

Project - RAG LLM

Aim

To develop a structured pipeline for retrieving historical information about India from Wikipedia, processing it into meaningful embeddings, storing it in a vector database, and building a Retrieval-Augmented Generation (RAG) chatbot that answers user queries accurately using FAISS and the Gemini API.

Acceptance Criteria

1. Data Scraping

- Successfully extract and store historical Wikipedia pages related to Indian history in JSON format.
- Ensure that only relevant Wikipedia links are collected while avoiding irrelevant pages like "Category" and "File" links.
- Implement rate-limiting mechanisms to prevent getting blocked by Wikipedia.

2. Chunking, Embedding, and Storing in FAISS

- Successfully load the **BERT model** and generate embeddings for historical text.
- Implement **semantic chunking** to divide long texts into manageable segments.
- Store all embeddings and metadata efficiently in the FAISS vector database.
- Ensure the FAISS index is properly saved and can be reloaded for future queries.

3. Query Processing and Answer Generation

- Convert user queries into embeddings and retrieve the most relevant historical documents from FAISS.
- Generate responses using the **Gemini API**, incorporating retrieved historical content.
- Ensure that responses are **factually accurate** and **contextually relevant** to user queries.
- Provide a user-friendly **Streamlit interface** for seamless interaction.

1 - Scrapping the Data - History of India

Source - Wikipedia

Purpose

- Automate the extraction of historical Wikipedia pages starting from the W "History of India" page.
- Collect **titles and main content** from relevant Wikipedia pages.
- Store the extracted information in a **structured JSON format** for future analysis.

Technologies Used

- **Python**: Main programming language.
- **Requests**: To send HTTP requests and fetch webpage content.
- **BeautifulSoup**: For parsing and extracting data from Wikipedia pages.
- **tqdm**: To display progress while scraping multiple pages.
- **JSON**: To store structured data in a readable format.

Code Workflow

The script follows a **systematic process** to extract and store data:

Step 1: Fetch Wikipedia Links

- Start from the W "History of India" Wikipedia page.
- Identify and collect **valid Wikipedia links** related to history.
- **Filter out** irrelevant links such as files, categories, and special pages.

Step 2: Scrape Wikipedia Content

- For each valid link:
 - Extract the **page title**.
 - Extract the **main content** (excluding navigation and citations).
 - Store the data in a structured format.

Step 3: Store Data

- Convert the collected data into a **structured JSON file** (wikipedia_history.json).
- Ensure that Unicode characters are stored correctly.

Step 4: Prevent Rate Limiting

- Introduce **random delays (1-3 seconds)** between requests to avoid getting blocked by Wikipedia.

Code Explanation

```
7
8 WIKI_BASE_URL = "https://en.wikipedia.org"
9
```

- Defines the base Wikipedia URL.

Extract Valid Wikipedia Links

```
10 # Function to extract valid Wikipedia links
11 def get_wiki_links(url, visited):
12     """Extracts valid Wikipedia history-related links from a given page."""
13     response = requests.get(url)
14     soup = BeautifulSoup(response.text, "html.parser")
15     links = set()
16
17     for link in soup.find_all("a", href=True):
18         href = link["href"]
19         if href.startswith("/wiki/") and ":" not in href and href not in visited:
20             full_url = WIKI_BASE_URL + href
21             links.add(full_url)
22
23     return links
24
```

- Fetches all links from the given Wikipedia page.
- Ensures only **valid** Wikipedia article links are considered.
- Filters out **special pages** like Category:, File:, Help: .

Scrape Wikipedia Content

```
25 # Function to scrape content from a Wikipedia page
26 def scrape_wikipedia_page(url):
27     """Scrapes the title and main content of a Wikipedia page."""
28     response = requests.get(url)
29     soup = BeautifulSoup(response.text, "html.parser")
30
31     # Extract the title
32     title = soup.find("h1").text.strip()
33
34     # Extract the main content
35     content = []
36     for p in soup.find_all("p"):
37         text = p.text.strip()
38         if text:
39             content.append(text)
40
41     return {"title": title, "content": " ".join(content)}
42
```

- Extracts the **title** and **main textual content** of a Wikipedia page.
- Ignores **irrelevant sections** like tables, sidebars, and references.

Main Execution Flow

```
42
43 # Starting from the Wikipedia page - history of India
44 start_url = "https://en.wikipedia.org/wiki/History_of_India"
45 visited_pages = set()
46 scraped_data = []
47
48 # Get initial links from the main History of India page
49 wiki_links = get_wiki_links(start_url, visited_pages)
50
51 print(f"✅ Found {len(wiki_links)} potential pages to scrape.")
52
```

- Starts from the **History of India** page.
- Collects valid Wikipedia links.

Scrape Multiple Pages

- Scrapes **up to 1000** Wikipedia pages.
- **Handles exceptions** gracefully (e.g., network failures).
- **Random delays** prevent getting blocked by Wikipedia.

Save Data in JSON

- Stores the scraped data in a **structured JSON format**.
- Ensures **Unicode compatibility**.

```

1  {
2  {
3      "title": "Arthashastra",
4      "content": "Divisions Sama vedic Yajur vedic Atharva vedic Vaishnava puranas Shaiva puranas Shakta puranas Kautilya's Arthashastra (Sanskrit: अर्थशास्त्र, IAST: Kautīliyam Arthasāstram; <See R
5  },
6  {
7      "title": "Soanian",
8      "content": "Fertile Crescent: Europe: Africa: Siberia: The Soanian culture is a prehistoric technological culture from the Siwalik Hills, Pakistan.[1][4] It is named after the Soan Valley
9  },
10 {
11     "title": "Maharana Pratap",
12     "content": "Pratap Singh I (9 May 1540 – 19 January 1597), popularly known as Maharana Pratap (IPA: [maːɦaːˈɽaːqəː pɾaːt̪aːp] ⓘ), was king of the Kingdom of Mewar, in north-western India in
13 },
14 {
15     "title": "Rajput Regiment",
16     "content": "The Rajput Regiment is one of the oldest infantry regiments of the Indian Army. The regiment traces its history back to 1778, when the 24th Regiment of Bengal Native Infantry v
17 },
18 {
19     "title": "Ashrafi",
20     "content": "Ashrafi (Arabic: أَشْرَافِي) is a gold coin which originated in the Muslim World, and which was later widely adopted as currency in regions under Muslim rule in the Middle East, H
21 },
22 {
23     "title": "Golconda",
24     "content": "Golconda is a fortified citadel and ruined city located on the western outskirts of Hyderabad, Telangana, India.[1][2] The fort was originally built by Kakatiya ruler Pratāpa
25 },

```

BERT (Bidirectional Encoder Representations from Transformers) is a transformer-based model developed by Google that has significantly advanced natural language understanding tasks. By capturing context from both directions, BERT produces embeddings that encapsulate the semantic meaning of text, making it ideal for tasks like semantic search.

b. Loading the BERT Model and Tokenizer

The BERT model and its corresponding tokenizer are loaded using Hugging Face's `transformers` library. The specific model used is `sentence-transformers/all-MiniLM-L6-v2`, a variant optimized for producing sentence embeddings efficiently.

```

chunking_embedding_vectorize.py > process_json
1  '''Hugging Face BERT Embedding'''
2  import json
3  import faiss
4  import numpy as np
5  import os
6  from abc import ABC, abstractmethod
7  from langchain.text_splitter import RecursiveCharacterTextSplitter
8  from langchain.schema import Document
9  from dotenv import load_dotenv
10 from transformers import AutoModel, AutoTokenizer
11 import torch
12
13 # Load Hugging Face Model
14 model_name = "sentence-transformers/all-MiniLM-L6-v2"
15 tokenizer = AutoTokenizer.from_pretrained(model_name)
16 model = AutoModel.from_pretrained(model_name)
17

```

Generating Embeddings

A function `get_embedding` is defined to convert input text into embeddings. The text is tokenized and passed through the BERT model, and mean pooling is applied to obtain a fixed-size vector representation.

```

def get_embedding(text):
    tokens = tokenizer(text, return_tensors="pt", padding=True, truncation=True, max_length=512)
    with torch.no_grad():
        embeddings = model(**tokens).last_hidden_state.mean(dim=1) # Mean pooling
    return embeddings.squeeze().numpy().tolist()

```

Text Chunking

To manage lengthy documents, a `SemanticChunker` class is implemented using LangChain's `RecursiveCharacterTextSplitter`. This class divides text into manageable chunks, ensuring that each segment is within the model's processing limits.

```

# Semantic Chunker (LangChain)
class SemanticChunker(BaseChunker):
    def __init__(self, chunk_size=512, chunk_overlap=50):
        self.text_splitter = RecursiveCharacterTextSplitter(
            chunk_size=chunk_size, chunk_overlap=chunk_overlap, separators=["\n\n", "\n", " ", ""]
        )

    def chunk(self, text: str):
        return self.text_splitter.split_text(text)

```

FAISS Vector Database Management

The `FAISSVectorDB` class manages the indexing and retrieval of document embeddings. It handles adding new documents, saving and loading the index, and storing metadata.

```

1 # FAISS Vector Database Class
2 class FAISSVectorDB:
3     def __init__(self, index_file="faiss_index.bin", metadata_file="documents_metadata.json"):
4         self.index_file = index_file
5         self.metadata_file = metadata_file
6         self.index = None
7         self.docs = []
8
9         self.load_index()
10        self.load_metadata()
11
12    def add_documents(self, docs):
13        embeddings = np.array([get_embedding(doc.page_content) for doc in docs], dtype=np.float32)
14
15        if self.index is None:
16            self.index = faiss.IndexFlatL2(embeddings.shape[1]) # Initialize FAISS
17
18        faiss.normalize_L2(embeddings)
19        self.index.add(embeddings)
20        self.docs.extend(docs)
21
22        self.save_index()
23        self.save_metadata()
24
25    def save_index(self):
26        if self.index is not None:
27            faiss.write_index(self.index, self.index_file)
28
29    def load_index(self):
30        if os.path.exists(self.index_file):
31            self.index = faiss.read_index(self.index_file)
32            print("✅ FAISS index loaded.")
33        else:
34            print("⚠️ No FAISS index found. A new one will be created.")
35
36    def save_metadata(self):
37        print(f"💡 Saving {len(self.docs)} metadata entries to {self.metadata_file}...") # Debugging print
38        with open(self.metadata_file, "w", encoding="utf-8") as f:
39            json.dump([{"title": doc.metadata["title"], "content": doc.page_content} for doc in self.docs], f, indent=4)
40        print("✅ Metadata file saved successfully!")
41
42    def load_metadata(self):
43        if os.path.exists(self.metadata_file):
44            with open(self.metadata_file, "r", encoding="utf-8") as f:
45                self.docs = [Document(page_content=doc["content"], metadata={"title": doc["title"]}) for doc in json.load(f)]
46            print("✅ Document metadata loaded.")
47        else:
48            print("⚠️ No metadata found. It will be created.")

```

f. Processing JSON Data

The `process_json` function reads a JSON file containing documents, chunks the content, and adds the resulting segments to the FAISS index.

```

1 # Function to Process JSON Data
2 def process_json(json_file, chunker, vector_db):
3     with open(json_file, "r", encoding="utf-8") as file:
4         data = json.load(file)
5
6     docs = []
7     for entry in data:
8         title = entry["title"]
9         content = entry["content"]
10        chunks = chunker.chunk(content)
11
12        for chunk in chunks:
13            doc = Document(page_content=chunk, metadata={"title": title})
14            docs.append(doc)
15
16    print(f"✅ Adding {len(docs)} documents to FAISS")
17    vector_db.add_documents(docs)
18
19    print("💡 Manually saving metadata...")
20    vector_db.save_metadata()

```

Execution Flow

The main execution block of the script performs the following steps:

- **Initialization:** It initializes the `SemanticChunker` with a chunk size of 512 and an overlap of 50 characters. Simultaneously, it sets up the `FAISSVectorDB` with specified file paths for the FAISS index and metadata storage.
- **Model Verification:** A test embedding is generated using the `get_embedding` function with the input "Hello, test query" to ensure that the Hugging Face BERT model is functioning correctly.
- **Index Verification:** The script checks if the FAISS index is already built by verifying the existence of the index file and the number of vectors (`ntotal`) in the index.
- **Index Construction:** If the FAISS index is absent or empty, the script processes the `wikipedia_history.json` file. It uses the `SemanticChunker` to divide the content into manageable chunks and adds these chunks to the FAISS index.
- **Metadata Saving:** After processing, the script manually saves the metadata to ensure synchronization between the index and its associated metadata.
- **Completion Message:** Finally, the script outputs the total number of vectors present in the FAISS index, indicating the completion of the

indexing process.

This structured approach ensures that the text data is efficiently processed, embedded, and indexed, facilitating rapid similarity searches using the FAISS library.

3- Answering the Query:

Retrieval-Augmented Generation (RAG) Chatbot Overview

Retrieval-Augmented Generation (RAG) is a technique that enhances Large Language Models (LLMs) by integrating an information retrieval system. Instead of relying solely on pre-trained knowledge, the chatbot retrieves relevant documents from a database and uses them to generate accurate, fact-based responses. This approach reduces misinformation and improves contextual accuracy.

System Components

1. Hugging Face BERT Embeddings

The chatbot uses the **sentence-transformers/all-MiniLM-L6-v2** model from Hugging Face to generate numerical representations (embeddings) of text. These embeddings capture the semantic meaning of the text, enabling effective similarity searches.

Code Snippet: Initializing the BERT Model

```
Query.py > ...
1  import os
2  import json
3  import faiss
4  import numpy as np
5  import streamlit as st
6  from google import genai
7  from dotenv import load_dotenv
8  from transformers import AutoTokenizer, AutoModel
9  import torch
10
11
12
13 # Initializing Hugging Face BERT Model for embeddings
14 model_name = "sentence-transformers/all-MiniLM-L6-v2"
15 tokenizer = AutoTokenizer.from_pretrained(model_name)
16 model = AutoModel.from_pretrained(model_name)
17
```

Function to Generate Embeddings

```
# Function to generate embeddings using Hugging Face BERT
def get_embedding(text):
    tokens = tokenizer(text, return_tensors="pt", padding=True, truncation=True, max_length=512)
    with torch.no_grad():
        embeddings = model(**tokens).last_hidden_state.mean(dim=1)
    return embeddings.squeeze().numpy()
```

This function converts input text into a vector representation using BERT.

2. FAISS Vector Database

FAISS is a library for fast searching of high-dimensional embeddings. The chatbot stores document embeddings in a FAISS index and uses it to find the most relevant documents based on similarity.

Code Snippet: Loading the FAISS Index

```
# Loading FAISS index
index_file = "faiss_index.bin"
if not os.path.exists(index_file):
    raise FileNotFoundError(f"FAISS index file '{index_file}' not found.")
index = faiss.read_index(index_file) #index(1000,384)
```

3. Processing User Queries

When a user submits a query, it is converted into an embedding and compared with stored document embeddings to find the most relevant information.

Code Snippet: Query Processing

```

42 # Function to process user query
43 def process_query(query, top_k=5):
44     query_embedding = np.array([get_embedding(query)], dtype=np.float32)
45     faiss.normalize_L2(query_embedding)
46
47     # Search FAISS index
48     '''Distance will be a list of distance of vectors with the query vector
49     | indices are index of the those documents
50     |'''
51     distances, indices = index.search(query_embedding, top_k)
52
53     # Retrieving relevant documents
54     results = []
55     for idx in indices[0]:
56         if idx < len(documents_metadata): #If Index exists
57             doc = documents_metadata[idx]
58             results.append(doc)
59
60     return results
61

```

This function:

- Converts the query into an embedding
- Normalizes it for better comparison
- Searches for the top-k closest matches in the FAISS index
- Retrieves metadata of matching documents

4. Generating Responses with Gemini API

The chatbot sends the retrieved document content along with the user's query to **Google's Gemini API**, which generates a response.

Code Snippet: Generating AI Response

```

2 # Function to generate response using Gemini API
3 def generate_response(query):
4     relevant_docs = process_query(query, top_k=5)
5
6     if not relevant_docs:
7         return "No relevant information found in the knowledge base."
8
9     #Preparing context for LLM
10    context = "\n".join([doc["content"][:500] for doc in relevant_docs])
11
12    #Prompt to the Gemini
13    prompt = f"Based on the following information:\n{context}\n\nAnswer this question:\n{query}"
14
15    response = client.models.generate_content(
16        model="gemini-2.0-flash", contents=prompt
17    )
18
19    return response.text.strip()
20

```

This function:

- Retrieves the most relevant documents
- Prepares a **prompt** with retrieved context
- Sends the prompt to **Gemini API**
- Returns the generated response

5. Streamlit User Interface

The chatbot interface is built using **Streamlit**, providing real-time interaction.

Code Snippet: Setting Up the Streamlit Interface

```

81 # Streamlit UI
82 st.set_page_config(page_title="RAG Chatbot by Nileshe Mehra", page_icon=":speech_balloon:")
83
84 st.title("RAG Chatbot")
85 st.caption("Developed by Nileshe Mehra")
86
87
88 if "chat_history" not in st.session_state:
89     st.session_state.chat_history = []
90
91 # Display chat history
92 for chat in st.session_state.chat_history:
93     with st.chat_message(chat["role"]):
94         st.markdown(chat["content"])
95
96 # User input
97 if "end_chat" not in st.session_state or not st.session_state.end_chat:
98     user_query = st.chat_input("Enter your query:")
99     if user_query:
100         # Display user message
101         st.chat_message("user").markdown(user_query)
102         # Generate and display assistant response
103         response = generate_response(user_query)
104         st.chat_message("assistant").markdown(response)
105         # Update chat history
106         st.session_state.chat_history.append({"role": "user", "content": user_query})
107         st.session_state.chat_history.append({"role": "assistant", "content": response})
108
109 # End chat button
110 if st.button("End Chat"):
111     st.session_state.end_chat = True
112     st.write("Thank you for using the chatbot. This service is developed by Nileshe Mehra.")
113

```

This part of the code:

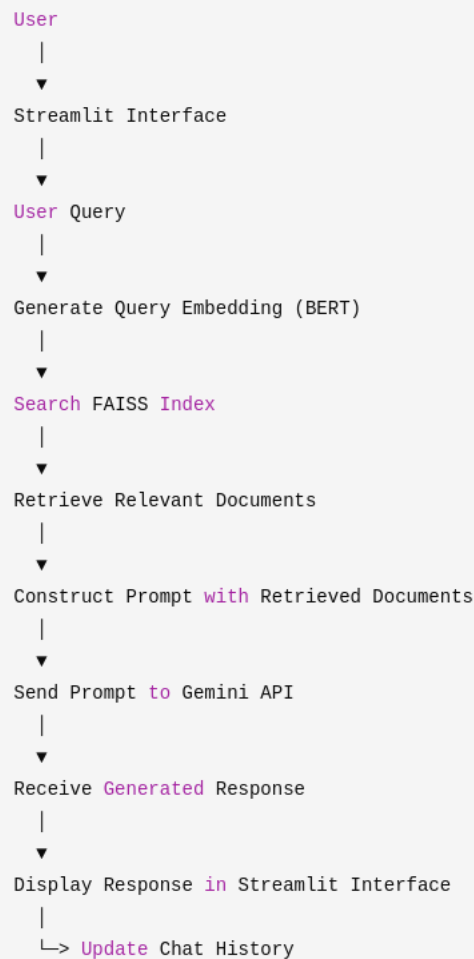
- Displays chat history
- Takes user input
- Calls the **generate_response()** function to fetch the response
- Updates the chat history dynamically

Execution Flow

1. The FAISS index and document metadata are loaded.
2. The chatbot interface is set up using Streamlit.
3. The user enters a query.
4. The query is converted into an embedding.
5. The FAISS index is searched for similar documents.
6. The most relevant documents are used as context for the **Gemini API**.
7. The chatbot generates a response and displays it to the user.
8. The chat history is updated.
9. The user can continue chatting or end the session.

Code Flow Diagram

Below is a flow diagram illustrating the interaction between the various components of the chatbot system:



Outcome

The project successfully automates Wikipedia data extraction, processes text using BERT embeddings, and stores it in a FAISS vector database for efficient retrieval. The RAG-based chatbot enhances responses with relevant historical context, leveraging the Gemini API. With a user-friendly Streamlit interface, users can quickly access accurate historical information in a scalable and efficient manner.