# Retreival Augmented Generation from web data

Md Khalid Siddiqui

## Introduction

## Md Khalid Siddiqui

- Masters Student
- Technology and Innovation Management
- Hochschule Harz
- Faculty of Automation and Computer Science

#### Interest domain:

Large Language Model Applications, Prompt Engineering, Deep Learning particularly NLP

#### **Completed Projects:**

- RAG using web data (Current Presentation)
- <u>Deep Learning Model Deployment on Hugging face using Gradio(model generation and deployment)</u>
- Prompt Engineering chat bot project on Huggingface



This notebook makes a question answering chain with a specified website as a context data.

## Setting up

Install dependencies

```
!pip install langchain
!pip install pinecone-client
!pip install openai
!pip install tiktoken
!pip install nest_asyncio
```

```
Set up OpenAl API key
    import os
    os.environ["OPENAI_API_KEY"] = "sk-0eXpSaxE
                                                           hidden
Set up Pinecone API keys
    import pinecone
    # initialize pinecone
    pinecone.init(
                                  hidden
        api_key="78f1
                                                        find at app.pinecone.io
        environment="asia-southeast1-gcp-free" # next to api key in console
```

## Index

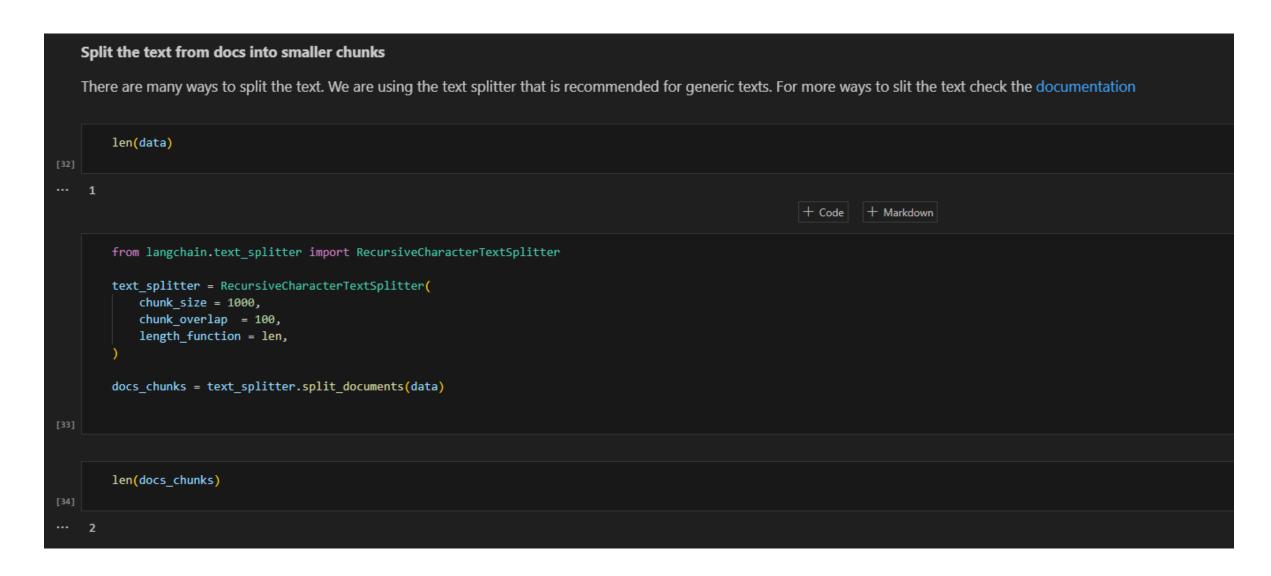
#### Load data from Web

Extends from the WebBaseLoader, this will load a sitemap from a given URL, and then scrape and load all the pages in the sitemap, returning each page as a document.

The scraping is done concurrently, using WebBaseLoader. There are reasonable limits to concurrent requests, defaulting to 2 per second.

Link to the documentation

```
!pip install unstructured
from langchain.document_loaders import UnstructuredURLLoader
urls = [
    "https://www.naeco.blue/"
loader = UnstructuredURLLoader(urls=urls)
data = loader.load()
```



## Create embeddings

```
from langchain.embeddings.openai import OpenAIEmbeddings

embeddings = OpenAIEmbeddings()
[36]
```

## Creating a vectorstore

A vectorstore stores Documents and associated embeddings, and provides fast ways to look up relevant Documents by embeddings.

There are many ways to create a vectorstore. We are going to use Pinecone. For other types of vectorstores visit the documentation

First you need to go to Pinecone and create an index there. Then type the index name in "index\_name"

```
from langchain.vectorstores import Pinecone
index_name = "hsharzirina"

# #create a new index
docsearch = Pinecone.from_documents(docs_chunks, embeddings, index_name=index_name)

# if you already have an index, you can load it like this
#docsearch = Pinecone.from_existing_index(index_name, embeddings)
[37]
```

Vectorstore is ready. Let's try to query our docsearch with similarity search

```
query = " Who is Felix Mertens?"
  data = docsearch.similarity_search(query)
  print(data)

[Document(page_content='Business Case\n\nProblem\n\nVolatilität der erneuerbaren Ene
```

## Making a question answering chain

The question answering chain will enable us to generate the answer based on the relevant context chunks. See the documentation for more explanation.

Additionally, we can return the source documents used to answer the question by specifying an optional parameter when constructing the chain. For more information visit the documentation

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

llm=OpenAI()

qa_with_sources = RetrievalQA.from_chain_type(llm=llm, chain_type="stuff", retriever-docsearch.as_retriever(), return_source_documents=True)

query = "Who is Felix Mertens?"
result = qa_with_sources({"query": query})
result["result"]

' Felix Mertens is the founder and CEO of NAECO Blue Management.'
```

 Output source documents that were found for the query

```
result["source_documents"]
```

[41]

[Document(page\_content='Business Case\n\nProblem\n\nVolatilität der erner Document(page\_content='English\n\nDeutsch\n\nNEWSNEWSCareer\n\nLegalLegation |

```
query = "What does Naeco Blue do?"
    result = qa_with_sources({"query": query})
    result["result"]

[42]
... ' Naeco Blue develops digital solutions that aim to minimiz
```

# Thankyou

Md Khalid Siddiqui

Questions: academia.mdkhalid@gmail.com

Contact: +491635346918