokens
8 Sedona: Sedona is a city in the United States located in \
9 the counties of Yavapai and Coconino, Arizona. It is also
10 the title of a 2012 film, a company, and a 2015 single b
11 y Houndmouth.
12
13 Total LLM token usage: 191 tokens
14 INFO:llama_index.indices.query.tree.leaf_query:> Starting\
15 query: What is California?
16 California: California is a U.S. state in the United Stat\
17 es of America.
18
19 Total LLM token usage: 138 tokens
20 INFO:llama_index.indices.query.tree.leaf_query:> Starting\
21 query: When was Bill Clinton president?
22 Bill Clinton: Bill Clinton was the 42nd President of the \

```

```
23 United States from 1993 to 2001.  
24  
25 Total LLM token usage: 159 tokens  
26 INFO:llama_index.indices.query.tree.leaf_query:> Starting\  
27 query: Who is Donald Trump?  
28 Donald Trump: Donald Trump is the 45th President of the U\  
29 nited States, serving from 2017 to 2021.
```

Using LLMs To Organize Information in Our Google Drives

My digital life consists of writing, working as an AI practitioner, and learning activities that I justify with my self-image of a “gentleman scientist.” Cloud storage like GitHub, Google Drive, Microsoft OneDrive, and iCloud are central to my activities.

About ten years ago I spent two months of my time writing a system in Clojure that was planned to be my own custom and personal DropBox, augmented with various NLP tools and a FireFox plugin to send web clippings directly to my personal system. To be honest, I stopped using my own project after a few months because the time it took to organize my information was a greater opportunity cost than the value I received.

In this chapter I am going to walk you through parts of a new system that I am developing for my own personal use to help me organize my material on Google Drive (and eventually other cloud services). Don’t be surprised if the completed project is an additional example in a future edition of this book!

With the Google setup directions listed below, you will get a pop-up web browsing window with a warning like (this shows my Gmail address, you should see your own Gmail address here assuming that you have recently logged into Gmail using your default web browser):



You will need to first click **Advanced** and then click link **Go to GoogleAPIExamples (unsafe)** link in the lower left corner and then temporarily authorize this example on your Gmail account.

Setting Up Requirements.

You need to create a credential at <https://console.cloud.google.com/cloud-resource-manager> (copied from the [PyDrive documentation](#), changing application type to “Desktop”):

- Search for ‘Google Drive API’, select the entry, and click ‘Enable’.
- Select ‘Credentials’ from the left menu, click ‘Create Credentials’, select ‘OAuth client ID’.
- Now, the product name and consent screen need to be set → click ‘Configure consent screen’ and follow the instructions. Once finished:
- Select ‘Application type’ to be Desktop application.
- Enter an appropriate name.
- Input `http://localhost:8080` for ‘Authorized JavaScript origins’.
- Input `http://localhost:8080/` for ‘Authorized redirect URIs’.
- Click ‘Save’.
- Click ‘Download JSON’ on the right side of Client ID to download `client_secret_.json`. Copy the downloaded JSON credential file to the example directory `google_drive_llm` for this chapter.

Write Utility To Fetch All Text Files From Top Level Google Drive Folder

For this example we will just authenticate our test script with Google, and copy all top level text files with names ending with “.txt” to the local file system in subdirectory **data**. The code is in the directory `google_drive_llm` in file `fetch_txt_files.py` (edited to fit page width):

```
1 from pydrive.auth import GoogleAuth
2 from pydrive.drive import GoogleDrive
3 from pathlib import Path
4
5 # good GD search docs:
6 # https://developers.google.com/drive/api/guides/search-f\
7 iles
8
9 # Authenticate with Google
10 gauth = GoogleAuth()
11 gauth.LocalWebserverAuth()
```

```

12 drive = GoogleDrive(gauth)
13
14 def get_txt_files(dir_id='root'):
15     " get all plain text files with .txt extension in top\
16     level Google Drive directory "
17
18     file_list = drive.ListFile({'q': f"'{dir_id}' in pare\
19 nts and trashed=false"}).GetList()
20     for file1 in file_list:
21         print('title: %s, id: %s' % (file1['title'], file\
22 l['id']))
23     return [[file1['title'], file1['id'], file1.GetConten\
24 tString()]]
25         for file1 in file_list
26             if file1['title'].endswith(".txt")]
27
28 def create_test_file():
29     " not currently used, but useful for testing. "
30
31     # Create GoogleDriveFile instance with title 'Hello.t\
32 xt':
33     file1 = drive.CreateFile({'title': 'Hello.txt'})
34     file1.SetContentString('Hello World!')
35     file1.Upload()
36
37 def test():
38     fl = get_txt_files()
39     for f in fl:
40         print(f)
41         file1 = open("data/" + f[0], "w")
42         file1.write(f[2])
43         file1.close()
44
45 if __name__ == '__main__':
46     test()

```

For testing I just have one text file with the file extension “.txt” on my Google Drive so my output from running this script looks like the following listing. I edited the output to change my file IDs and to only print a few lines of the debug printout of file titles.

```

1 $ python fetch_txt_files.py
2 Your browser has been opened to visit:
3
4     https://accounts.google.com/o/oauth2/auth?client_id=5\
5 29311921932-xsmj3hhiplr0dhqjln13fo4rrtvoslo8.apps.googleu
6 sercontent.com&redirect_uri=http%3A%2F%2Flocalhost%3B6180

```

```

7 %2F&scope=https%3A%2F%2Fwww.googleapis.com%2Fauth%2Fdrive
8 &access_type=offline&response_type=code
9
10 Authentication successful.
11
12 title: testdata, id: 1TZ9bnL5XYQvKACJw8VoKwDVJ8jeCszJ
13 title: sports.txt, id: 18RN4ojvURWt5yoKNtDdAJbh4fvmRpzwb
14 title: Anaconda blog article, id: 1kpLaYQA4Ao8ZbdFaXU209h\
15 g-z0tv1xA7YQO4L8y8NbU
16 title: backups_2023, id: 1-k_r1HTfuZRWN7vwWWSYqfssl0C96J2x
17 title: Work notes, id: 1fDyHyZtKI-0oRNabA_P41LltYjGoek21
18 title: Sedona Writing Group Contact List, id: 1zK-5v9OQUf\
19 y8Sw33nTCl9vnL822hLlw
20 ...
21 ['sports.txt', '18RN4ojvURWt5yoKNtDdAJbh4fvmRpzwb', 'Spor\
22 t is generally recognised as activities based in physical
23 athleticism or physical dexterity.[3] Sports are usually
24 governed by rules to ensure fair competition and consist
25 ent adjudication of the winner.\n\n"Sport" comes from the
26 Old French desport meaning "leisure", with the oldest de
27 finition in English from around 1300 being "anything huma
28 ns find amusing or entertaining".[4]\n\nOther bodies advo
29 cate widening the definition of sport to include all phys
30 ical activity and exercise. For instance, the Council of
31 Europe include all forms of physical exercise, including
32 those completed just for fun.\n\n']

```

Generate Vector Indices for Files in Specific Google Drive Directories

The example script in the last section should have created copies of the text files in your home Google Documents directory that end with “.txt”. Here, we use the same LlamaIndex test code that we used in a previous chapter. The test script `index_and_QA.py` is listed here:

```

1 # make sure you set the following environment variable is\
2 set:
3 # OPENAI_API_KEY
4
5 from llama_index import GPTSimpleVectorIndex, SimpleDirec\
6 toryReader
7 documents = SimpleDirectoryReader('data').load_data()
8 index = GPTSimpleVectorIndex(documents)
9
10 # save to disk
11 index.save_to_disk('index.json')
12 # load from disk

```

```

13 index = GPTSimpleVectorIndex.load_from_disk('index.json')
14
15 # search for a document
16 print(index.query("What is the definition of sport?"))

```

For my test file, the output looks like:

```

1 $ python index_and_QA.py
2 INFO:llama_index.token_counter.token_counter:> [build_ind\
3 ex_from_documents] Total LLM token usage: 0 tokens
4 INFO:llama_index.token_counter.token_counter:> [build_ind\
5 ex_from_documents] Total embedding token usage: 111 token
6 s
7 INFO:llama_index.token_counter.token_counter:> [query] To\
8 tal LLM token usage: 202 tokens
9 INFO:llama_index.token_counter.token_counter:> [query] To\
10 tal embedding token usage: 7 tokens
11
12 Sport is generally recognised as activities based in phys\
13 ical athleticism or physical dexterity that are governed
14 by rules to ensure fair competition and consistent adjudi
15 cation of the winner. It is anything humans find amusing
16 or entertaining, and can include all forms of physical ex
17 ercise, even those completed just for fun.

```

It is interesting to see how the query result is rewritten in a nice form, compared to the raw text in the file **sports.txt** on my Google Drive:

```

1 $ cat data/sports.txt
2 Sport is generally recognised as activities based in phys\
3 ical athleticism or physical dexterity.[3] Sports are usu
4 ally governed by rules to ensure fair competition and con
5 sistent adjudication of the winner.
6
7 "Sport" comes from the Old French desport meaning "leisur\
8 e", with the oldest definition in English from around 130
9 0 being "anything humans find amusing or entertaining".[4
10 ]
11
12 Other bodies advocate widening the definition of sport to\
13 include all physical activity and exercise. For instance
14 , the Council of Europe include all forms of physical exe
15 rcise, including those completed just for fun.

```

Google Drive Example Wrap Up

If you already use Google Drive to store your working notes and other documents, then you might want to expand the simple example in this chapter to build your own query system for your documents. In addition to Google Drive, I also use Microsoft Office 365 and OneDrive in my work and personal projects.

I haven't written my own connectors yet for OneDrive but this is on my personal to-do list using the Microsoft library <https://github.com/OneDrive/onedrive-sdk-python>.

Using Zapier Integrations With GMail and Google Calendar

Zapier is a service for writing integrations with hundreds of cloud services. Here we will write some demos for writing automatic integrations with GMail and Google Calendar.

Using the Zapier service is simple. You need to register the services you want to interact with on the Zapier developer web site and then you can express how you want to interact with services using natural language prompts.

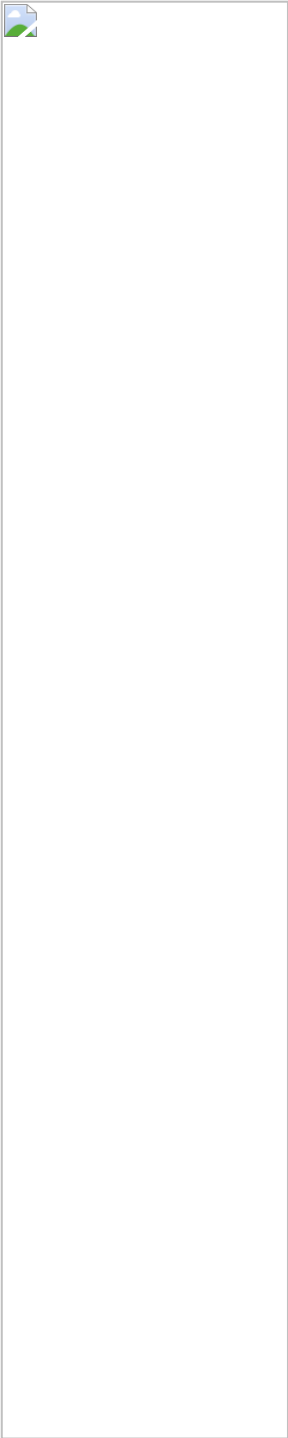
Set Up Development Environment

You will need a developer key for [Zapier Natural Language Actions API](#). Go to this linked web page and look for “Dev App” in the “Provider Name” column. If a key does not exist, you'll need to set up an action to create a key. Click “Set up Actions” and follow the instructions. Your key will be in the Personal API Key column for the “Dev App.” Click to reveal and copy your key. You can [read the documentation](#).

When I set up my Zapier account I set up three Zapier Natural Language Actions:

- Gmail: Fine Email
- Gmail: Send Email
- Google Calendar: Find Event

If you do the same then you will see the Zapier registered actions:



Sending a Test GMail

In the following example replace **TEST_EMAIL_ADDRESS** with an email address that you can use for testing.

```
1 from langchain.llms import OpenAI
2 from langchain.agents import initialize_agent
3 from langchain.agents.agent_toolkits import ZapierToolkit
4 from langchain.utilities.zapier import ZapierNLAWrapper
5
6 llm = OpenAI(temperature=0)
7 zapier = ZapierNLAWrapper()
8 toolkit = ZapierToolkit.from_zapier_nla_wrapper(zapier)
9 agent = initialize_agent(toolkit.get_tools(), llm, agent=\
10 "zero-shot-react-description", verbose=True)
11
12 agent.run("Send an Email to TEST_EMAIL_ADDRESS via gmail \
13 that is a pitch for hiring Mark Watson as a consultant fo
14 r deep learning and large language models")
```

Here is the sample output:

```
1 $ python send_gmail.py
2
3
4 > Entering new AgentExecutor chain...
5 I need to use the Gmail: Send Email tool
6 Action: Gmail: Send Email
7 Action Input: Send an email to TEST_EMAIL_ADDRESS with th\
8 e subject "Pitch for Hiring Mark Watson as a Consultant f
9 or Deep Learning and Large Language Models" and the body
10 "Dear Mark Watson, I am writing to you to pitch the idea
11 of hiring you as a consultant for deep learning and large
12 language models. I believe you have the expertise and ex
13 perience to help us achieve our goals. Please let me know
14 if you are interested in discussing further. Thank you f
15 or your time."
16 Cc: not enough information provided in the instruction, m\
17 issing Cc
18 Observation: {"labelIds": "SENT"}
19 Thought: I now know the final answer
20 Final Answer: An email has been sent to TEST_EMAIL_ADDRES\
21 S with the subject "Pitch for Hiring Mark Watson as a Con
22 sultant for Deep Learning and Large Language Models" and
23 the body "Dear Mark Watson, I am writing to you to pitch
```

```

24 the idea of hiring you as a consultant for deep learning
25 and large language models. I believe you have the experi
26 se and experience to help us achieve our goals. Please le
27 t me know if you are interested in discussing further. Th
28 ank you for your time."
29
30 > Finished chain.

```

Google Calendar Integration Example

Assuming that you configured the Zapier Natural Language Action “Google Calendar: Find Event” then the same code we used to send an email in the last section works for checking calendar entries, you just need to change the natural language prompt:

```

1 from langchain.llms import OpenAI
2 from langchain.agents import initialize_agent
3 from langchain.agents.agent_toolkits import ZapierToolkit
4 from langchain.utilities.zapier import ZapierNLAWrapper
5
6 llm = OpenAI(temperature=0)
7 zapier = ZapierNLAWrapper()
8 toolkit = ZapierToolkit.from_zapier_nla_wrapper(zapier)
9 agent = initialize_agent(toolkit.get_tools(), llm,
10                          agent="zero-shot-react-descripti\
11 on", verbose=True)
12
13 agent.run("Get my Google Calendar entries for tomorrow")

```

And the output looks like:

```

1 $ python get_google_calendar.py
2
3 > Entering new AgentExecutor chain...
4 I need to find events in my Google Calendar
5 Action: Google Calendar: Find Event
6 Action Input: Find events in my Google Calendar tomorrow
7 Observation: {"location": "Greg to call Mark on (928) XXX\
8 -ZZZZ", "kind": "calendar#event", "end__dateTime": "2023-
9 03-23T10:00:00-07:00", "status": "confirmed", "end__dateT
10 ime_pretty": "Mar 23, 2023 10:00AM", "htmlLink": "https:/
11 /zpr.io/WWWWWWWWW"}
12 Thought: I now know the final answer
13 Final Answer: I have an event in my Google Calendar tomor\

```

```
14 row at 10:00AM.  
15  
16 > Finished chain.
```

I edited this output to remove some private information.

Natural Language SQLite Database Queries With LangChain

The LangChain library support for SQLite databases uses the Python library SQLAlchemy for database connections. This abstraction layer allows LangChain to use the same logic and models for other relational databases.

I have a long work history of writing natural language interfaces for relational databases that I will review in the chapter wrap up. For now, I invite you to be amazed at how simple it is to write the LangChain scripts for querying a database in natural language.

We will use the SQLite sample database from the SQLite Tutorial web site:

```
1 https://www.sqlitetutorial.net/sqlite-sample-database/
```

This database has 11 tables. The above URI has documentation for this database so please take a minute to review the [table schema diagram and text description](#).

This example is derived from the [LangChain documentation](#). We use three classes from the langchain library:

- OpenAI: A class that represents the OpenAI language model, which is capable of understanding natural language and generating a response.
- SQLiteDatabase: A class that represents a connection to an SQL database.
- SQLiteDatabaseChain: A class that connects the OpenAI language model with the SQL database to allow natural language querying.

The temperature parameter set to 0 in this example. The temperature parameter controls the randomness of the generated output. A lower value (like 0) makes the model's output more deterministic and focused, while a higher value introduces more randomness (or "creativity"). The run method of the db_chain object translates the natural language query into an appropriate SQL query, execute it on the connected database, and then returns the result converting the output into natural language.

```

1 # SQLite NLP Query Demo Script
2
3 from langchain import OpenAI, SQLDatabase
4 from langchain import SQLDatabaseChain
5
6 db = SQLDatabase.from_uri("sqlite:///chinook.db")
7 llm = OpenAI(temperature=0)
8
9 db_chain = SQLDatabaseChain(llm=llm, database=db, verbose\
10 =True)
11
12 db_chain.run("How many employees are there?")
13 db_chain.run("What is the name of the first employee?")
14 db_chain.run("Which customer has the most invoices?")
15 db_chain.run("List all music genres in the database")

```

The output (edited for brevity) shows the generated SQL queries and the query results:

```

1 $ python sqlite_chat_test.py
2
3 > Entering new SQLDatabaseChain chain...
4 How many employees are there?
5  SELECT COUNT(*) FROM employees;
6 SQLResult: [(8,)]
7 Answer: There are 8 employees.
8 > Finished chain.
9
10 > Entering new SQLDatabaseChain chain...
11 What is the name of the first employee?
12  SELECT FirstName, LastName FROM employees WHERE Employee\
13 Id = 1;
14 SQLResult: [('Andrew', 'Adams')]
15 Answer: The first employee is Andrew Adams.
16 > Finished chain.
17
18 > Entering new SQLDatabaseChain chain...
19 Which customer has the most invoices?
20  SELECT customers.FirstName, customers.LastName, COUNT(in\
21 voices.InvoiceId) AS NumberOfInvoices FROM customers INNE
22 R JOIN invoices ON customers.CustomerId = invoices.Custome
23 rId GROUP BY customers.CustomerId ORDER BY NumberOfInvoi
24 ces DESC LIMIT 5;
25 SQLResult: [('Luis', 'Goncalves', 7), ('Leonie', 'Kohler'\
26 , 7), ('Francois', 'Tremblay', 7), ('Bjorn', 'Hansen', 7)
27 , ('Frantisek', 'Wichterlova', 7)]
28 Answer: Luis Goncalves has the most invoices with 7.

```

```

29 > Finished chain.
30
31 > Entering new SQLiteDatabaseChain chain...
32 List all music genres in the database
33 SQLQuery: SELECT Name FROM genres
34 SQLResult: [('Rock',), ('Jazz',), ('Metal',), ('Alternati\
35 ve & Punk',), ('Rock And Roll',), ('Blues',), ('Latin',),
36 ('Reggae',), ('Pop',), ('Soundtrack',), ('Bossa Nova',),
37 ('Easy Listening',), ('Heavy Metal',), ('R&B/Soul',), ('
38 Electronica/Dance',), ('World',), ('Hip Hop/Rap',), ('Sci
39 ence Fiction',), ('TV Shows',), ('Sci Fi & Fantasy',), ('
40 Drama',), ('Comedy',), ('Alternative',), ('Classical',),
41 ('Opera',)]
42 Answer: Rock, Jazz, Metal, Alternative & Punk, Rock And R\
43 oll, Blues, Latin, Reggae, Pop, Soundtrack, Bossa Nova, E
44 asy Listening, Heavy Metal, R&B/Soul, Electronica/Dance,
45 World, Hip Hop/Rap, Science Fiction, TV Shows, Sci Fi & F
46 antasy, Drama, Comedy, Alternative, Classical, Opera
47 > Finished chain.

```

Natural Language Database Query Wrap Up

I had an example I wrote for the first two editions of my [Java AI book](#) (I later removed this example because the code was too long and too difficult to follow). I later reworked this example in Common Lisp and used both versions in several consulting projects in the late 1990s and early 2000s.

The last book I wrote [Practical Python Artificial Intelligence Programming](#) used an OpenAI example https://github.com/openai/openai-cookbook/blob/main/examples/Backtranslation_of_SQL_queries.py that relatively simple code (relative to my older hand-written Java and Common Lisp code) for a NLP database interface.

Compared to the elegant support for NLP database queries in LangChain, the previous examples have limited power and required a lot more code. As I write this in March 2023, it is a good feeling that for the rest of my career, NLP database access is now a solved problem!

Examples Using Hugging Face Open Source Models

To start with you will need to create a free account on the [Hugging Face Hub](#) and get an API key and install:

```
1 pip install --upgrade huggingface_hub
```

You need to set the foillowing environment variable to your Hugging Face Hub access token:

```
1 HUGGINGFACEHUB_API_TOKEN
```

So far in this book we have been using the OpenAI LLM wrapper:

```
1 from langchain.llms import OpenAI
```

Here we will use the alternative Hugging Face wrapper class:

```
1 from langchain import HuggingFaceHub
```

The LangChain library hides most of the details of using both APIs. This is a really good thing. I have had a few discussions on social tech media with people who object to the non open source nature of OpenAI. While I like the convenience of using OpenAI's APIs, I always like to have alternatives for proprietary technology I use.

The Hugging Face Hub endpoint in LangChain connects to the Hugging Face Hub and runs the models via their free inference endpoints. We need a Hugging Face account and API key to use these endpoints³. There exists two Hugging Face LLM wrappers, one for a local pipeline and one for a model hosted on Hugging Face Hub. Note that these wrappers only work for models that support the text2text-generation and text-generation tasks. Text2text-generation refers to the task of generating a text sequence from another text sequence. For example, generating a summary of a long article. Text-generation refers to the task of generating a text sequence from scratch.

Using LangChain as a Wrapper for Hugging Face Prediction Model APIs

We will start with a simple example using the prompt text support in LangChain. The following example is in the script `simple_example.py`:

```
1 from langchain import HuggingFaceHub, LLMChain
2 from langchain.prompts import PromptTemplate
3
4 hub_llm = HuggingFaceHub(
5     repo_id='google/flan-t5-xl',
6     model_kwargs={'temperature':1e-6})
```

```
7 )
8
9 prompt = PromptTemplate(
10     input_variables=["name"],
11     template="What year did {name} get elected as preside\
12 nt?",
13 )
14
15 llm_chain = LLMChain(prompt=prompt, llm=hub_llm)
16
17 print(llm_chain.run("George Bush"))
```

By changing just a few lines of code, you can run many of the examples in this book using the Hugging Face APIs in place of the OpenAI APIs.

The LangChain documentation lists the source code for a wrapper to use local Hugging Face embeddings [here](#).

Creating a Custom LlamaIndex Hugging Face LLM Wrapper Class That Runs on Your Laptop

We will be downloading the Hugging Face model **facebook/opt-1.3b** that is a 2.6 gigabyte file. This model is downloaded the first time it is requested and is then cached in `~/.cache/huggingface/hub` for later reuse.

This example is modified from an example for custom LLMs in the [LlamaIndex documentation](#). Note that I have used a much smaller model in this example and reduced the prompt and output text size.

```
1 # Derived from example:
2 # https://gpt-index.readthedocs.io/en/latest/how_to/cus\
3 tom_llms.html
4
5 import time
6 import torch
7 from langchain.llms.base import LLM
8 from llama_index import SimpleDirectoryReader, LangchainE\
9 mbedding
10 from llama_index import GPTListIndex, PromptHelper
11 from llama_index import LLMPredictor
12 from transformers import pipeline
13
14 max_input_size = 512
```



```
15 num_output = 64
16 max_chunk_overlap = 10
17 prompt_helper = PromptHelper(max_input_size, num_output, \
18 max_chunk_overlap)
19
20 class CustomLLM(LLM):
21     model_name = "facebook/opt-1.3b"
22     # I am not using a GPU, but you can add device="cuda:\
23 0"
24     # to the pipeline call if you have a local GPU or
25     # are running this on Google Colab:
26     pipeline = pipeline("text-generation", model=model_na\
27 me,
28                         model_kwargs={"torch_dtype":torch\
29 .bfloat16})
30
31     def _call(self, prompt, stop = None):
32         prompt_length = len(prompt)
33         response = self.pipeline(prompt, max_new_tokens=n\
34 um_output)
35         first_response = response[0]["generated_text"]
36         # only return newly generated tokens
37         returned_text = first_response[prompt_length:]
38         return returned_text
39
40     @property
41     def _identifying_params(self):
42         return {"name_of_model": self.model_name}
43
44     @property
45     def _llm_type(self):
46         return "custom"
47
48 time1 = time.time()
49
50 # define our LLM
51 llm_predictor = LLMPredictor(llm=CustomLLM())
52
53 # Load the your data
54 documents = SimpleDirectoryReader('../data_small').load_d\
55 ata()
56 index = GPTListIndex(documents, llm_predictor=llm_predict\
57 or,
58                       prompt_helper=prompt_helper)
59
60 time2 = time.time()
61 print(f"Time to load model from disk: {time2 - time1} sec\
```

```
62 onds.")
63
64 # Query and print response
65 response = index.query("What is the definition of sport?")
66 print(response)
67
68 time3 = time.time()
69 print(f"Time for query/prediction: {time3 - time2} second\
70 s.")
```

When running on my M1 MacBook Pro using only the CPU (no GPU or Neural Engine configuration) we can read the model from disk quickly but it takes a while to process queries:

```
1 $ python hf_transformer_local.py
2 INFO:llama_index.token_counter.token_counter:> [build_ind\
3 ex_from_documents] Total LLM token usage: 0 tokens
4 INFO:llama_index.token_counter.token_counter:> [build_ind\
5 ex_from_documents] Total embedding token usage: 0 tokens
6 Time to load model from disk: 1.5303528308868408 seconds.
7 INFO:llama_index.token_counter.token_counter:> [query] To\
8 tal LLM token usage: 182 tokens
9 INFO:llama_index.token_counter.token_counter:> [query] To\
10 tal embedding token usage: 0 tokens
11
12 "Sport" comes from the Old French desport meaning "leisur\
13 e", with the oldest definition in English from around 130
14 0 being "anything humans find amusing or entertaining".[4
15 ]
16 Time for query/prediction: 228.8184850215912 seconds.
```

Even though my M1 MacBook does fairly well when I configure TensorFlow and PyTorch to use the Apple Silicon GPUs and Neural Engines, I usually do my model development using Google Colab.

Let's rerun the last example on Colab:



Using a standard Colab GPU, the query/prediction time is much faster. Here is a [link to my Colab notebook](#) if you would prefer to run this example on Colab instead of on your laptop.

Using Large Language Models to Write Recipes

If you ask the ChatGPT web app to write a recipe using a user supplied ingredient list and a description it does a fairly good job at generating recipes. For the example in this chapter I am taking a different approach:

- Use the recipe and ingredient files from my web app <http://cookingspace.com> to create context text, given a user prompt for a recipe.
- Treat this as a text prediction problem.
- Format the response for display.

This approach has an advantage (for me!) that the generated recipes will be more similar to the recipes I enjoy cooking since the context data will be derived from my own recipe files.

Preparing Recipe Data

I am using the JSON Recipe files from my web app <http://cookingspace.com>. The following Python script converts my JSON data to text descriptions, one per file:

```
1 import json
2
3 def process_json(fpath):
4     with open(fpath, 'r') as f:
5         data = json.load(f)
6
7     for d in data:
8         with open(f"text_data/{d['name']}.txt", 'w') as f:
9             f.write("Recipe name: " + d['name'] + '\n\n')
10            f.write("Number of servings: " +
11                    str(d['num_served']) + '\n\n')
12            ingredients = [" " + str(ii['amount']) +
13                           " " + ii['units'] + " " +
14                           ii['description']
15                           for ii in d['ingredients']]
16            f.write("Ingredients:\n" +
```

```

17         "\n".join(ingrediants) + '\n\n')
18     f.write("Directions: " +
19           ' '.join(d['directions']) + '\n')
20
21 if __name__ == "__main__":
22     process_json('data/vegetarian.json')
23     process_json('data/desert.json')
24     process_json('data/fish.json')
25     process_json('data/meat.json')
26     process_json('data/misc.json')

```

Here is a listing of one of the shorter generated recipe files (i.e., text recipe data converted from raw JSON recipe data from my CookingSpace.com web site):

```

1 Recipe name: Black Bean Dip
2
3 Number of servings: 6
4
5 Ingredients:
6   2 cup Refried Black Beans
7   1/4 cup Sour cream
8   1 teaspoon Ground cumin
9   1/2 cup Salsa
10
11 Directions: Use either a food processor or a mixing bowl \
12 and hand mixer to make this appetizer. Blend the black be
13 ans and cumin for at least one minute until the mixture i
14 s fairly smooth. Stir in salsa and sour cream and lightly
15 mix. Serve immediately or store in the refrigerator.

```

I have generated 41 individual recipe files that will be used for the remainder of this chapter.

In the next section when we use a LLM to generate a recipe, the directions are numbered steps and the formatting is different than my original recipe document files.

A Prediction Model Using the OpenAI text-davinci-002 Model

Here we use the **DirectoryLoader** class that we have used in previous examples to load and then create an embedding index.

Here is the listing for the script **recipe_generator.py**:

```

1 from langchain.text_splitter import CharacterTextSplitter
2 from langchain.vectorstores import Chroma
3 from langchain.embeddings import OpenAIEmbeddings
4 from langchain.document_loaders import DirectoryLoader
5 from langchain import OpenAI, VectorDBQA
6
7 embeddings = OpenAIEmbeddings()
8
9 loader = DirectoryLoader('./text_data/', glob="**/*.txt")
10 documents = loader.load()
11 text_splitter = CharacterTextSplitter(chunk_size=2500,
12                                     chunk_overlap=0)
13
14 texts = text_splitter.split_documents(documents)
15
16 docsearch = Chroma.from_documents(texts, embeddings)
17
18 qa = VectorDBQA.from_chain_type(llm=OpenAI(temperature=0,
19                                     model_name=
20                                     "text-davinci-002"),
21                                 chain_type="stuff",
22                                 vectorstore=docsearch)
23
24 def query(q):
25     print(f"\n\nRecipe creation request: {q}\n")
26     print(f"{qa.run(q)}\n\n")
27
28 query("Create a new recipe using Broccoli and Chicken")
29 query("Create a recipe using Beans, Rice, and Chicken")

```

This generated two recipes. Here is the output for the first request:

```

1 $ python recipe_generator.py
2 Running Chroma using direct local API.
3 Using DuckDB in-memory for database. Data will be transie\
4 nt.
5
6 Recipe creation request: Create a new recipe using both B\
7 roccoli and Chicken
8
9 Recipe name: Broccoli and Chicken Teriyaki
10 Number of servings: 4
11
12 Ingredients:
13 1 cup broccoli
14 1 pound chicken meat

```

```
15 2 tablespoons soy sauce
16 1 tablespoon honey
17 1 tablespoon vegetable oil
18 1 clove garlic, minced
19 1 teaspoon rice vinegar
20
21 Directions:
22
23 1. In a large bowl, whisk together soy sauce, honey, vege\
24 table oil, garlic, and rice vinegar.
25 2. Cut the broccoli into small florets. Add the broccoli \
26 and chicken to the bowl and toss to coat.
27 3. Preheat a grill or grill pan over medium-high heat.
28 4. Grill the chicken and broccoli for 5-7 minutes per sid\
29 e, or until the chicken is cooked through and the broccol
30 i is slightly charred.
31 5. Serve immediately.
```

If you examine the text recipe files I indexed you see that the prediction model merged information from multiple training data recipes while creating new original text for directions that is loosely based on the directions that I wrote and information encoded in the OpenAI text-davinci-002 model.

Here is the output for the second request:

```
1 Recipe creation request: Create a recipe using Beans, Ric\
2 e, and Chicken
3
4 Recipe name: Beans and Rice with Chicken
5 Number of servings: 4
6 Ingredients:
7 1 cup white rice
8 1 cup black beans
9 1 chicken breast, cooked and shredded
10 1/2 teaspoon cumin
11 1/2 teaspoon chili powder
12 1/4 teaspoon salt
13 1/4 teaspoon black pepper
14 1 tablespoon olive oil
15 1/2 cup salsa
16 1/4 cup cilantro, chopped
17
18 Directions:
19 1. Cook rice according to package instructions.
20 2. In a medium bowl, combine black beans, chicken, cumin,\
```

```
21  chili powder, salt, and black pepper.  
22 3. Heat olive oil in a large skillet over medium heat. Ad\  
23 d the bean mixture and cook until heated through, about 5  
24 minutes.  
25 4. Stir in salsa and cilantro. Serve over cooked rice.
```

Cooking Recipe Generation Wrap Up

Cooking is one of my favorite activities (in addition to hiking, kayaking, and playing a variety of musical instruments). I originally wrote the [CookingSpace.com](#) web app to scratch a personal itch: due to a medical issue I had to closely monitor and regulate my vitamin K intake. I used the US Government's USDA Nutrition Database to estimate the amounts of vitamins and nutrients in some recipes that I use.

When I wanted to experiment with generative models, backed by my personal recipe data, to create recipes, having available recipe data from my previous project as well as tools like OpenAI APIs and LangChain made this experiment simple to set up and run. It is a common theme in this book that it is now relatively easy to create personal projects based on our data and our interests.

Book Wrap Up

This book has been fun to write but it has also somewhat frustrating.

It was fun because I have never been as excited by new technology as I have by LLMs and utility software like LangChain and LlamaIndex for building personalized applications.

This book was frustrating in the sense that it is now so very easy to build applications that just a few years would have been impossible to write. Usually when I write books I have two criteria: I only write about things that I am personally interested in and use, and I also hope to figure out non-obvious edge cases and make easier for my readers to use new tech. Here my frustration is writing about something that it is increasingly simple to do so I feel like my value is diminished.

All that said I hope, dear reader, that you found this book to be worth your time reading.

What am I going to do next? Although I am not fond of programming in JavaScript (although I find Typescript to be somewhat better), I want to explore the possibilities of writing an open source Persistent Storage Web App Personal Knowledge Base Management System. I

might get pushback on this but I would probably make it Apple Safari specific so I can use Apple's CloudKit JS to make its use seamless across macOS, iPadOS, and iOS. If I get the right kind of feedback on social media I might write a book around this project.

Thank you for reading my book!

Best regards, Mark Watson