

# LangChain and LlamaIndex Projects Lab Book: Hooking Large Language Models Up to the Real World

by Mark Watson







# LangChain and LlamaIndex Projects Lab Book: Hooking Large Language Models Up to the Real World

Using GPT-4, ChatGPT, and Hugging Face Models in  
Applications.

Mark Watson

This book is for sale at <http://leanpub.com/langchain>

This version was published on 2024-02-01



This is a [Leanpub](#) book. Leanpub empowers authors and publishers with the Lean Publishing process. [Lean Publishing](#) is the act of publishing an in-progress ebook using lightweight tools and many iterations to get reader feedback, pivot until you have the right book and build traction once you do.



This work is licensed under a [Creative Commons Attribution-NonCommercial-ShareAlike 4.0 International License](#)

## Also By **Mark Watson**

[Safe For Humans AI](#)

[Das LangChain und LlamaIndex Projects Lab Buch: Large Language Models für die Echte Welt](#)

[Practical Python Artificial Intelligence Programming](#)

[Practical Artificial Intelligence Development With Racket](#)

[Practical Artificial Intelligence Programming With Clojure](#)

[Artificial Intelligence Using Swift](#)

[A Lisp Programmer Living in Python-Land: The Hy Programming Language](#)

[Haskell Tutorial and Cookbook](#)

[Loving Common Lisp, or the Savvy Programmer's Secret Weapon](#)

[Practical Artificial Intelligence Programming With Java](#)

# Contents

<b>Preface</b>	<b>1</b>
About the Author	1
Book Cover	2
Acknowledgements	2
Requirements for Running and Modifying Book Examples	3
Issues and Workarounds for Using the Material in this Book	4
<b>Large Language Model Overview</b>	<b>5</b>
Big Tech Businesses vs. Small Startups Using Large Language Models	5
<b>Getting Started With LangChain</b>	<b>7</b>
Installing Necessary Packages	7
Creating a New LangChain Project	8
Basic Usage and Examples	8
Creating Embeddings	11
Using LangChain Vector Stores to Query Documents	12
Example Using LangChain Integrations: Using Server APIs for Google Search	14
LangChain Overview Wrap Up	15
<b>Overview of LlamaIndex</b>	<b>16</b>
Using LlamaIndex to Search Local Documents Using GPT-3	16
Using LlamaIndex for Question Answering from a List of Web Sites	18
LlamaIndex/GPT-Index Case Study Wrap Up	19
<b>Retrieval Augmented Generation (RAG) Applications</b>	<b>21</b>
<b>Using Google's Knowledge Graph APIs With LangChain</b>	<b>22</b>
Setting Up To Access Google Knowledge Graph APIs	22
<b>Using DBpedia and WikiData as Knowledge Sources</b>	<b>27</b>
Using DBpedia as a Data Source	28
Using Wikidata as a Data Source	32
<b>Using LLMs To Organize Information in Our Google Drives</b>	<b>36</b>
Setting Up Requirements.	37
Write Utility To Fetch All Text Files From Top Level Google Drive Folder	38

Generate Vector Indices for Files in Specific Google Drive Directories . . . . .	40
Google Drive Example Wrap Up . . . . .	41
<b>Using Zapier Integrations With GMail and Google Calendar . . . . .</b>	<b>42</b>
Set Up Development Environment . . . . .	42
Sending a Test GMail . . . . .	43
Google Calendar Integration Example . . . . .	44
<b>Natural Language SQLite Database Queries With LangChain . . . . .</b>	<b>46</b>
Natural Language Database Query Wrap Up . . . . .	48
<b>Examples Using Hugging Face Open Source Models . . . . .</b>	<b>49</b>
Using LangChain as a Wrapper for Hugging Face Prediction Model APIs . . . . .	49
Creating a Custom LlamaIndex Hugging Face LLM Wrapper Class That Runs on Your Laptop	50
<b>Running Local LLMs Using Llama.cpp and LangChain . . . . .</b>	<b>54</b>
Installing Llama.cpp with a Llama2-13b-orca Model . . . . .	54
Python Example . . . . .	54
<b>Running Local LLMs Using Ollama . . . . .</b>	<b>57</b>
Simple Use of a local Mistral Model Using LangChain . . . . .	57
Minimal Example Using Ollama with the Mistral Open Model for Retrieval Augmented	
Queries Against Local Documents . . . . .	58
<b>Wrap Up for Running Local LLMs Using Ollama . . . . .</b>	<b>61</b>
<b>Using Large Language Models to Write Recipes . . . . .</b>	<b>62</b>
Preparing Recipe Data . . . . .	62
A Prediction Model Using the OpenAI text-embedding-3-large Model . . . . .	63
Cooking Recipe Generation Wrap Up . . . . .	66
<b>LangChain Agents . . . . .</b>	<b>67</b>
Overview of LangChain Tools . . . . .	67
Overview of ReAct Library for Implementing Reading in LMS Applications . . . . .	68
LangChain Agent Tool Example Using DBPedia SPARQL Queries . . . . .	69
LangChain Agent Tools Wrap Up . . . . .	75
<b>More Useful Libraries for Working with Unstructured Text Data . . . . .</b>	<b>76</b>
EmbedChain Wrapper for LangChain Simplifies Application Development . . . . .	76
Kor Library . . . . .	77
<b>Book Wrap Up . . . . .</b>	<b>80</b>

# Preface

I have been working in the field of artificial intelligence since 1982 and without a doubt Large Language Models (LLMs) like GPT-4 represent the greatest advance in practical AI technology that I have experienced. Infrastructure projects like LangChain and LlamaIndex make it simpler to use LLMs and provide some level of abstraction to facilitate switching between LLMs. This book will use the [LangChain](#)<sup>1</sup> and [LlamaIndex](#)<sup>2</sup> projects along with the OpenAI GPT-3.5, GPT-4, ChatGPT APIs, and local models run on your computer to solve a series of interesting problems.

As I write this new edition in January 2024, I mostly run local LLMs but I still use OpenAI APIs, as well as APIs from Anthropic for the Claude 2 model.

Harrison Chase started the LangChain project in October 2022 and as I wrote the first book in February 2023 the GitHub repository for LangChain 171 contributors and as I write this new edition LangChain has 2100+ contributors on GitHub. 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](https://github.com/jerryliu/gpt_index)<sup>3</sup> 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](#)<sup>4</sup> and star

---

<sup>1</sup><https://github.com/langchain-ai/langchain>

<sup>2</sup>[https://github.com/run-llama/llama\\_index](https://github.com/run-llama/llama_index)

<sup>3</sup>[https://github.com/run-llama/llama\\_index](https://github.com/run-llama/llama_index)

<sup>4</sup><https://leanpub.com/u/markwatson>

my git repositories that you find useful on [GitHub](#)<sup>5</sup>. You can also interact with me on social media on [Mastodon](#)<sup>6</sup> and [Twitter](#)<sup>7</sup>. 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.

---

<sup>5</sup><https://github.com/mark-watson?tab=repositories&q=&type=public>

<sup>6</sup>[https://mastodon.social/@mark\\_watson](https://mastodon.social/@mark_watson)

<sup>7</sup>[https://twitter.com/mark\\_l\\_watson](https://twitter.com/mark_l_watson)





Figure 1. Mark and Carol Watson

I would also like to thank the following readers who reported errors or typos in this book: Armando Flores, Peter Solimine, and David Rupp.

## 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](#)<sup>8</sup>.

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](#)<sup>9</sup>.

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**<sup>10</sup> account for the GMail and Google Calendar examples.

After reading though this book, you can review the website [LangChainHub](#)<sup>11</sup> which contains prompts, chains and agents that are useful for building LLM applications.

## Issues and Workarounds for Using the Material in this Book

The libraries that I use in this book are frequently updated and sometimes the documentation or code links change, invalidating links in this book. I will try to keep everything up to date. Please report broken links to me.

In some cases you will need to use specific versions or libraries for some of the code examples.

Because the Python code listings use colorized text you may find that copying code from this eBook may drop space characters. All of the code listings are in the GitHub repository for this book so you should clone the repository to experiment with the example code.

---

<sup>8</sup><https://cloud.google.com/enterprise-knowledge-graph/docs/search-api>

<sup>9</sup><https://pypi.org/project/google-search-results/>

<sup>10</sup><https://zapier.com>

<sup>11</sup><https://github.com/hwchase17/langchain-hub>

# Large Language Model Overview

Large language models<sup>1</sup> are a subset of artificial intelligence that use deep learning and neural networks to process natural language. Transformers<sup>2</sup> 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<sup>3</sup> LLMs are BERT, GPT-3, T5 and Megatron-LM<sup>4</sup>.

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<sup>5</sup> 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>

---

<sup>1</sup><https://blogs.nvidia.com/blog/2022/10/10/llms-ai-horizon/>

<sup>2</sup><https://www.linkedin.com/pulse/chatgpt-tip-iceberg-paul-golding>

<sup>3</sup><https://factored.ai/transformer-based-language-models/>

<sup>4</sup>[https://en.wikipedia.org/wiki/Transformer\\_\(machine\\_learning\\_model\)](https://en.wikipedia.org/wiki/Transformer_(machine_learning_model))

<sup>5</sup>[https://en.wikipedia.org/wiki/BERT\\_\(Language\\_model\)](https://en.wikipedia.org/wiki/BERT_(Language_model))

## 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](#)<sup>6</sup> 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.

---

<sup>6</sup><https://iblnews.org/aws-partners-with-hugging-face-an-ai-startup-rival-to-chatgpt-working-on-open-source-models/>

# Getting Started With LangChain

[LangChain](#)<sup>1</sup> 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 memory<sup>2</sup>. 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:

---

<sup>1</sup>[https://api.python.langchain.com/en/stable/langchain\\_api\\_reference.html](https://api.python.langchain.com/en/stable/langchain_api_reference.html)



```
1 pip install -U langchain llama_index openai langchain-openai
2 pip install -U langchain-community langchainhub
3 pip install -U kor pydrive pandas rdflib
4 pip install -U 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 [Langchain Quickstart Guide documentation](#)<sup>2</sup> 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-4 API from OpenAI because it is general purpose but is more expensive than the older GPT-3.5 APIs. You will need to [sign up](#)<sup>3</sup> 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:

---

<sup>2</sup>[https://python.langchain.com/docs/get\\_started/quickstart](https://python.langchain.com/docs/get_started/quickstart)

<sup>3</sup><https://platform.openai.com/account/api-keys>

```

1  $ python
2  >>> from langchain_openai import OpenAI
3  >>> llm = OpenAI(temperature=0.8)
4  >>> s = llm("John got into his new sports car, and he drove it")
5  >>> s
6  ' to work\n\nJohn started up his new sports car and drove it to work. He had a huge \
7  smile on his face as he drove, excited to show off his new car to his colleagues. Th
8  e wind blowing through his hair, and the sun on his skin, he felt a sense of freedom
9  and joy as he cruised along the road. He arrived at work in no time, feeling refres
10 hed and energized.'
11 >>> s = llm("John got into his new sports car, and he drove it")
12 >>> s
13 " around town\n\nJohn drove his new sports car around town, enjoying the feeling of \
14 the wind in his hair. He stopped to admire the view from a scenic lookout, and then
15 sped off to the mall to do some shopping. On the way home, he took a detour down a w
16 inding country road, admiring the scenery and enjoying the feel of the car's powerfu
17 l engine. By the time he arrived back home, he had a huge smile 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` in the directory `langchain_getting_started` 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_openai 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 transportation such as walking, \
6  driving, biking, or taking public transportation. Depending on the location of the s
7  tore, you may need to look up directions or maps to determine the best route to take
8  .
9
10 How do I hang a picture on the wall?:
11
12 1. Find a stud in the wall, or use two or three wall anchors for heavier pictures.
13 2. Measure and mark the wall where the picture hanger will go.
14 3. Pre-drill the holes and place wall anchors if needed.
15 4. Hammer the picture hanger into the holes.
16 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_openai 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](https://python.langchain.com/en/latest/reference/modules/embeddings.html)<sup>4</sup>. 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 ] on darwin
3 Type "help", "copyright", "credits" or "license" for more information.
4 >>> from langchain_openai import OpenAIEmbeddings
5 >>> embeddings = OpenAIEmbeddings(model="text-embedding-3-large")
6 >>> text = "Mary has blond hair and John has brown hair. Mary lives in town and John\
7   lives in the country."
8 >>> doc_embeddings = embeddings.embed_documents([text])
9 >>> doc_embeddings
10 [[0.007727328687906265, 0.0009025644976645708, -0.0033224383369088173, -0.0179449208\
11 0807686, -0.017969949170947075, 0.028506645932793617, -0.013414892368018627, 0.00466
12 76816418766975, -0.0024965214543044567, -0.02662956342101097,
13 ...]]
14 >>> query_embedding = embeddings.embed_query("Does John live in the city?")
15 >>> query_embedding
16 [0.028048301115632057, 0.011499025858938694, -0.00944007933139801, -0.02080961130559\
17 4444, -0.023904507979750633, 0.018750663846731186, -0.01626438833773136, 0.018129095
18 435142517,
19 ...]
20 >>>
```

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](https://python.langchain.com/en/latest/reference/modules/vectorstore.html)<sup>5</sup>. You need to install a few libraries:

---

<sup>4</sup><https://python.langchain.com/en/latest/reference/modules/embeddings.html>

<sup>5</sup><https://python.langchain.com/en/latest/reference/modules/vectorstore.html>



```
1 pip install chroma
2 pip install chromadb
3 pip install unstructured pdf2image pytesseract
```

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 VectorDBQA
6 from langchain_openai import OpenAI
7
8 embeddings = OpenAIEmbeddings(model="text-embedding-3-large")
9
10 loader = DirectoryLoader('../data/', glob="**/*.txt")
11 documents = loader.load()
12 text_splitter = CharacterTextSplitter(chunk_size=2500, 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(),
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 laboratory?")
27 query("What is Austrian School of Economics?")
28 query("Why do people engage in sports?")
29 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 specified 1000
3 Running Chroma using direct local API.
4 Using DuckDB in-memory for database. Data will be transient.
5 Query: What kinds of equipment are in a chemistry laboratory?
6 Answer: A chemistry lab would typically include glassware, such as beakers, flasks,\
7 and test tubes, as well as other equipment such as scales, Bunsen burners, and ther
8 mometers.
9 Query: What is Austrian School of Economics?
10 Answer: The Austrian School is a school of economic thought that emphasizes the spo\
11 ntaneous organizing power of the price mechanism. Austrians hold that the complexity
12 of subjective human choices makes mathematical modelling of the evolving market ext
13 remely difficult and advocate a "laissez faire" approach to the economy. Austrian Sc
14 hool economists advocate the strict enforcement of voluntary contractual agreements
15 between economic agents, and hold that commercial transactions should be subject to
16 the smallest possible imposition of forces they consider to be (in particular the sm
17 allest possible amount of government intervention). The Austrian School derives its
18 name from its predominantly Austrian founders and early supporters, including Carl M
19 enger, Eugen von Bohm-Bawerk and Ludwig von Mises.
20 Query: Why do people engage in sports?
21 Answer: People engage in sports for leisure and entertainment, as well as for physici\
22 cal exercise and athleticism.
23 Query: What is the effect of body chemistry on exercise?
24 Answer: Body chemistry can affect the body's response to exercise, as certain hormo\
25 nes and enzymes produced by the body can affect the energy levels and muscle perform
26 ance. Chemicals in the body, such as adenosine triphosphate (ATP) and urea, can affe
27 ct the body's energy production and muscle metabolism during exercise. Additionally,
28 the body's levels of electrolytes, vitamins, and minerals can affect exercise perfo
29 rmance.
30 Exiting: Cleaning up .chroma directory
```

## Example Using LangChain Integrations: Using Server APIs for Google Search

The example shown here is in the directory `from_langchain_docs` in the source file `search_simple.py`. The relevant LangChain Integrations documentation page is [https://python.langchain.com/docs/integrations/tools/google\\_serper](https://python.langchain.com/docs/integrations/tools/google_serper).

```
1  # make sure SERPER_API_KEY is set in your environment
2
3  from langchain_community.utilities import GoogleSerperAPIWrapper
4  search_helper = GoogleSerperAPIWrapper()
5
6  def search(query):
7      return search_helper.run(query)
8
9  print(search("What is the capital of Arizona?"))
10 #print(search("Sedona Arizona?"))
```

You will need a Server API key from <https://serper.dev>. Currently you can get a free key for 2500 API calls. After that the paid tier currently starts at \$50 for 50K API calls and these credits must be used within a 6 month period.

## 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](https://github.com/hwchase17/langchain-hub)<sup>6</sup>.

---

<sup>6</sup><https://github.com/hwchase17/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](#)<sup>1</sup>.

[LlamaIndex](#)<sup>2</sup> 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](#)<sup>3</sup>. 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.

---

<sup>1</sup>[https://gpt-index.readthedocs.io/en/latest/how\\_to/customization/embeddings.html](https://gpt-index.readthedocs.io/en/latest/how_to/customization/embeddings.html)

<sup>2</sup>[https://github.com/jerryliu/gpt\\_index/blob/main/docs/index.rst](https://github.com/jerryliu/gpt_index/blob/main/docs/index.rst)

<sup>3</sup><https://gpt-index.readthedocs.io/en/latest/index.html>

## 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  # make sure you set the following environment variable is set:
2  #   OPENAI_API_KEY
3
4  from llama_index import GPTVectorStoreIndex, SimpleDirectoryReader
5  from llama_index import StorageContext, load_index_from_storage
6
7  documents = SimpleDirectoryReader('../data').load_data()
8  # index = GPTListIndex(documents) # llama_index < 0.5
9  index = GPTVectorStoreIndex.from_documents(documents)
10 engine = index.as_query_engine()
11 print(engine.query("what are key economic indicators?"))
12
13 # save to disk
14 index.storage_context.persist(persist_dir='./cache')
15
16 # load from disk
17 storage_context = StorageContext.from_defaults(persist_dir='./cache')
18 index = load_index_from_storage(storage_context)
19 engine = index.as_query_engine()
20
21 # search for a document
22 print(engine.query("effect of body chemistry on exercise?"))
```

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



```

1  $ python test_from_docs.py
2
3  Key economic indicators are measures of economic activity that are used to assess th\
4  e health of an economy. Examples of key economic indicators include gross domestic p
5  roduct (GDP), unemployment rate, inflation rate, consumer price index (CPI), balance
6  of trade, and industrial production.
7
8  The effect of body chemistry on exercise depends on the type of exercise being perfo\
9  rmed. For aerobic exercise, the body needs to be able to produce enough energy to su
10 stain the activity. This requires the body to have adequate levels of oxygen, glucos
11 e, and other nutrients in the bloodstream. For anaerobic exercise, the body needs to
12 be able to produce energy without relying on oxygen. This requires the body to have
13 adequate levels of lactic acid, creatine phosphate, and other energy-producing mole
14 cules. In both cases, the body's chemistry must be balanced in order for the exercis
15 e to be effective.

```

## 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 URIs 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/examples/data_connectors/WebPageD\
3  emo.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) # llama_index < 0.5
13     index = GPTListIndex.from_documents(documents)
14     engine = index.as_query_engine()
15     for question in questions:

```

```

16         print(f"\n== QUESTION: {question}\n")
17         response = engine.query(question)
18         print(f"== RESPONSE: {response}")
19
20 if __name__ == "__main__":
21     url_list = ["https://markwatson.com"]
22     query_website(url_list, "How many patents does Mark have?",
23                   "How many books has Mark written?")

```

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  $ python web_page_QA.py
2  documents:
3  [Document(text='Mark Watson's Artificial Intelligence Books and Blog\n'
4             'Author of 20+ books on AI and I have 50+ US patents. Here I '
5             'talk about both technology and provide links to read my '
6             'published books online.\n'
7             "By registering you agree to Substack's Terms of Service, our "
8             'Privacy Policy, and our Information Collection Notice',
9             doc_id='c6ed3ddd-b160-472b-9c7d-fd0446762cc8',
10            embedding=None,
11            doc_hash='126a335ed76d5d79ad94470dd4dd06ab5556b9f1347e5be25ecc595e0290ab57\
12 ',
13            extra_info=None)]
14
15 == QUESTION: How many patents does Mark have?
16
17 == RESPONSE:
18 Mark Watson has 50+ US patents.
19
20 == QUESTION: How many books has Mark written?
21
22 == RESPONSE:
23 Mark has written 20+ books on AI.

```

## 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.

# Retrieval Augmented Generation (RAG) Applications

TBD: this chapter will be complete by March 1, 2023

Note to readers: this chapter was added February 2024 and supersedes some of the older material in later chapters on indexing and chatting about data from local text and PDF documents as well as data from web sites.

Retrieval Augmented Generation (RAG) Applications work by pre-processing a user's query to search for indexed document fragments that are semantically similar to the user's query. These fragments are concatenated together as context text that is attached to the user's query and then passed of to a LLM model. The LLM can preferentially use information in this context text as well as innate knowledge stored in the LLM to process user queries.

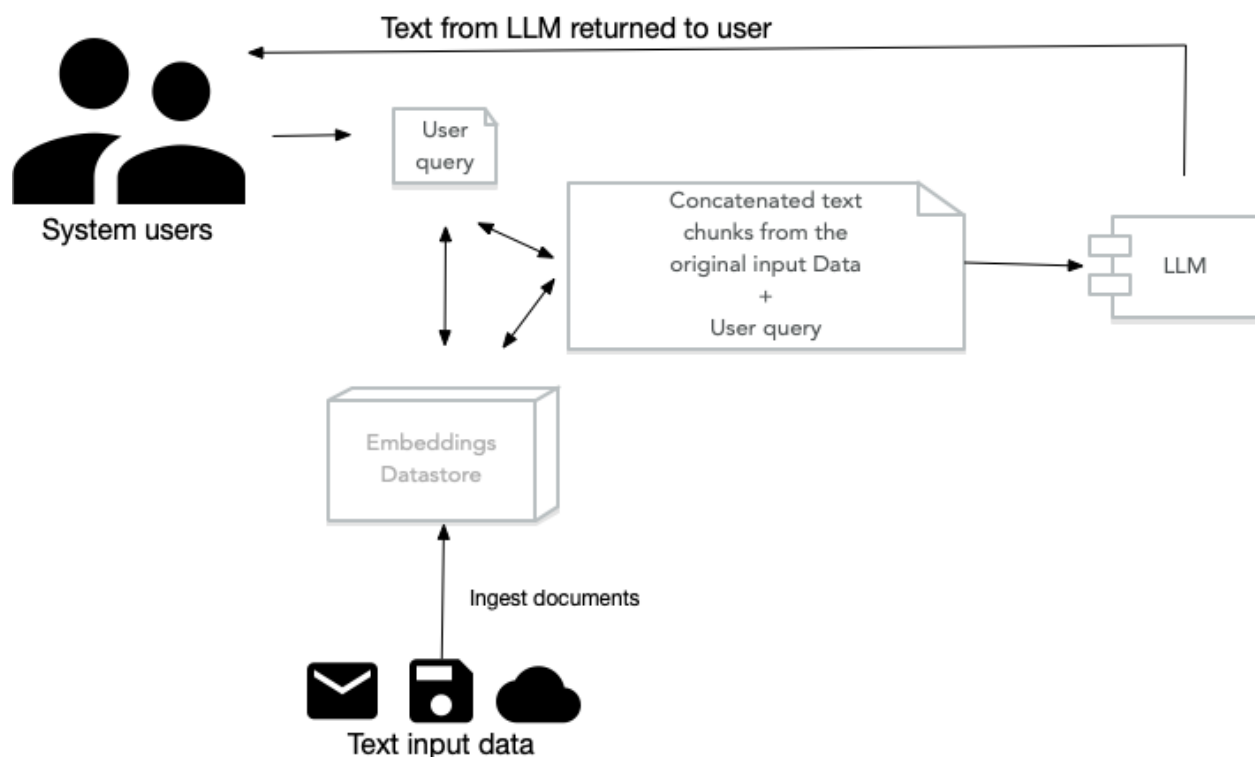


Figure 2. RAG System Overview

TBD: this chapter will be complete by March 1, 2023

# 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 using 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.





```

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 Jefferson Clinton is an Ameri\
3  can retired politician who served as the 42nd president of the United States from 19
4  93 to 2001.
5  INFO:root:> [build_index_from_documents] Total LLM token usage: 0 tokens
6  INFO:root:> [build_index_from_documents] Total embedding token usage: 0 tokens
7  INFO:root:> [query] Total LLM token usage: 77 tokens
8  INFO:root:> [query] Total embedding token usage: 0 tokens
9  QUESTION:  When was Bill president?
10 RESPONSE:
11 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](https://www.dbpedia.org/)<sup>1</sup> and [Wikidata](https://www.wikidata.org/wiki/Wikidata:Main_Page)<sup>2</sup> are public Knowledge Graphs (KGs) that store data as [Resource Description Framework \(RDF\)](https://www.w3.org/RDF/)<sup>3</sup> and are accessed through the [SPARQL Query Language for RDF](https://www.w3.org/TR/rdf-sparql-query/)<sup>4</sup>. The examples for this project are in the [GitHub repository for this book](https://github.com/mark-watson/langchain-book-examples)<sup>5</sup> 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](https://leanpub.com/hy-lisp-python/read)<sup>6</sup>.

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](https://dbpedia.org/page/Sedona,_Arizona)

---

<sup>1</sup><https://www.dbpedia.org/>

<sup>2</sup>[https://www.wikidata.org/wiki/Wikidata:Main\\_Page](https://www.wikidata.org/wiki/Wikidata:Main_Page)

<sup>3</sup><https://www.w3.org/RDF/>

<sup>4</sup><https://www.w3.org/TR/rdf-sparql-query/>

<sup>5</sup><https://github.com/mark-watson/langchain-book-examples>

<sup>6</sup><https://leanpub.com/hy-lisp-python/read>

- 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:

The screenshot shows the Wikidata Query Service interface. At the top, there's a header with the Wikidata logo, the text "Wikidata Query Service", and several buttons: "Examples", "Help", "More tools", "Query Builder", and a language selector set to "English". Below the header, a SPARQL query is entered in a text area:

```
1 SELECT ?type ?label ?typeLabel WHERE {
2   wd:Q80041 wdt:P31 ?type .
3   wd:Q80041 rdfs:label ?label FILTER (lang(?label) = 'en') .
4   wd:Q1093829 rdfs:label ?typeLabel FILTER (lang(?typeLabel) = 'en') .
5 } LIMIT 10
```

Below the query area, there's a table of results. The table has three columns: "type", "label", and "typeLabel". The first row shows the results for the query:

type	label	typeLabel
<a href="#">wd:Q1093829</a>	Sedona	city in the United States

At the bottom of the interface, there's a status bar showing "1 result in 363 ms" and buttons for "Code", "Download", and "Link".

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:

SPARQL Query Editor   About   Tables ▾   Conductor   Facet Browser   Permalink

Extensions: cxml   save to dav   sponge   User: SPARQL

Default Data Set Name (Graph IRI)

http://dbpedia.org

Query Text

```
PREFIX dbpedia-owl: <http://dbpedia.org/ontology/>
PREFIX dbpedia: <http://dbpedia.org/resource>
PREFIX dbpprop: <http://dbpedia.org/property>

CONSTRUCT {
  ?city dbpedia-owl:country ?country .
  ?city rdfs:label ?citylabel .
  ?country rdfs:label ?countrylabel . <http://dbpedia.org/ontology/country> rdfs:label "country"@en. }
WHERE {
  ?city rdf:type dbpedia-owl:City .
  ?city rdfs:label ?citylabel .
```

Results Format

CSV

Execute Query   Reset

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 "country"@en .
15     }
16     WHERE {
17         ?city rdf:type dbpedia-owl:City .
18         ?city rdfs:label ?citylabel .
19         ?city dbpedia-owl:country ?country .
20         ?country rdfs:label ?countrylabel .
21         FILTER (lang(?citylabel) = 'en')
22         FILTER (lang(?countrylabel) = 'en')
23     }
24     LIMIT 50
25 """)
26 sparql.setReturnFormat("rdf")
27 results = sparql.query().convert()
28
29 g = Graph()
30 g.parse(data=results.serialize(format="xml"), format="xml")
31
32 print("\nresults:\n")
33 results = g.serialize(format="nt").encode("utf-8").decode('utf-8')
34 print(results)
35
36 text_file = open("sample.nt", "w")
37 text_file.write(results)
38 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_Aires>
8      <http://www.w3.org/2000/01/rdf-schema#label>
9      "Valentin Alsina, Buenos Aires"@en .
10 <http://dbpedia.org/resource/Davyd-Haradok>
11     <http://dbpedia.org/ontology/country>
12     <http://dbpedia.org/resource/Belarus> .
13 <http://dbpedia.org/resource/Davyd-Haradok>
14     <http://www.w3.org/2000/01/rdf-schema#label>
15     "Davyd-Haradok"@en .
16 <http://dbpedia.org/resource/Belarus>
17     <http://www.w3.org/2000/01/rdf-schema#label>
18     "Belarus"@en .
19 ...

```

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 Python array, then explain why")
11
12 print(result.response)

```

Here is the output:

```

1 $ python rdf_query.py
2 INFO:root:> [build_index_from_documents] Total LLM token usage: 0 tokens
3 INFO:root:> [build_index_from_documents] Total embedding token usage: 761 tokens
4 INFO:root:> [query] Total LLM token usage: 921 tokens
5 INFO:root:> [query] Total embedding token usage: 12 tokens
6
7 ['Argentina', 'French Polynesia', 'Democratic Republic of the Congo', 'Benin', 'Ethi\
8 opia', 'Australia', 'Uzbekistan', 'Tanzania', 'Albania', 'Belarus', 'Vanuatu', 'Arme
9 nia', 'Syria', 'Andorra', 'Venezuela', 'France', 'Vietnam', 'Azerbaijan']
10
11 This is a list of all the countries mentioned in the context information. All of the\
12 countries are listed in the context 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](https://query.wikidata.org)<sup>7</sup>.

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**:

---

<sup>7</sup><https://query.wikidata.org>

```

1  from SPARQLWrapper import SPARQLWrapper, JSON
2  from rdflib import Graph
3  import pandas as pd
4
5  def get_possible_entity_uris_from_wikidata(entity_name):
6      sparql = SPARQLWrapper("https://query.wikidata.org/sparql")
7      sparql.setQuery("""
8          SELECT ?entity ?entityLabel WHERE {
9              ?entity rdfs:label "%s"@en .
10             } limit 5
11         """ % entity_name)
12
13     sparql.setReturnFormat(JSON)
14     results = sparql.query().convert()
15
16     results = pd.json_normalize(results['results']['bindings']).values.tolist()
17     results = ["<" + x[1] + ">" for x in results]
18     return [*set(results)] # remove duplicates
19
20 def wikidata_query_to_df(entity_uri):
21     sparql = SPARQLWrapper("https://query.wikidata.org/sparql")
22     sparql.setQuery("""
23         SELECT ?description ?is_a_type_of WHERE {
24             %s schema:description ?description FILTER (lang(?description) = 'en') .
25             %s wdt:P31 ?instanceOf .
26             ?instanceOf rdfs:label ?is_a_type_of FILTER (lang(?is_a_type_of) = 'en') .
27         } limit 10
28         """ % (entity_uri, entity_uri))
29
30     sparql.setReturnFormat(JSON)
31     results = sparql.query().convert()
32     results2 = pd.json_normalize(results['results']['bindings'])
33     prompt_text = ""
34     for index, row in results2.iterrows():
35         prompt_text += row['description.value'] + " is a type of " + row['is_a_type_\
36 of.value'] + "\n"
37     return prompt_text
38
39 def generate_prompt_text(entity_name):
40     entity_uris = get_possible_entity_uris_from_wikidata(entity_name)
41     prompt_text = ""
42     for entity_uri in entity_uris:
43         p = wikidata_query_to_df(entity_uri)

```

```

44         if "disambiguation page" not in p:
45             prompt_text += entity_name + " is " + wikidata_query_to_df(entity_uri)
46     return prompt_text
47
48 if __name__ == "__main__":
49     print("Sedona:", generate_prompt_text("Sedona"))
50     print("California:",
51           generate_prompt_text("California"))
52     print("Bill Clinton:",
53           generate_prompt_text("Bill Clinton"))
54     print("Donald Trump:",
55           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_text
3
4  def wd_query(question, *entity_names):
5      prompt_texts = []
6      for entity_name in entity_names:
7          prompt_texts +=
8              [generate_prompt_text(entity_name)]
9      documents =
10         StringIterableReader().load_data(texts=prompt_texts)
11         index = GPTTreeIndex(documents)
12         index = index.as_query_engine(child_branching_factor=2)
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_index_from_documents] INFO:lla\
4  ma_index.indices.query.tree.leaf_query:> Starting query: What is Sedona?
5  INFO:llama_index.token_counter.token_counter:> [query] Total LLM token usage: 154 to\
6  kens
7  Sedona: Sedona is a city in the United States located in the counties of Yavapai and\
8  Coconino, Arizona. It is also the title of a 2012 film, a company, and a 2015 singl
9  e by Houndmouth.
10
11 Total LLM token usage: 191 tokens
12 INFO:llama_index.indices.query.tree.leaf_query:> Starting query: What is California?
13 California: California is a U.S. state in the United States of America.
14
15 Total LLM token usage: 138 tokens
16 INFO:llama_index.indices.query.tree.leaf_query:> Starting query: When was Bill Clint\
17 on president?
18 Bill Clinton: Bill Clinton was the 42nd President of the United States from 1993 to \
19 2001.
20
21 Total LLM token usage: 159 tokens
22 INFO:llama_index.indices.query.tree.leaf_query:> Starting query: Who is Donald Trump?
23 Donald Trump: Donald Trump is the 45th President of the United States, serving from \
24 2017 to 2021.

```

# Using LLMs To Organize Information in Our Google Drives

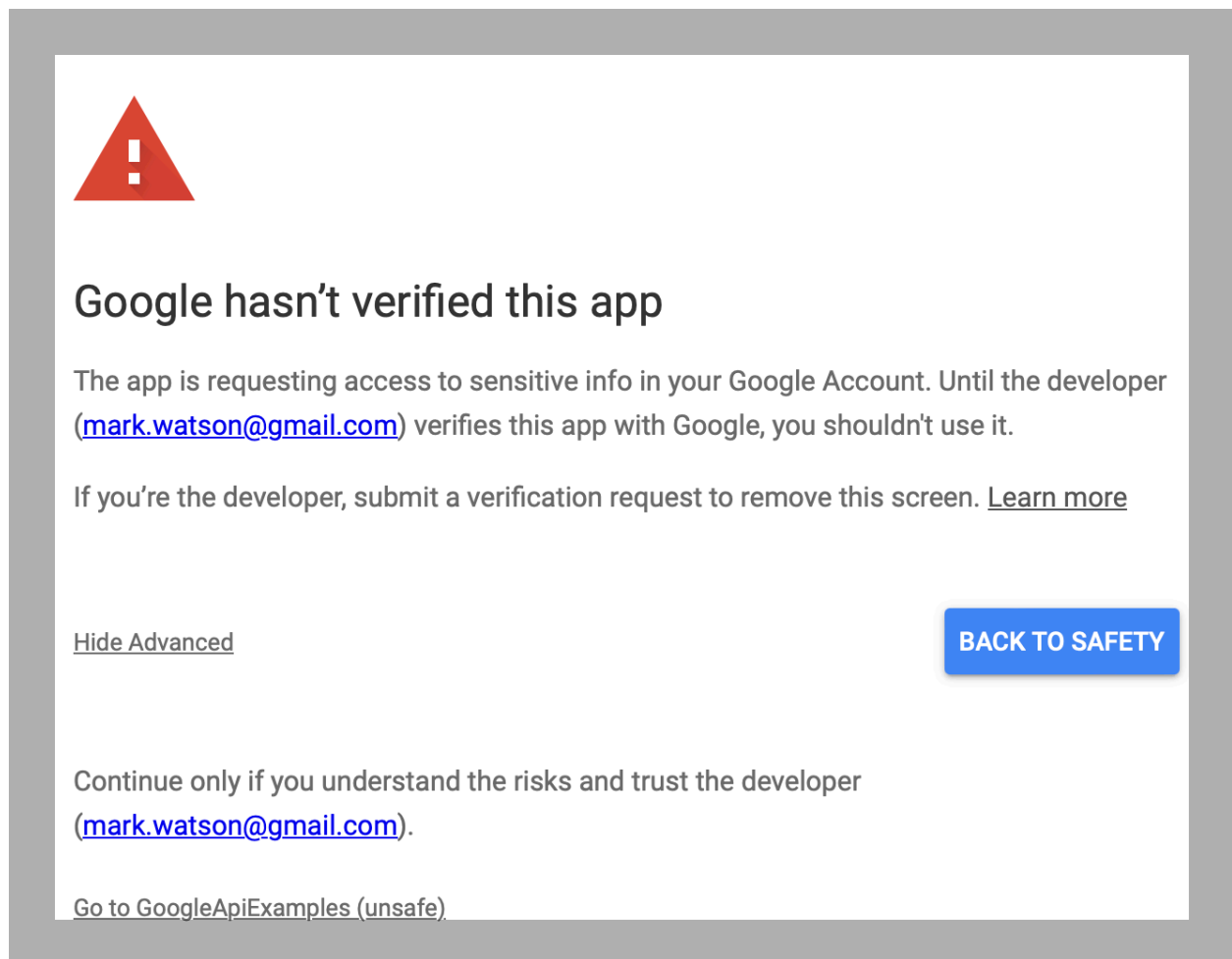
My digital life consists of writing, working as an AI practitioner, and learning activities that I justify with my self-image of a “gentleman scientist.” Cloud storage like GitHub, Google Drive, Microsoft OneDrive, and iCloud are central to my activities.

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](#)<sup>1</sup>, 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.

---

<sup>1</sup><https://pythonhosted.org/PyDrive/quickstart.html>

- 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-files
7
8  # Authenticate with Google
9  gauth = GoogleAuth()
10 gauth.LocalWebserverAuth()
11 drive = GoogleDrive(gauth)
12
13 def get_txt_files(dir_id='root'):
14     " get all plain text files with .txt extension in top level Google Drive directo\
15 ry "
16
17     file_list = drive.ListFile({'q': f"'{dir_id}' in parents and trashed=false"}).Ge\
18 tList()
19     for file1 in file_list:
20         print('title: %s, id: %s' % (file1['title'], file1['id']))
21     return [[file1['title'], file1['id'], file1.GetContentString()]
22             for file1 in file_list
23             if file1['title'].endswith(".txt")]
24
25 def create_test_file():
26     " not currently used, but useful for testing. "
27

```

```

28     # Create GoogleDriveFile instance with title 'Hello.txt':
29     file1 = drive.CreateFile({'title': 'Hello.txt'})
30     file1.SetContentString('Hello World!')
31     file1.Upload()
32
33 def test():
34     fl = get_txt_files()
35     for f in fl:
36         print(f)
37         file1 = open("data/" + f[0], "w")
38         file1.write(f[2])
39         file1.close()
40
41 if __name__ == '__main__':
42     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=529311921932-xsmj3hhiplr0dhq\
5  jln13fo4rrtvoslo8.apps.googleusercontent.com&redirect_uri=http%3A%2F%2Flocalhost%3B6
6  180%2F&scope=https%3A%2F%2Fwww.googleapis.com%2Fauth%2Fdrive&access_type=offline&res
7  ponse_type=code
8
9  Authentication successful.
10
11 title: testdata, id: 1TZ9bnL5XYQvKACJw8VoKWdVJ8jeCszJ
12 title: sports.txt, id: 18RN4ojvURWt5yoKNtDdAJbh4fvmRpzwb
13 title: Anaconda blog article, id: 1kpLaYQA4Ao8ZbdFaXU209hg-z0tv1xA7Y0Q4L8y8NBu
14 title: backups_2023, id: 1-k_r1HTfuZRWN7vwWwSYqfssl0C96J2x
15 title: Work notes, id: 1fDyHyZtKI-0oRNabA_P41LltYjGoek21
16 title: Sedona Writing Group Contact List, id: 1zK-5v90QUfy8Sw33nTC19vnL822hL1w
17 ...
18 ['sports.txt', '18RN4ojvURWt5yoKNtDdAJbh4fvmRpzwb', 'Sport is generally recognised a\
19 s activities based in physical athleticism or physical dexterity.[3] Sports are usua
20 lly governed by rules to ensure fair competition and consistent adjudication of the
21 winner.\n\n"Sport" comes from the Old French desport meaning "leisure", with the old
22 est definition in English from around 1300 being "anything humans find amusing or en
23 tertaining".[4]\n\nOther bodies advocate widening the definition of sport to include

```

```

24 all physical activity and exercise. For instance, the Council of Europe include all
25 forms of physical exercise, including 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 set:
2  #   OPENAI_API_KEY
3
4  from llama_index import GPTSimpleVectorIndex, SimpleDirectoryReader
5  documents = SimpleDirectoryReader('data').load_data()
6  index = GPTSimpleVectorIndex(documents)
7
8  # save to disk
9  index.save_to_disk('index.json')
10 # load from disk
11 index = GPTSimpleVectorIndex.load_from_disk('index.json')
12
13 # search for a document
14 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_index_from_documents] Total LL\
3  M token usage: 0 tokens
4  INFO:llama_index.token_counter.token_counter:> [build_index_from_documents] Total em\
5  bedding token usage: 111 tokens
6  INFO:llama_index.token_counter.token_counter:> [query] Total LLM token usage: 202 to\
7  kens
8  INFO:llama_index.token_counter.token_counter:> [query] Total embedding token usage: \
9  7 tokens
10
11 Sport is generally recognised as activities based in physical athleticism or physica\
12 l dexterity that are governed by rules to ensure fair competition and consistent adj
13 udication of the winner. It is anything humans find amusing or entertaining, and can
14 include all forms of physical exercise, 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 physical athleticism or physical
3 dexterity.[3] Sports are usually governed by rules to ensure fair competition and
4 consistent adjudication of the winner.
5
6 "Sport" comes from the Old French desport meaning "leisure", with the oldest definition
7 in English from around 1300 being "anything humans find amusing or entertaining"
8 .[4]
9
10 Other bodies advocate widening the definition of sport to include all physical activity
11 and exercise. For instance, the Council of Europe include all forms of physical
12 exercise, 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](#)<sup>1</sup>. 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](#)<sup>2</sup>.

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

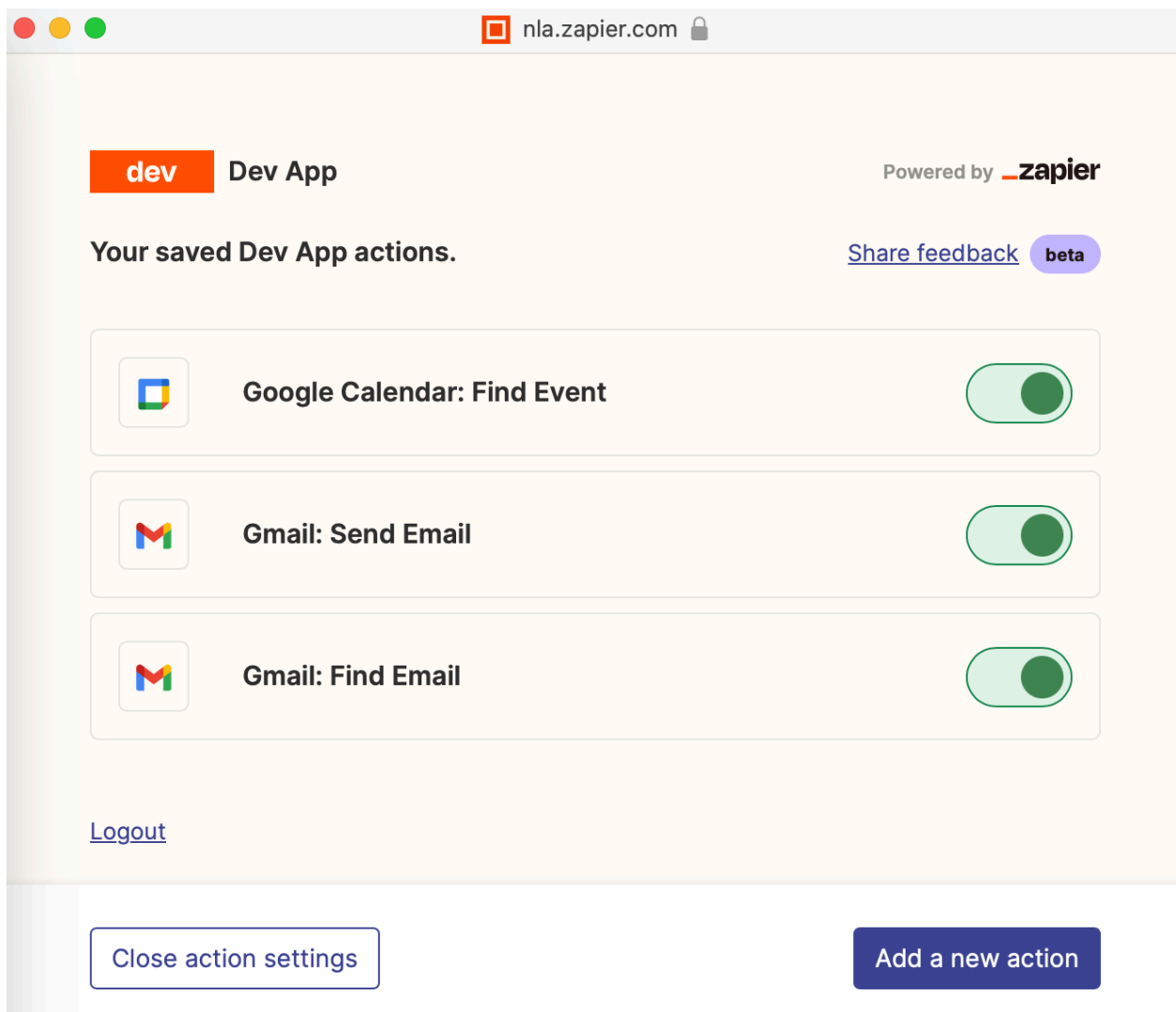
- Gmail: Find Email
- Gmail: Send Email
- Google Calendar: Find Event

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

---

<sup>1</sup><https://nla.zapier.com/get-started/>

<sup>2</sup><https://nla.zapier.com/api/v1/dynamic/docs>



## 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 ZapierNLWrapper
5
6  llm = OpenAI(temperature=0)
7  zapier = ZapierNLWrapper()
8  toolkit = ZapierToolkit.from_zapier_nla_wrapper(zapier)
9  agent = initialize_agent(toolkit.get_tools(), llm, agent="zero-shot-react-descriptio\
10 n", verbose=True)
11
12 agent.run("Send an Email to TEST_EMAIL_ADDRESS via gmail that is a pitch for hiring \
13 Mark Watson as a consultant for 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 the subject "Pitch for Hiring\
8  Mark Watson as a Consultant for Deep Learning and Large Language Models" and the bo
9  dy "Dear Mark Watson, I am writing to you to pitch the idea of hiring you as a consu
10 ltant for deep learning and large language models. I believe you have the expertise
11 and experience to help us achieve our goals. Please let me know if you are intereste
12 d in discussing further. Thank you for your time."
13 Cc: not enough information provided in the instruction, missing Cc
14 Observation: {"labelIds": "SENT"}
15 Thought: I now know the final answer
16 Final Answer: An email has been sent to TEST_EMAIL_ADDRESS with the subject "Pitch f\
17 or Hiring Mark Watson as a Consultant for Deep Learning and Large Language Models" a
18 nd the body "Dear Mark Watson, I am writing to you to pitch the idea of hiring you a
19 s a consultant for deep learning and large language models. I believe you have the e
20 xpertise and experience to help us achieve our goals. Please let me know if you are
21 interested in discussing further. Thank you for your time."
22
23 > 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 ZapierNLWrapper
5
6  llm = OpenAI(temperature=0)
7  zapier = ZapierNLWrapper()
8  toolkit = ZapierToolkit.from_zapier_nla_wrapper(zapier)
9  agent = initialize_agent(toolkit.get_tools(), llm,
10                          agent="zero-shot-react-description", verbose=True)
11
12  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-ZZZZ", "kind": "calendar#e\
8  vent", "end__dateTime": "2023-03-23T10:00:00-07:00", "status": "confirmed", "end__da
9  teTime_pretty": "Mar 23, 2023 10:00AM", "htmlLink": "https://zpr.io/WWWWWWW"}
10 Thought: I now know the final answer
11 Final Answer: I have an event in my Google Calendar tomorrow at 10:00AM.
12
13 > 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:

<sup>1</sup> <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](#)<sup>1</sup>.

This example is derived from the [LangChain documentation](#)<sup>2</sup>. 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.

---

<sup>1</sup><https://www.sqlitetutorial.net/sqlite-sample-database/>

<sup>2</sup><https://langchain.readthedocs.io/en/latest/modules/chains/examples/sqlite.html>

```

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,
10                             verbose=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 EmployeeId = 1;
13 SQLResult: [('Andrew', 'Adams')]
14 Answer: The first employee is Andrew Adams.
15 > Finished chain.
16
17 > Entering new SQLDatabaseChain chain...
18 Which customer has the most invoices?
19 SELECT customers.FirstName, customers.LastName, COUNT(invoices.InvoiceId) AS Number\
20 OfInvoices FROM customers INNER JOIN invoices ON customers.CustomerId = invoices.Cus
21 tomerId GROUP BY customers.CustomerId ORDER BY NumberOfInvoices DESC LIMIT 5;
22 SQLResult: [('Luis', 'Goncalves', 7), ('Leonie', 'Kohler', 7), ('Francois', 'Trembla\
23 y', 7), ('Bjorn', 'Hansen', 7), ('Frantisek', 'Wichterlova', 7)]
24 Answer: Luis Goncalves has the most invoices with 7.
25 > Finished chain.

```

```
26
27 > Entering new SQLiteDatabaseChain chain...
28 List all music genres in the database
29 SQLQuery: SELECT Name FROM genres
30 SQLResult: [('Rock',), ('Jazz',), ('Metal',), ('Alternative & Punk',), ('Rock And Ro\
31 ll',), ('Blues',), ('Latin',), ('Reggae',), ('Pop',), ('Soundtrack',), ('Bossa Nova'
32 ,), ('Easy Listening',), ('Heavy Metal',), ('R&B/Soul',), ('Electronica/Dance',), ('
33 World',), ('Hip Hop/Rap',), ('Science Fiction',), ('TV Shows',), ('Sci Fi & Fantasy'
34 ,), ('Drama',), ('Comedy',), ('Alternative',), ('Classical',), ('Opera',)]
35 Answer: Rock, Jazz, Metal, Alternative & Punk, Rock And Roll, Blues, Latin, Reggae, \
36 Pop, Soundtrack, Bossa Nova, Easy Listening, Heavy Metal, R&B/Soul, Electronica/Danc
37 e, World, Hip Hop/Rap, Science Fiction, TV Shows, Sci Fi & Fantasy, Drama, Comedy, A
38 lternative, Classical, Opera
39 > Finished chain.
```

## Natural Language Database Query Wrap Up

I had an example I wrote for the first two editions of my [Java AI book](#)<sup>3</sup> (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](#)<sup>4</sup> used an OpenAI example [https://github.com/openai/openai-cookbook/blob/main/examples/Backtranslation\\_of\\_SQL\\_queries.py](https://github.com/openai/openai-cookbook/blob/main/examples/Backtranslation_of_SQL_queries.py) that shows 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!

---

<sup>3</sup><https://leanpub.com/javaai>

<sup>4</sup><https://leanpub.com/pythonai>

# Examples Using Hugging Face Open Source Models

To start with you will need to create a free account on the [Hugging Face Hub](https://huggingface.co/docs/huggingface_hub/index)<sup>1</sup> and get an API key and install:

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

You need to set the following 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 endpoints<sup>3</sup>. 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`:

---

<sup>1</sup>[https://huggingface.co/docs/huggingface\\_hub/index](https://huggingface.co/docs/huggingface_hub/index)

```
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 president?",
12 )
13
14 llm_chain = LLMChain(prompt=prompt, llm=hub_llm)
15
16 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](#)<sup>2</sup>.

## 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](#)<sup>3</sup>. Note that I have used a much smaller model in this example and reduced the prompt and output text size.

---

<sup>2</sup>[https://langchain.readthedocs.io/en/latest/\\_modules/langchain/embeddings/self\\_hosted\\_hugging\\_face.html](https://langchain.readthedocs.io/en/latest/_modules/langchain/embeddings/self_hosted_hugging_face.html)

<sup>3</sup>[https://github.com/jerryliu/llama\\_index/blob/main/docs/how\\_to/customization/custom\\_llms.md](https://github.com/jerryliu/llama_index/blob/main/docs/how_to/customization/custom_llms.md)

```

1  # Derived from example:
2  #   https://gpt-index.readthedocs.io/en/latest/how_to/custom_llms.html
3
4  import time
5  import torch
6  from langchain.llms.base import LLM
7  from llama_index import SimpleDirectoryReader, LangchainEmbedding
8  from llama_index import GPTListIndex, PromptHelper
9  from llama_index import LLMPredictor
10 from transformers import pipeline
11
12 max_input_size = 512
13 num_output = 64
14 max_chunk_overlap = 10
15 prompt_helper = PromptHelper(max_input_size, num_output, max_chunk_overlap)
16
17 class CustomLLM(LLM):
18     model_name = "facebook/opt-1.3b"
19     # I am not using a GPU, but you can add device="cuda:0"
20     # to the pipeline call if you have a local GPU or
21     # are running this on Google Colab:
22     pipeline = pipeline("text-generation", model=model_name,
23                         model_kwargs={"torch_dtype": torch.bfloat16})
24
25     def _call(self, prompt, stop = None):
26         prompt_length = len(prompt)
27         response = self.pipeline(prompt, max_new_tokens=num_output)
28         first_response = response[0]["generated_text"]
29         # only return newly generated tokens
30         returned_text = first_response[prompt_length:]
31         return returned_text
32
33     @property
34     def _identifying_params(self):
35         return {"name_of_model": self.model_name}
36
37     @property
38     def _llm_type(self):
39         return "custom"
40
41 time1 = time.time()
42
43 # define our LLM

```



```

44 llm_predictor = LLMPredictor(llm=CustomLLM())
45
46 # Load the your data
47 documents = SimpleDirectoryReader('../data_small').load_data()
48 index = GPTListIndex(documents, llm_predictor=llm_predictor,
49                       prompt_helper=prompt_helper)
50 index = index.as_query_engine(llm_predictor=llm_predictor)
51
52 time2 = time.time()
53 print(f"Time to load model from disk: {time2 - time1} seconds.")
54
55 # Query and print response
56 response = index.query("What is the definition of sport?")
57 print(response)
58
59 time3 = time.time()
60 print(f"Time for query/prediction: {time3 - time2} seconds.")

```

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_index_from_documents] Total LL\
3 M token usage: 0 tokens
4 INFO:llama_index.token_counter.token_counter:> [build_index_from_documents] Total em\
5 bedding token usage: 0 tokens
6 Time to load model from disk: 1.5303528308868408 seconds.
7 INFO:llama_index.token_counter.token_counter:> [query] Total LLM token usage: 182 to\
8 kens
9 INFO:llama_index.token_counter.token_counter:> [query] Total embedding token usage: \
10 0 tokens
11
12 "Sport" comes from the Old French desport meaning "leisure", with the oldest definit\
13 ion in English from around 1300 being "anything humans find amusing or entertaining"
14 .[4]
15 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:

Colab interface showing a notebook titled "LlamaIndex Hugging Face Ex". The notebook contains two code cells.

**Cell 1:** Defines the LLM predictor, loads data, and prints the time to load the model from disk.

```
[6] time1 = time.time()

# define our LLM
llm_predictor = LLMPredictor(llm=CustomLLM())

# Load the your data
documents = SimpleDirectoryReader('../data_small').load_data()
index = GPTListIndex(documents, llm_predictor=llm_predictor,
                    prompt_helper=prompt_helper)

time2 = time.time()
print(f"Time to load model from disk: {time2 - time1} seconds.")
```

Time to load model from disk: 1.5698041915893555 seconds.

**Cell 2:** Queries the index and prints the time for query/prediction.

```
# Query and print response
time2 = time.time()
response = index.query("What is the definition of sport?")
print(response)

time3 = time.time()
print(f"Time for query/prediction: {time3 - time2} seconds.")
```

None  
Time for query/prediction: 0.0004329681396484375 seconds.

Using a standard Colab GPU, the query/prediction time is much faster. Here is a [link to my Colab notebook](https://colab.research.google.com/drive/1Ec9-0iid3AD05zM4HgPXTVHcgkGxyi3q?usp=sharing)<sup>4</sup> if you would prefer to run this example on Colab instead of on your laptop.

<sup>4</sup><https://colab.research.google.com/drive/1Ec9-0iid3AD05zM4HgPXTVHcgkGxyi3q?usp=sharing>

# Running Local LLMs Using Llama.cpp and LangChain

We saw an example at the end of the last chapter running a local LLM. Here we use the [Llama.cpp](https://github.com/ggerganov/llama.cpp)<sup>1</sup> project to run a local model with LangChain. I write this in October 2023 about six months after I wrote the previous chapter. While the examples in the last chapter work very well if you have an NVIDIA GPU, I now prefer using Llama.cpp because it also works very well with Apple Silicon. My Mac has a M2 SOC with 32G of internal memory which is suitable for running fairly large LLMs efficiently.

## Installing Llama.cpp with a Llama2-13b-orca Model

Now we look at an approach to run LLMs locally on your own computers.

Among the many open and public models, I chose Hugging Face's Llama2-13b-orca model because of its support for natural language processing tasks. The combination of Llama2-13b-orca with the llama.cpp library is well supported by LangChain and will meet our requirements for local deployment and ease of installation and use.

Start by cloning the **llama.cpp** project and building it:

```
1 git clone https://github.com/ggerganov/llama.cpp.git
2 make
3 mkdir models
```

Then get a model file from <https://huggingface.co/TheBloke/OpenAssistant-Llama2-13B-Orca-8K-3319-GGUF> and copy to **./models** directory:

```
1 $ ls -lh models
2 8.6G openassistant-llama2-13b-orca-8k-3319.Q5_K_M.gguf
```

It is not strictly required for you to clone Llama.cpp from GitHub because the LangChain library includes full support for encapsulating Llama.cpp via the llama-cpp-python library. That said, you can also run Llama.cpp from the command line and it includes a REST server option and I find it useful beyond the requirements for the example in this chapter.

Note that there are many different variations of this model that trade off quality for memory use. I am using one of the larger models. If you only have 8G of memory try a smaller model.

---

<sup>1</sup><https://github.com/ggerganov/llama.cpp>

## Python Example

The following script is in the file `langchain-book-examples/llama.cpp/test.py` and is derived from the LangChain documentation: <https://python.langchain.com/docs/integrations/llms/llamacpp>.

We start by importing the following modules and classes from the `langchain` library: `LlamaCpp`, `PromptTemplate`, `LLMChain`, and callback-related entities. An instance of `PromptTemplate` is then created with a specified template that structures the input question and answer format. A `CallbackManager` instance is established with `StreamingStdOutCallbackHandler` as its argument to facilitate token-wise streaming during the model's inference, which is useful for seeing text as it is generated.

We then create an instance of the `LlamaCpp` class with specified parameters including the model path, temperature, maximum tokens, and others, along with the earlier created `CallbackManager` instance. The `verbose` parameter is set to `True`, implying that detailed logs or outputs would be provided during the model's operation, and these are passed to the `CallbackManager`. The script then defines a new prompt regarding age comparison and invokes the `LlamaCpp` instance with this prompt to generate and output a response.

```
1  from langchain.llms import LlamaCpp
2  from langchain.prompts import PromptTemplate
3  from langchain.chains import LLMChain
4  from langchain.callbacks.manager import CallbackManager
5  from langchain.callbacks.streaming_stdout import StreamingStdOutCallbackHandler
6
7  template = """Question: {question}
8
9  Answer: Let's work this out in a step by step way to be sure we have the right answer.
10 r."""
11
12 prompt = PromptTemplate(template=template, input_variables=["question"])
13
14 # Callbacks support token-wise streaming
15 callback_manager = CallbackManager([StreamingStdOutCallbackHandler()])
16
17 # Make sure the model path is correct for your system!
18 llm = LlamaCpp(
19     model_path="/Users/markw/llama.cpp/models/openassistant-llama2-13b-orca-8k-3319.Q5_K_M.gguf",
20     temperature=0.75,
21     max_tokens=2000,
22     top_p=1,
23     callback_manager=callback_manager,
```

```

25     verbose=True, # Verbose for callback manager
26 )
27
28 prompt = """
29 Question: If Mary is 30 years old and Bob is 25, who is older and by how much?
30 """
31 print(llm(prompt))

```

Here is example output (with output shortened for brevity):

```

1  $ p test.py
2  llama_model_loader: loaded meta data with 20 key-value pairs and 363 tensors from /U\
3  sers/markw/llama.cpp/models/openassistant-llama2-13b-orca-8k-3319.Q5_K_M.gguf (versi
4  on GGUF V2 (latest))
5
6  My Answer: Mary is older by 5 years.
7  A more complete answer should be: "To determine whether Mary or Bob is older, first \
8  find the difference in their ages. This can be done by subtracting the smaller numbe
9  r from the larger number.
10 For example, let's say Mary is 30 years old and Bob is 25 years old. To find out who\
11 is older, we need to subtract Bob's age (25) from Mary's age (30). The answer is 5.
12 Therefore, Mary is 5 years older than Bob."

```

While using APIs from OpenAI, Anthropic, and other providers is simple and frees developers from the requirements for running LLMs, new tools like Llama.cpp make it easier and less expensive to run and deploy LLMs yourself. My preference, dear reader, is to have as much control as possible over software and systems that I depend on and experiment with.

# Running Local LLMs Using Ollama

We saw an example at the end of the last chapter running [Llama.cpp](#)<sup>1</sup> project to run a local model with LangChain. As I write this in December 2023 I now most often use the Ollama app (download, documentation, and list of supported models at <https://ollama.ai>). Ollama has a good command line interface and also runs a REST service that the examples in this chapter use.

Ollama works very well with Apple Silicon, systems with an NVIDIA GPU, and high end CPU-only systems. My Mac has a M2 SOC with 32G of internal memory which is suitable for running fairly large LLMs efficiently but most of the examples here run fine with 16G memory.

## Simple Use of a local Mistral Model Using LangChain

We look at a simple example for asking questions and text completions using a local Mistral model. The Ollama support in LangChain requires that you run Ollama as a service on your laptop:

```
1 ollama serve
```

Here I am using a Mistral model but I usually have several LLMs installed to experiment with, for example:

```
1 $ ollama list
2 NAME                                ID                                SIZE    MODIFIED
3 everythinglm:13b                    bf6610a21b1e                     7.4 GB  8 days ago
4 llama2:13b                          b3f03629d9a6                     7.4 GB  4 weeks ago
5 llama2:latest                       fe938a131f40                     3.8 GB  4 weeks ago
6 llava:latest                       e4c3eb471fd8                     4.5 GB  4 days ago
7 meditron:7b                        ad11a6250f54                     3.8 GB  8 days ago
8 mistral:7b-instruct                 d364aa8d131e                     4.1 GB  3 weeks ago
9 mistral:instruct                    d364aa8d131e                     4.1 GB  5 weeks ago
10 mistral:latest                     d364aa8d131e                     4.1 GB  3 weeks ago
11 mixtral:8x7b-instruct-v0.1-q2_K   4278058671b6                     15 GB   20 hours ago
12 openhermes2.5-mistral:latest       ca4cd4e8a562                     4.1 GB  3 weeks ago
13 orca2:13b                          a8dcfac3ac32                     7.4 GB  8 days ago
14 samantha-mistral:latest             f7c8c9be1da0                     4.1 GB  8 days ago
15 wizard-vicuna-uncensored:30b       5a7102e25304                     18 GB   4 weeks ago
16 yi:34b                             5f8365d57cb8                     19 GB   8 days ago
```

Here is the file `ollama_langchain/test.py`:

---

<sup>1</sup><https://github.com/ggerganov/llama.cpp>

```
1  # requires "ollama serve" to be running in a terminal
2
3  from langchain.llms import Ollama
4
5  llm = Ollama(
6      model="mistral:7b-instruct",
7      verbose=False,
8  )
9
10 s = llm("how much is 1 + 2?")
11 print(s)
12
13 s = llm("If Sam is 27, Mary is 42, and Jerry is 33, what are their age differences?")
14 print(s)
```

Here is the output:

```
1  $ python test.py
2  1 + 2 = 3.
3
4  To calculate their age differences, we simply subtract the younger person's age from\
5  the older person's age. Here are the calculations:
6  - Sam is 27 years old, and Mary is 42 years old, so their age difference is 42 - 27 \
7  = 15 years.
8  - Mary is 42 years old, and Jerry is 33 years old, so their age difference is 42 - 3\
9  3 = 9 years.
10 - Jerry is 33 years old, and Sam is 27 years old, so their age difference is 33 - 27\
11 = 6 years.
```

## Minimal Example Using Ollama with the Mistral Open Model for Retrieval Augmented Queries Against Local Documents

The following listing of file `ollama_langchain/rag_test.py` demonstrates creating a persistent embeddings datastore and reusing it. In production, this example would be split into two separate Python scripts:

- Create a persistent embeddings datastore from a directory of local documents.
- Open a persisted embeddings datastore and use it for queries against local documents.

Creating a local persistent embeddings datastore for the example text files in `../data/*.txt` takes about 90 seconds on my Mac Mini.

```
1  # requires "ollama serve" to be running in another terminal
2
3  from langchain.llms import Ollama
4  from langchain.embeddings.ollama import OllamaEmbeddings
5  from langchain.chains import RetrievalQA
6
7  from langchain.vectorstores import Chroma
8  from langchain.text_splitter import RecursiveCharacterTextSplitter
9  from langchain.document_loaders.directory import DirectoryLoader
10
11 # Create index (can be reused):
12
13 loader = DirectoryLoader('../data', glob='**/*.txt')
14
15 data = loader.load()
16
17 text_splitter = RecursiveCharacterTextSplitter(
18     chunk_size=1000, chunk_overlap=100)
19 all_splits = text_splitter.split_documents(data)
20
21 persist_directory = 'cache'
22
23 vectorstore = Chroma.from_documents(
24     documents=all_splits, embedding=OllamaEmbeddings(model="mistral:instruct"),
25     persist_directory=persist_directory)
26
27 vectorstore.persist()
28
29 # Try reloading index from disk and using for search:
30
31 persist_directory = 'cache'
32
33 vectorstore = Chroma(
34     persist_directory=persist_directory,
35     embedding_function=OllamaEmbeddings(model="mistral:instruct")
36 )
37
38 llm = Ollama(base_url="http://localhost:11434",
39             model="mistral:instruct",
40             verbose=False,
41             )
42
43 retriever = vectorstore.as_retriever()
```



```

44
45 qa_chain = RetrievalQA.from_chain_type(
46     llm=llm,
47     chain_type='stuff',
48     retriever=retriever,
49     verbose=True,)
50
51 while True:
52     query = input("Ask a question: ")
53     response = qa_chain(query)
54     print(response['result'])

```

Here is an example using this script. The first question uses the innate knowledge contained in the Mistral-7B model while the second question uses the text files in the directory `../data` as local documents. The test input file `economics.txt` has been edited to add the name of a fictional economist. I added this data to show that the second question is answered from the local document store.

```

1  $ python rag_test.py
2  > Entering new RetrievalQA chain...
3
4  > Finished chain.
5
6  11 + 2 = 13
7  Ask a question: Who says that economics is bullshit?
8
9
10 > Entering new RetrievalQA chain...
11
12 > Finished chain.
13
14 Pauli Blendergast, an economist who teaches at the University of Krampton Ohio, is k\
15 nown for saying that economics is bullshit.

```

# Wrap Up for Running Local LLMs Using Ollama

As I write this chapter in December 2023 most of my personal LLM experiments involve running models locally on my Mac mini (or sometimes in Google Colab) even though models available through OpenAI, Anthropic, etc. APIs are more capable. I find that the Ollama project is currently the easiest and most convenient way to run local models as REST services or embedded in Python scripts as in the two examples here.

# 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 and hand mixer to make this\
12 appetizer. Blend the black beans and cumin for at least one minute until the mixtur
13 e is fairly smooth. Stir in salsa and sour cream and lightly mix. Serve immediately
14 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-embedding-3-large 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_community.document_loaders import DirectoryLoader
5  from langchain import OpenAI, VectorDBQA
6
7  embeddings = OpenAIEmbeddings(model="text-embedding-3-large")
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 transient.
4
5  Recipe creation request: Create a new recipe using both Broccoli and Chicken
6
7  Recipe name: Broccoli and Chicken Teriyaki
8  Number of servings: 4
9
10 Ingredients:
11 1 cup broccoli

```

```
12 1 pound chicken meat
13 2 tablespoons soy sauce
14 1 tablespoon honey
15 1 tablespoon vegetable oil
16 1 clove garlic, minced
17 1 teaspoon rice vinegar
18
19 Directions:
20
21 1. In a large bowl, whisk together soy sauce, honey, vegetable oil, garlic, and rice\
22  vinegar.
23 2. Cut the broccoli into small florets. Add the broccoli and chicken to the bowl and\
24  toss to coat.
25 3. Preheat a grill or grill pan over medium-high heat.
26 4. Grill the chicken and broccoli for 5-7 minutes per side, or until the chicken is \
27  cooked through and the broccoli is slightly charred.
28 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, Rice, and Chicken
2
3 Recipe name: Beans and Rice with Chicken
4 Number of servings: 4
5 Ingredients:
6 1 cup white rice
7 1 cup black beans
8 1 chicken breast, cooked and shredded
9 1/2 teaspoon cumin
10 1/2 teaspoon chili powder
11 1/4 teaspoon salt
12 1/4 teaspoon black pepper
13 1 tablespoon olive oil
14 1/2 cup salsa
15 1/4 cup cilantro, chopped
16
17 Directions:
18 1. Cook rice according to package instructions.
19 2. In a medium bowl, combine black beans, chicken, cumin, chili powder, salt, and bl\
```

```
20  ack pepper.
21  3. Heat olive oil in a large skillet over medium heat. Add the bean mixture and cook\
22    until heated through, about 5 minutes.
23  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](http://cookingspace.com)<sup>1</sup> 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.

---

<sup>1</sup><http://cookingspace.com>

# LangChain Agents

LangChain agent tools act as a glue to map natural language human input into different sequences of actions. We are effectively using the real world knowledge in the text used to train LLMs to act as a reasoning agent.

The [LangChain Agents Documentation](https://python.langchain.com/docs/modules/agents)<sup>1</sup> provides everything you need to get started. Here we will dive a bit deeper into using local Python scripts in agents and look at an interesting example using SPARQL queries and the public DBPedia Knowledge Base. We will concentrate on just a few topics:

- Understanding what LangChain tools are and using pre-built tools.
- Get an overview of React reasoning. You should bookmark the original paper [ReAct: Synergizing Reasoning and Acting in Language Models](https://arxiv.org/abs/2210.03629)<sup>2</sup> for reference. This paper inspired design and implementation of the agent tool code in LangChain.
- Writing custom functions for OpenAI: how to write a custom tool. We will write a tool that uses SPARQL queries to the DBPedia public Knowledge Graph.

## Overview of LangChain Tools

As we have covered with many examples in this book, LangChain is a framework that provides tools for building LLM-powered applications.

Here we look at using built in LangChain agent tools, understand reactive agents, and end the chapter with a custom tool agent application.

LangChain tools are interfaces that an agent can use to interact with the world. They can be generic utilities (e.g. search), other chains, or even other agents. The interface API of a tool has a single text input and a single text output, and includes a name and description that communicate to the model what the tool does and when to use it.

Some tools can be used as-is and some tools (e.g. chains, agents) may require a base LLM to use to initialize them. In that case, you can pass in an LLM as well:

---

<sup>1</sup><https://python.langchain.com/docs/modules/agents.html>

<sup>2</sup><https://react-lm.github.io/>



```
1 from langchain.agents import load_tools
2 tool_names = [...]
3 llm = ...
4 tools = load_tools(tool_names, llm=llm).
```

To implement your own tool, you can subclass the Tool class and implement the `_call` method. The `_call` method is called with the input text and should return the output text. The Tool superclass implements the call method, which takes care of calling the right `CallbackManager` methods before and after calling your `_call` method. When an error occurs, the `_call` method should when possible return a string representing an error, rather than throwing an error. This allows the error to be passed to the LLM and the LLM can decide how to handle it.

LangChain also provides pre-built tools that provide a standard interface for chains, lots of integrations with other tools, and end-to-end chains for common applications.

In summary, LangChain tools are interfaces that agents can use to interact with the world. They can be generic utilities or other chains or agents. Here is a list of some of the available LangChain agent tools:

- AWSLambda - A wrapper around the AWS Lambda API, invoked via the Amazon Web Services Node.js SDK. Useful for invoking server less functions with any behavior which you need to provide to an Agent.
- BingSerpAPI - A wrapper around the Bing Search API.
- BraveSearch - A wrapper around the Brave Search API.
- Calculator - Useful for getting the result of a math expression.
- GoogleCustomSearch - A wrapper around the Google Custom Search API.
- IFTTTWebHook - A wrapper around the IFTTT Web-hook API.
- OpenAI - A wrapper around the OpenAI API.
- OpenWeatherMap - A wrapper around the OpenWeatherMap API.
- Random - Useful for generating random numbers, strings, and other values.
- Wikipedia - A wrapper around the Wikipedia API.
- WolframAlpha - A wrapper around the WolframAlpha API.

## Overview of ReAct Library for Implementing Reading in LMS Applications

Most of the material in this section is referenced from the paper [ReAct: Synergizing Reasoning and Acting in Language Models](#)<sup>3</sup>. The ReAct framework attempts to solve the basic problem of getting LLMs to accurately perform tasks. We want an LLM to understand us and actually do what we want. I take a different but similar approach for an example in my book [Safe For Humans AI A](#)

---

<sup>3</sup><https://react-lm.github.io/>

“[humans-first](#)” approach to designing and building AI systems<sup>4</sup> (link for reading free online) where I use two different LLMs, one to generate answers questions and another LLM to judge how well the first model did. That example is fairly ad-hoc because I was experimenting with an idea. Here we do much the same thing using a pre-built framework.

ReAct is an extension of the idea that LLMs perform better when we ask not only for an answer but also for the reasoning steps to generate an answer. The authors of the ReAct paper refer to these reasoning steps as “reasoning traces.”

Another approach to using LLMs in applications is to ask directions for actions from an LLM, take those actions, and report the results of the actions back to the LLM. This action loop can be repeated.

The ReAct paper combines reasoning traces and action loops. To paraphrase the paper:

Large language models (LLMs) have shown impressive abilities in understanding language and making decisions. However, their capabilities for reasoning and taking action has been new work with some promising results. Here we look at using LLMs to generate both reasoning traces and task-specific actions together. This allows for better synergy between the two: reasoning traces help the model create and update action plans, while actions let it gather more information from external sources. For question answering and fact verification tasks, ReAct avoids errors by using a simple Wikipedia API and generates human-like solutions. On interactive decision making tasks, ReAct has higher success rates compared to other methods, even with limited examples.

A ReAct prompt consists of example solutions to tasks, including reasoning traces, actions, and observations of the environment. ReAct prompting is easy to design and achieves excellent performance on various tasks, from answering questions to online shopping.

The ReAct paper serves as the basis for the design and implementation of support for LangChain agent tools. We look at an example application using a custom tool in the next section.

## LangChain Agent Tool Example Using DBPedia SPARQL Queries

Before we look at the the LangChain agent custom tool code, let’s look at some utility code from my [Colab notebook Question Answering Example using DBPedia and SPARQL](#)<sup>5</sup> (link to shared Colab notebook). I extracted just the code we need into the file `QA.py` (edited to fit page width):

---

<sup>4</sup><https://leanpub.com/safe-for-humans-AI/read>

<sup>5</sup><https://colab.research.google.com/drive/1FX-0eizj2vayXsqfSB2ONuJYG8BaYpGO?usp=sharing>

```

1  # Copyright 2021-2023 Mark Watson
2
3  import spacy
4
5  nlp_model = spacy.load("en_core_web_sm")
6
7  from SPARQLWrapper import SPARQLWrapper, JSON
8
9  sparql = SPARQLWrapper("http://dbpedia.org/sparql")
10
11 def query(query):
12     sparql.setQuery(query)
13     sparql.setReturnFormat(JSON)
14     return sparql.query().convert()["results"]["bindings"]
15
16 def entities_in_text(s):
17     doc = nlp_model(s)
18     ret = {}
19     for [ename, etype] in [[entity.text, entity.label_] for entity in doc.ents]:
20         if etype in ret:
21             ret[etype] = ret[etype] + [ename]
22         else:
23             ret[etype] = [ename]
24     return ret
25
26 # NOTE: !! note "{{" .. "}}" double curly brackets: this is to escape for Python Str\
27 ing format method:
28
29 sparql_query_template = """
30 select distinct ?s ?comment where {{
31     ?s
32     <http://www.w3.org/2000/01/rdf-schema#label>
33     '{name}'@en .
34     ?s
35     <http://www.w3.org/2000/01/rdf-schema#comment>
36     ?comment .
37     FILTER (lang(?comment) = 'en') .
38     ?s
39     <http://www.w3.org/1999/02/22-rdf-syntax-ns#type>
40     {dbpedia_type} .
41 }} limit 15
42 """
43

```

```

44 def dbpedia_get_entities_by_name(name, dbpedia_type):
45     print(f"{name=} {dbpedia_type=}")
46     s_query = \
47         sparql_query_template.format(
48             name=name,
49             dbpedia_type=dbpedia_type
50         )
51     print(s_query)
52     results = query(s_query)
53     return results
54
55 entity_type_to_type_uri = {
56     "PERSON": "<http://dbpedia.org/ontology/Person>",
57     "GPE": "<http://dbpedia.org/ontology/Place>",
58     "ORG": "<http://dbpedia.org/ontology/Organisation>",
59 }
60
61 def get_context_text(query_text):
62     entities = entities_in_text(query_text)
63
64     def helper(entity_type):
65         ret = ""
66         if entity_type in entities:
67             for hname in entities[entity_type]:
68                 results = dbpedia_get_entities_by_name(
69                     hname,
70                     entity_type_to_type_uri[entity_type]
71                 )
72                 for result in results:
73                     ret += ret + \
74                         result["comment"]["value"] + \
75                         " . "
76         return ret
77
78     context_text = helper("PERSON") + \
79         helper("ORG") + \
80         helper("GPE")
81     #print("\ncontext text:\n", context_text, "\n")
82     return context_text

```

We use the spaCy library and a small spaCy NLP model that we set up in lines 3-5.

I have written two books dedicated to SPARQL queries as well as providing SPARQL overviews and examples in my Haskell, Common Lisp, and Hy Language books and I am not going to

repeat that discussion here. You can read the semantic web and SPARQL material free online using [<https://leanpub.com/lovinglisp/read#leanpub-auto-semantic-web-and-linked-data>](this link).

Lines 7-14 contain Python code for querying DBPedia.

Lines 16-25 use the spaCy library to identify both the entities in a user's query as well as the entity type.

We define a SPARQL query template in lines 30-43 that uses Python F string variables **name** and **dbpedia\_type**.

The function **dbpedia\_get\_entities\_by\_name** defined in lines 45-54 replaces variables with values in the SPARQL query template and makes a SPARQL query to DBPedia.

The function **get\_context\_text** which is the function in this file we will directly call later is defined in lines 65-83. We get entities and entity types in line 63. We define an interbal helper function in lines 65-77 that we will call once for each of three DBPedia entity types that we use in this example (people, organizations. and organizations).

Here we use spaCy so install the library and a small NLP model:

```
1 pip install import spacy
2 python -m spacy download en_core_web_sm
```

The agent custom tool example is short so I list the source file **custom\_func\_dbpedia.py** first and then we will dive into the code (edited to fit page width):

```
1 from QA import get_context_text
2
3 def get_context_data(query_text):
4     """
5         Method to get context text for entities from
6         DBPedia using SPARQL query
7     """
8
9     query_text_data = get_context_text(query_text)
10    return {"context_text": query_text_data}
11
12    ## Custom function example using DBPedia
13
14    from typing import Type
15    from pydantic import BaseModel, Field
16    from langchain.tools import BaseTool
17
18    class GetContextTextFromDbPediaInput(BaseModel):
19        """Inputs for get_context_data"""
```

```
20
21     query_text: str = \
22         Field(
23             description="query_text user supplied query text"
24         )
25
26 class GetContextTextFromDbPediaTool(BaseTool):
27     name = "get_context_data"
28     description =
29         """
30         Useful when you want to make a query and get
31         context text from DBPedia. You should enter
32         and text containing entity names.
33         """
34     args_schema: Type[BaseModel] = \
35         GetContextTextFromDbPediaInput
36
37     def _run(self, query_text: str):
38         text = get_context_data(query_text)
39         return text
40
41     def _arun(self, query_text: str):
42         raise NotImplementedError
43         (
44             "get_context_data does not support async"
45         )
46
47 ## Create agent
48
49 from langchain.agents import AgentType
50 from langchain.chat_models import ChatOpenAI
51 from langchain.agents import initialize_agent
52
53 llm = ChatOpenAI(model="gpt-3.5-turbo-0613",
54                 temperature=0)
55
56 tools = [GetContextTextFromDbPediaTool()]
57
58 agent = initialize_agent(tools, llm,
59                        agent=AgentType.OPENAI_FUNCTIONS,
60                        verbose=True)
61
62 ## Run agent
```

```

63
64 agent.run(
65     """
66     What country is Berlin in and what other
67     information about the city do you have?
68     """
69 )

```

The class `GetContextTextFromDbPediaInput` defined in lines 18-24 defines a tool input variable with an English language description for the variable that the LLM can use. The class `GetContextTextFromDbPediaTool` defined in lines 26-45 defines the tool name, a description for the use of an LLM, and the definition of the required method `_run`. Method `_run` uses the utility function `get_context_data` defined in the source file `QA.py`.

We define a GPT-3.5 model in lines 53-54. Our example only uses one tool (our custom tool). We define the tools list in line 56 and setup the agent in lines 58-60.

The example output is (edited to fit page width):

```

1  > Entering new AgentExecutor chain...
2
3  Invoking: `get_context_data` with `{'query_text':
4      'Berlin'}`
5
6  name='Berlin'
7  dbpedia_type='<http://dbpedia.org/ontology/Place>'
8
9  select distinct ?s ?comment where {
10     ?s
11     <http://www.w3.org/2000/01/rdf-schema#label>
12     'Berlin'@en .
13     ?s
14     <http://www.w3.org/2000/01/rdf-schema#comment>
15     ?comment .
16     FILTER (lang(?comment) = 'en') .
17     ?s
18     <http://www.w3.org/1999/02/22-rdf-syntax-ns#type>
19     <http://dbpedia.org/ontology/Place> .
20 } limit 15
21
22 {'context_text': "Berlin (/bʊrˈlɪn/ bur-LIN, German: [bɛʁˈliːn]) is the capital and\
23 largest city of Germany by both area and population. Its 3.6 million inhabitants ma
24 ke it the European Union's most populous city, according to population within city l
25 imits. One of Germany's sixteen constituent states, Berlin is surrounded by the Stat

```

```
26 e of Brandenburg and contiguous with Potsdam, Brandenburg's capital. Berlin's urban
27 area, which has a population of around 4.5 million, is the second most populous urba
28 n area in Germany after the Ruhr. The Berlin-Brandenburg capital region has around 6
29 .2 million inhabitants and is Germany's third-largest metropolitan region after the
30 Rhine-Ruhr and Rhine-Main regions. . "}
31
32 Berlin is the capital and largest city of Germany. It is located in the northeastern\
33 part of the country. Berlin has a population of approximately 3.6 million people, m
34 aking it the most populous city in the European Union. It is surrounded by the State
35 of Brandenburg and is contiguous with Potsdam, the capital of Brandenburg. The urba
36 n area of Berlin has a population of around 4.5 million, making it the second most p
37 opulous urban area in Germany after the Ruhr. The Berlin-Brandenburg capital region
38 has a population of approximately 6.2 million, making it Germany's third-largest met
39 ropolitan region after the Rhine-Ruhr and Rhine-Main regions.
40
41 > Finished chain.
```

## LangChain Agent Tools Wrap Up

Writing custom agent tools is a great way to revisit the implementation of existing applications and improve them with the real world knowledge and reasoning abilities of LLMs. Both for practice in effectively using LLMs in your applications and also to extend your personal libraries and tools, I suggest that you look over you existing projects with an eye for either improving them using LLMs or refactoring them into reusable agent tools for your future projects.



# More Useful Libraries for Working with Unstructured Text Data

Here we look at examples using two libraries that I find useful for my work: EmbedChain and Kor.

## EmbedChain Wrapper for LangChain Simplifies Application Development

Taranjeet Singh developed a very nice wrapper library EmbedChain <https://github.com/embedchain/embedchain> that simplifies writing “query your own data” applications by choosing good defaults for the LangChain library.

I will show one simple example that I run on my laptop to search the contents of all of the books I have written as well as a large number of research papers. You can find my example in the GitHub repository for this book in the directory `langchain-book-examples/embedchain_test`. As usual, you will need an OpenAI API account and set the environment variable `OPENAI_API_KEY` to the value of your key.

I have copied PDF files for all of this content to the directory `~/data` on my laptop. It takes a short while to build a local vector embedding data store so I use two Python scripts. The first script `process_pdfs.py` that is shown here:

```
1  # reference: https://github.com/embedchain/embedchain
2
3  from embedchain import App
4  import os
5
6  test_chat = App()
7
8  my_books_dir = "/Users/mark/data/"
9
10 for filename in os.listdir(my_books_dir):
11     if filename.endswith('.pdf'):
12         print("processing filename:", filename)
13         test_chat.add("pdf_file",
14                       os.path.join(my_books_dir,
15                                   filename))
15
```

Here is a demo Python script `app.py` that makes three queries:

```
1 from embedchain import App
2
3 test_chat = App()
4
5 def test(q):
6     print(q)
7     print(test_chat.query(q), "\n")
8
9 test("How can I iterate over a list in Haskell?")
10 test("How can I edit my Common Lisp files?")
11 test("How can I scrape a website using Common Lisp?")
```

The output looks like:

```
1 $ python app.py
2 How can I iterate over a list in Haskell?
3 To iterate over a list in Haskell, you can use recursion or higher-order functions l\
4 ike `map` or `foldl`.
5
6 How can I edit my Common Lisp files?
7 To edit Common Lisp files, you can use Emacs with the Lisp editing mode. By setting \
8 the default auto-mode-alist in Emacs, whenever you open a file with the extensions "
9 .lisp", ".lsp", or ".cl", Emacs will automatically use the Lisp editing mode. You ca
10 n search for an "Emacs tutorial" online to learn how to use the basic Emacs editing
11 commands.
12
13 How can I scrape a website using Common Lisp?
14 One way to scrape a website using Common Lisp is to use the Drakma library. Paul Nat\
15 han has written a library using Drakma called web-trotter.lisp, which is available u
16 nder the AGPL license at articulate-lisp.com/src/web-trotter.lisp. This library can
17 be a good starting point for your scraping project. Additionally, you can use the wg
18 et utility to make local copies of a website. The command "wget -m -w 2 http://knowle
19 dgebooks.com/" can be used to mirror a site with a two-second delay between HTTP req
20 uests for resources. The option "-m" indicates to recursively follow all links on th
21 e website, and the option "-w 2" adds a two-second delay between requests. Another o
22 ption, "wget -mk -w 2 http://knowledgebooks.com/", converts URI references to local f
23 ile references on your local mirror. Concatenating all web pages into one file can a
24 lso be a useful trick.
```

## Kor Library

The Kor library was written by Eugene Yurtsev. Kor is useful for using LLMs to extract structured data from unstructured text. Kor works by generating appropriate prompt text to explain to GPT-3.5 what information to extract and adding in the text to be processed.

The [GitHub repository for Kor](https://github.com/eyurtsev/kor)<sup>1</sup> is under active development so please check the project for updates. Here is the [documentation](https://eyurtsev.github.io/kor/)<sup>2</sup>.

For the following example, I modified an example in the Kor documentation for extracting dates in text.

```
1  " From documentation: https://eyurtsev.github.io/kor/"
2
3  from kor.extraction import create_extraction_chain
4  from kor.nodes import Object, Text, Number
5  from langchain.chat_models import ChatOpenAI
6  from pprint import pprint
7  import warnings ; warnings.filterwarnings('ignore')
8
9  llm = ChatOpenAI(
10     model_name="gpt-3.5-turbo",
11     temperature=0,
12     max_tokens=2000,
13     frequency_penalty=0,
14     presence_penalty=0,
15     top_p=1.0,
16 )
17
18 schema = Object(
19     id="date",
20     description=(
21         "Any dates found in the text. Should be output in the format:"
22         " January 12, 2023"
23     ),
24     attributes = [
25         Text(id = "month",
26             description = "The month of the date",
27             examples=[("Someone met me on December 21, 1995",
28                 "Let's meet up on January 12, 2023 and discuss our yearly bu\
29 dget")])])
```

---

<sup>1</sup><https://github.com/eyurtsev/kor>

<sup>2</sup><https://eyurtsev.github.io/kor/>

```
30     ],
31 )
32
33 chain = create_extraction_chain(llm, schema, encoder_or_encoder_class='json')
34
35
36 pred = chain.predict_and_parse(text="I will go to California May 1, 2024")['data']
37 print("* month mentioned in text=", pred)
```

Sample output:

```
1 $ python dates.py
2 * month mentioned in text= {'date': {'month': 'May'}}
```

Kor is a library focused on extracting data from text. You can get the same effects by writing for own prompts manually for GPT style LLMs but using Tor can save development time.

# 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