# Retrieval Augmented Generation

## Indexing data

Let's load and split the pdf file of [Element of Statistical Learning](#):

```
1.  from langchain.document_loaders import PyPDFLoader
2.  from langchain.text_splitter import RecursiveCharacterTextSplitter
3.
4.  file_path = '...'
5.
6.  loader = PyPDFLoader(file_path=file_path)
7.
8.  text_splitter = RecursiveCharacterTextSplitter(
9.      chunk_size=500,
10.     chunk_overlap=0
11. )
12.
13. data = loader.load_and_split(text_splitter=text_splitter)
14. data
```

```
[Document(page_content='Springer Series in Statistics\nTrevor Hastie\nRob
ert TibshiraniJerome FriedmanSpringer Series in Statistics\nThe Elements
of\nStatistical Learning\nData Mining, Inference, and Prediction\nThe Ele
ments of Statistical LearningDuring the past decade there has been an exp
losion in computation and information tech-', metadata={'source': '/User
s/damienbenveniste/Projects/Teaching/Introduction_Langchain/data/mixed_da
ta/element_of_SL.pdf', 'page': 0}),
 Document(page_content='nology. With it have come vast amounts of data in
a variety of fields such as medicine, biolo-gy, finance, and marketing. T
he challenge of understanding these data has led to the devel-opment of n
ew tools in the field of statistics, and spawned new areas such as data m
ining,machine learning, and bioinformatics. Many of these tools have comm
on underpinnings butare often expressed with different terminology. This
book describes the important ideas inthese areas in a common conceptual f
ramework.', metadata={'source': '/Users/damienbenveniste/Projects/Teachin
g/Introduction_Langchain/data/mixed_data/element_of_SL.pdf', 'page': 0}),
 Document(page_content='While the approach is statistical, theemphasis is
on concepts rather than mathematics. Many examples are given, with a libe
```

Let's install the [FAISS](#) package:

```
1.  %pip install faiss-cpu
```

And let's embed text with the OpenAI embeddings

```
1.  from langchain.embeddings.openai import OpenAIEmbeddings
2.
3.  embeddings = OpenAIEmbeddings(show_progress_bar=True)
4.  vector1 = embeddings.embed_query('How are you?')
5.  len(vector1)
6.
7.  > 1536
```

Let's embed the book data

# Retrieval Augmented Generation

```python
1. from langchain.vectorstores import FAISS
2.
3. index = FAISS.from_documents(data, embeddings)
```

We can search in that index:

```python
1. index.similarity_search_with_relevance_scores(
2.     "What is machine learning?"
3. )
```

[(Document(page_content='This is page 1\nPrinter: Opaque this\n1\nIntroduction\nStatistical learning plays a key role in many areas of science, finance and\nindustry. Here are some examples of learning problems:\n•Predict whether a patient, hospitalized due to a heart attac k, will\nhave a second heart attack. The prediction is to be based on de mo-\ngraphic, diet and clinical measurements for that patient.\n•Predict the price of a stock in 6 months from now, on the basis o f\ncompany performance measures and economic data.', metadata={'source': '/Users/damienbenveniste/Projects/Teaching/Introduction_Langchain/data/mixed_data/element_of_SL.pdf', 'page': 19}),
  0.7547787193298542),
 (Document(page_content='This is page 389\nPrinter: Opaque this\n11\nNeural Networks\n11.1 Introduction\nIn this chapter we describe a class of learning methods that w as developed\nseparately in different fields—statistics and artificial inte lligence—based\non essentially identical models. The central idea is to extr act linear com-\nbinations of the inputs as derived features, and then model t he target as\na nonlinear function of these features. The result is a power ful learning', metadata={'source': '/Users/damienbenveniste/Projects/Teaching/Introduction_Langchain/data/mixed_data/element_of_SL.pdf', 'page': 407}),
  0.7512151611414096),
 (Document(page_content='While the approach is statistical, theemphasis is on concepts rather than mathematics. Many examples are given, with a liberaluse of color graphics. It should be a valuable resource for statisti

We can use that index in a chain

```python
1. from langchain.chains import RetrievalQA
2. from langchain.chat_models import ChatOpenAI
3. from langchain.callbacks import StdOutCallbackHandler
4.
5. retriever = index.as_retriever()
6. retriever.search_kwargs['fetch_k'] = 20
7. retriever.search_kwargs['maximal_marginal_relevance'] = True
8. retriever.search_kwargs['k'] = 10
9.
10. llm = ChatOpenAI()
11.
12. chain = RetrievalQA.from_chain_type(
13.     llm=llm,
14.     retriever=retriever,
```

```
15.     verbose=True
16. )
17.
18. handler = StdOutCallbackHandler()
19.
20. chain.run(
21.     'What is machine learning?',
22.     callbacks=[handler]
23. )
```

*Machine learning is a field of study that involves the development of algorithms and models that can learn from data and make predictions or decisions without being explicitly programmed. It focuses on creating computer systems that can automatically learn and improve from experience, rather than being explicitly programmed for specific tasks. Machine learning algorithms analyze large amounts of data to identify patterns, make predictions, or learn from examples and feedback. It is widely used in various fields such as science, finance, and industry for tasks like predicting stock prices, medical diagnoses, and customer behavior analysis.*

---------------------------------------------------------------------------------------------------------------------------------

## Loading data into a Vector Database

We are going to load the data in [Pinecone](). Let's install the Python package

```
1. %pip install pinecone-client
```

And let's load the data into the database

```
1. import pinecone
2. from langchain.vectorstores import Pinecone
3.
4. pinecone.init(
5.     api_key=PINECONE_API_KEY,  # find at app.pinecone.io
6.     environment=PINECONE_ENV  # next to api key in console
7. )
8.
9. index_name = "langchain-demo"
10. db = Pinecone.from_documents(
11.     data,
12.     embeddings,
13.     index_name=index_name
14. )
```

And we can now augment an LLM with the database

```
1. chain = RetrievalQA.from_chain_type(
2.     llm=llm,
3.     retriever=db.as_retriever(),
4.     verbose=True
5. )
6.
7. chain.run(
```

```
 8.        'What is machine learning?',
 9.        callbacks=[handler]
10. )
```

---------------------------------------------------------------------------------------------------------------------------------

## Providing sources

I will show to provide the sources as we answer questions. Let's install the NewsAPI Python package:

```
1. %pip install newsapi-python
```

Let's get the news about "Artificial Intelligence" from the past week:

```
 1. from datetime import date, timedelta
 2. from newsapi import NewsApiClient
 3.
 4. newsapi = NewsApiClient(api_key=NEWS_API_KEY)
 5.
 6. today = date.today()
 7. last_week = today - timedelta(days=7)
 8.
 9. latest_news = newsapi.get_everything(
10.     q='artificial intelligence',
11.     from_param=last_week.strftime("%Y-%m-%d"),
12.     to=today.strftime("%Y-%m-%d"),
13.     sort_by='relevancy',
14.     language='en'
15. )
```

and let's create documents:

```
1. from langchain.docstore.document import Document
2. docs = [
3.     Document(
4.         page_content=article['title'] + '\n\n' + article['description'],
5.         metadata={
6.             'source': article['url'],
7.         }
8.     ) for article in latest_news['articles']
9. ]
```

Let's create a chain that provides the sources with the answers

```
1. from langchain.chains import create_qa_with_sources_chain
```

# Retrieval Augmented Generation

```python
2. from langchain.chains.combine_documents.stuff import StuffDocumentsChain
3. from langchain.prompts import PromptTemplate
4.
5. qa_chain = create_qa_with_sources_chain(llm)
6.
7. doc_prompt = PromptTemplate(
8.     template="Content: {page_content}\nSource: {source}",
9.     input_variables=["page_content", "source"],
10. )
11.
12. final_qa_chain = StuffDocumentsChain(
13.     llm_chain=qa_chain,
14.     document_variable_name="context",
15.     document_prompt=doc_prompt,
16. )
17.
18. index = FAISS.from_documents(docs, embedding=embeddings)
19.
20.
21. chain = RetrievalQA(
22.     retriever=index.as_retriever(),
23.     combine_documents_chain=final_qa_chain
24. )
```

Let's ask a question:

```python
1. question = """
2. What is the most important news about artificial intelligence from last week?
3. """
4.
5. answer = chain.run(question)
6. answer
```

```
1. {
2.   "answer": "The most important news about artificial intelligence from last
      week is the use of AI to train on the works of authors Stephen King and
      Margaret Atwood. These authors responded to the revelation that their work is
      being used to train AI. Additionally, AI took the stage at the Edinburgh
      Fringe festival, raising the question of whether AI can deliver a satisfying
      punchline. Furthermore, a tech expert from the University of Oxford
      highlighted the potential workplace threats of AI, including the possibility
      of AI becoming a monitoring boss. Finally, AI is being seen as a tool that can
      help companies connect with customers in a more personalized and efficient
      way.",
3.   "sources": [
4.     "https://www.theatlantic.com/newsletters/archive/2023/09/books-briefing-
      ai-stephen-king-margaret-atwood/675213/?utm_source=feed",
5.     "https://www.cnet.com/tech/ai-took-the-stage-at-the-worlds-largest-arts-
      festival-heres-what-happened/",
6.     "https://www.foxnews.com/tech/tech-expert-existential-fears-ai-are-
      overblown-sees-very-disturbing-workplace-threats",
7.     "https://www.techradar.com/pro/ai-could-help-companies-connect-with-
      customers-like-never-before"
8.   ]
9. }
```

# Retrieval Augmented Generation

--------------------------------------------------------------------------------------------------------

## Indexing a website

We are going to use Apify to crawl a website. Let's download the Python package

```
1. %pip install apify-client chromadb
```

and let's create a loader that will crawl the AiEdge Newsletter website:

```python
1. from langchain.utilities import ApifyWrapper
2. from langchain.document_loaders.base import Document
3.
4. apify = ApifyWrapper()
5.
6. loader = apify.call_actor(
7.     actor_id="apify/website-content-crawler",
8.     run_input={
9.         "startUrls": [{"url": "https://newsletter.theaiedge.io/"}],
10.        "aggressivePrune": True,
11.    },
12.    dataset_mapping_function=lambda item: Document(
13.        page_content=item["text"] or "", metadata={"source": item["url"]}
14.    ),
15. )
```

Let's index the website data:

```python
1. from langchain.indexes import VectorstoreIndexCreator
2.
3. text_splitter = RecursiveCharacterTextSplitter(
4.     chunk_size=500,
5.     chunk_overlap=0
6. )
7.
8. index = VectorstoreIndexCreator(
9.     text_splitter=text_splitter
10. ).from_loaders([loader])
```

Let's make a search on that index:

```python
1. query = "What is the main subject of the aiedge newsletter?"
2.
3. index.query_with_sources(query)
1. {
2.     'question': 'What is the main subject of the aiedge newsletter?',
3.     'answer': ' The main subject of the AiEdge newsletter is Machine Learning
   applications, Machine Learning System Design, MLOps, and the latest techniques
   and news about the field.\n',
4.     'sources': ''
```

```
5. }
```

Let's now pass that index to a chain

```
1. retriever = index.vectorstore.as_retriever()
2.
3. qa = RetrievalQA.from_chain_type(
4.     llm=llm,
5.     retriever=retriever,
6. )
7.
8. query = "What is the most recent article of the aiedge newsletter?"
9.
10. qa.run(
11.     query,
12.     callbacks=[handler]
13. )
```

*"I'm sorry, but I don't have access to the specific articles or the most recent content of the AiEdge Newsletter. As an AI language model, I don't have real-time access to current articles or newsletters. It would be best to subscribe to the newsletter and check the latest edition for the most recent article."*

---------------------------------------------------------------------------------------------------------------------------

## Indexing a GitHub repo

Let's install the Python package:

```
1. %pip install GitPython
```

Let's load a repo

```
1. from langchain.document_loaders import GitLoader
2.
3. loader = GitLoader(
4.     clone_url="https://github.com/langchain-ai/langchain",
5.     repo_path="./data/repo/",
6.     file_filter=lambda file_path: file_path.endswith(".py"),
7.     branch="master",
8. )
9.
10. documents = loader.load()
```

Let's resplit the documents for the Python language:

```
1. from langchain.text_splitter import Language
2.
3. python_splitter = RecursiveCharacterTextSplitter.from_language(
4.     language=Language.PYTHON,
5.     chunk_size=1000,
```

```
6.        chunk_overlap=200
7.  )
8.
9.  documents = python_splitter.split_documents(documents)
```
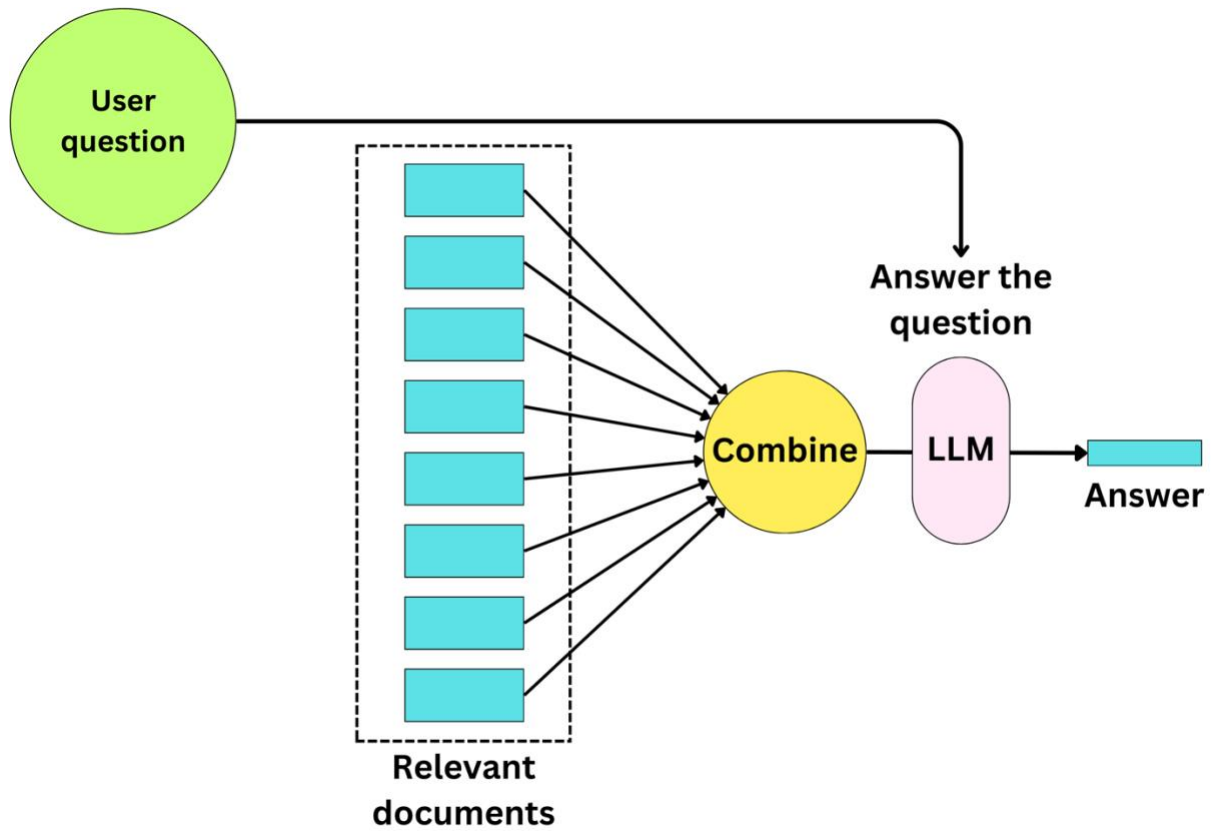
Let's index the data and create a chain

```
1.  index = FAISS.from_documents(documents, embeddings)
2.  retriever = index.as_retriever()
3.
4.  qa = RetrievalQA.from_chain_type(
5.      llm=llm,
6.      retriever=retriever,
7.  )
8.
9.  query = "What is a stuff chain?"
10.
11. qa.run(query, callbacks=[handler])
```
*'A stuff chain is a sequence of operations performed on a language model (LLM) to generate or process text. It typically consists of a language model chain (LLMChain) and a document chain (StuffDocumentsChain). The LLMChain is responsible for generating text based on a prompt, while the StuffDocumentsChain is used to process and manipulate documents or summaries. The specific details and functionality of a stuff chain can vary depending on the context and configuration.'*
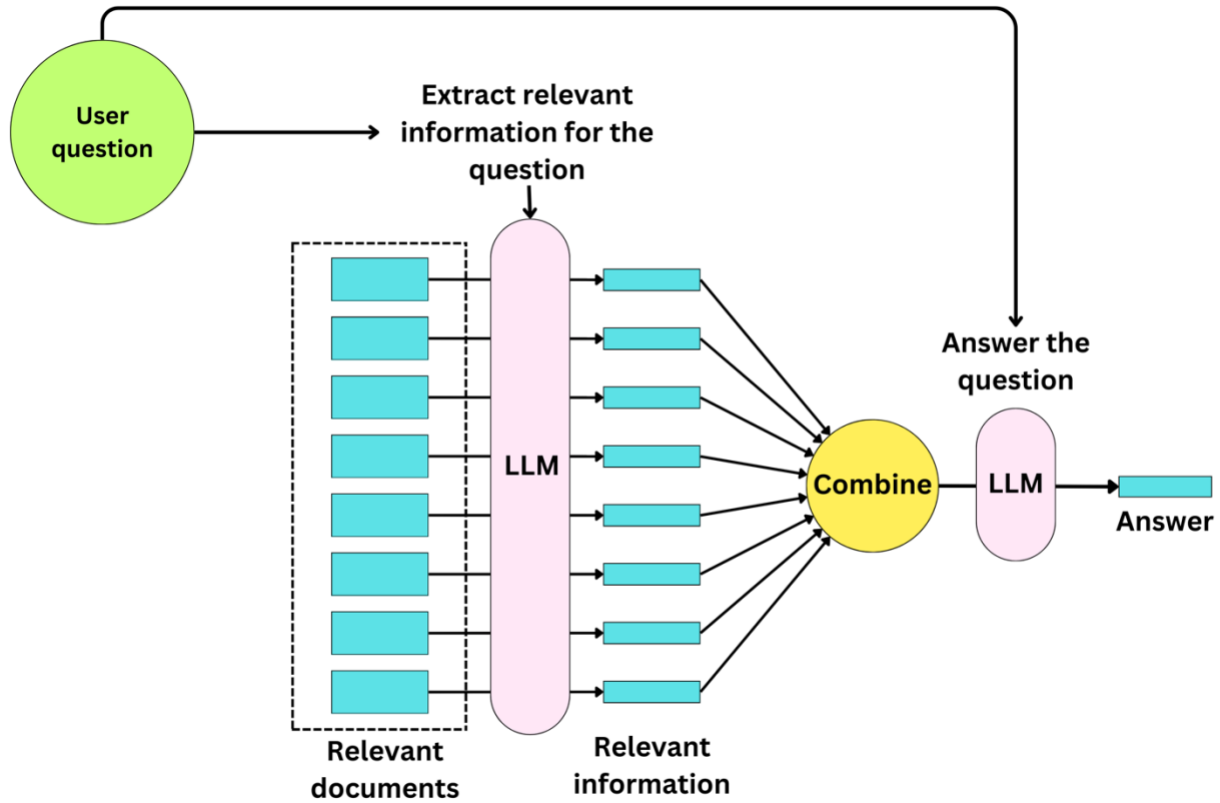
---------------------------------------------------------------------------------------------------------------------------------

The Stuff strategy

# Retrieval Augmented Generation



The Map-reduce strategy

# Retrieval Augmented Generation

User question

Extract relevant information for the question

LLM

Relevant documents

Relevant information

Combine

Answer the question

LLM

Answer

The Refine strategy

User question

Relevant documents

The Map-rerank strategy

# Retrieval Augmented Generation