

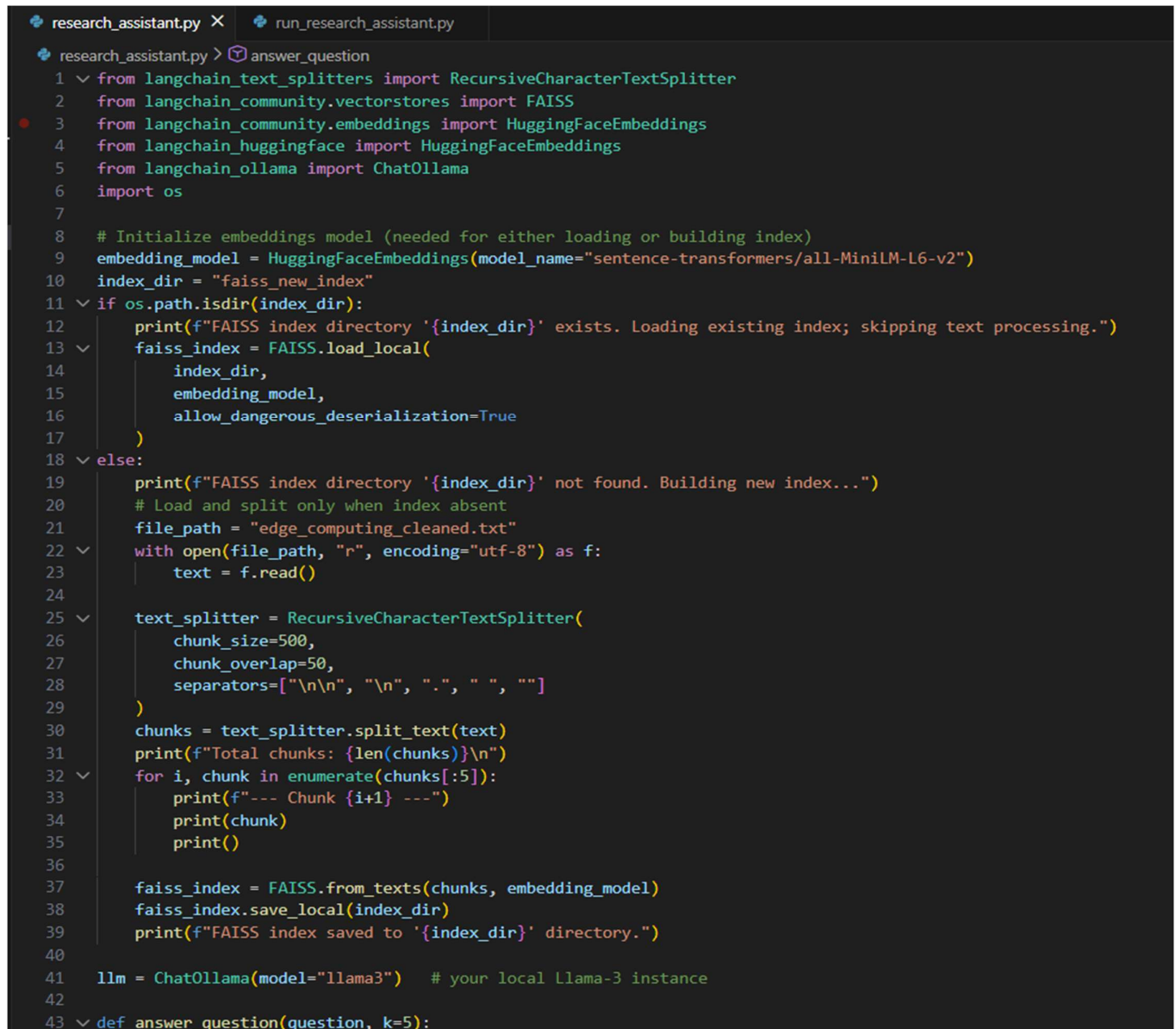
- **Introduction:**

This script implements a simple Retrieval-Augmented Generation pipeline. It either loads or builds a FAISS vector index using langchain from an edge computing text corpus **file name - edge_computing_cleaned.txt**, it retrieves the most relevant chunks (retrieving top 5 records) for a user question, and feeds those as context to a local large language model to produce a concise answer. To use the ollama, I need to install it on my local machine so that I can access it and can use ChatOllama to retrieve the data from file.

Fig1.1 and Fig1.2 Shows the main implementation of the RAG using langchain that shows how we upload the file, after that build a FAISS vector index to retrieve the relevant data and **Fig2.**, shows how we are calling the main file and what information we want.

To run this code, you must use - **python run_research_assistant.py** command (as shown in output) this will call **research_assistant.py** file that contains runnable implementations.

- **Code Screenshot:**



```
research_assistant.py X run_research_assistant.py
research_assistant.py > answer_question
1 from langchain_text_splitters import RecursiveCharacterTextSplitter
2 from langchain_community.vectorstores import FAISS
3 from langchain_community.embeddings import HuggingFaceEmbeddings
4 from langchain_huggingface import HuggingFaceEmbeddings
5 from langchain_ollama import ChatOllama
6 import os
7
8 # Initialize embeddings model (needed for either loading or building index)
9 embedding_model = HuggingFaceEmbeddings(model_name="sentence-transformers/all-MiniLM-L6-v2")
10 index_dir = "faiss_new_index"
11 if os.path.isdir(index_dir):
12     print(f"FAISS index directory '{index_dir}' exists. Loading existing index; skipping text processing.")
13     faiss_index = FAISS.load_local(
14         index_dir,
15         embedding_model,
16         allow_dangerous_deserialization=True
17     )
18 else:
19     print(f"FAISS index directory '{index_dir}' not found. Building new index...")
20     # Load and split only when index absent
21     file_path = "edge_computing_cleaned.txt"
22     with open(file_path, "r", encoding="utf-8") as f:
23         text = f.read()
24
25     text_splitter = RecursiveCharacterTextSplitter(
26         chunk_size=500,
27         chunk_overlap=50,
28         separators=["\n\n", "\n", ".", " ", ""],
29     )
30     chunks = text_splitter.split_text(text)
31     print(f"Total chunks: {len(chunks)}\n")
32     for i, chunk in enumerate(chunks[:5]):
33         print(f"--- Chunk {i+1} ---")
34         print(chunk)
35         print()
36
37     faiss_index = FAISS.from_texts(chunks, embedding_model)
38     faiss_index.save_local(index_dir)
39     print(f"FAISS index saved to '{index_dir}' directory.")
40
41     llm = ChatOllama(model="llama3") # your local Llama-3 instance
42
43 def answer_question(question, k=5):
```

Fig1.1: Implementation Code

```

43 def answer_question(question, k=5):
44     docs = faiss_index.similarity_search(question, k=k)
45     context = "\n\n".join([d.page_content for d in docs])
46     prompt = (
47         "Use the following context to answer concisely:\n"
48         f"{context}\n\n"
49         f"Question: {question}\n"
50         "Answer:"
51     )
52     response = llm.invoke(prompt)
53     return response.content

```

Fig1.2: Implementation Code

```

run_research_assistant.py > ...
1  import research_assistant as RA
2  if __name__ == "__main__":
3
4      query = "What is Edge Computing?"
5      print(f"\nGeneral Query: {query}\n")
6      answer = RA.answer_question(query)
7      print("Answer:\n")
8      print(answer)
9
10
11     query = "According to the file that you have just read, why it is important now a days?"
12     print(f"\nFollow-up Query: {query}\n")
13     answer = RA.answer_question(query)
14     print("Answer:\n")
15     print(answer)
16
17     query = "Give me a summary of the contents"
18     print(f"\nSummary Query: {query}\n")
19     answer = RA.answer_question(query)
20     print("Answer:\n")
21     print(answer)

```

Fig2: Demonstrating Usage

- **Output:**

```
PS C:\Users\achop\Documents\CIS583\project3> python .\run_research_assistant.py
2025-11-21 12:55:09.454366: I tensorflow/core/util/port.cc:153] oneDNN custom operations are on. You may see slightly different numerical results due to floating-point round-off errors from different c
computation orders. To turn them off, set the environment variable 'TF_ENABLE_ONEDNN_OPTS=0'.
WARNING:tensorflow:From C:\Users\achop\AppData\Local\Programs\Python\Python310\lib\site-packages\tf_keras\src\losses.py:2976: The name tf.losses.sparse_softmax_cross_entropy is deprecated. Please use t
f.compat.v1.losses.sparse_softmax_cross_entropy instead.

FAISS index directory 'faiss_new_index' exists. Loading existing index; skipping text processing.

General Query: What is Edge Computing?

Answer:

Edge computing is a way of processing data closer to where it's created, such as on phones or machines, instead of sending it all the way to a distant cloud or data center. This reduces delay, cuts down on internet usage, and enables faster real-time responses.

Follow-up Query: According to the file that you have just read, why it is important now a days?

Answer:

According to the text, edge computing is important nowadays because it enables devices to respond faster and process information in real-time, even in places with slow or unavailable cloud connections. This prevents network congestion, reduces latency, and increases data privacy, which is crucial for applications like self-driving cars and smart cities where rapid responses are critical.

Summary Query: Give me a summary of the contents

Answer:

The context discusses examples of AI agents in research papers, including:

* Sensors tracking soil moisture, temperature, and crop conditions, sending only critical updates to trigger irrigation systems.
* Trains equipped with sensors detecting mechanical faults or vibrations, analyzing data onboard and sending important alerts or maintenance requests.

In summary, the contents showcase practical applications of AI agents in monitoring and control systems.
PS C:\Users\achop\Documents\CIS583\project3> █
```

- **How does the attention + retrieval mechanism in Transformers compare to the sequence-based memory of RNNs/LSTMs?**

Sequential Based (RNN / LSTMs):

1. RNNs process inputs one step at a time in sequence.
2. They store memory in a hidden state passed from step to step.
3. LSTMs add gates to better control what to remember or forget.
4. Memory is compressed, so older information fades over time.
5. Training is sequential and slow, since each step depends on the previous one.

Transformers (Attention + retrieval):

1. Transformers process all tokens in parallel instead of step-by-step.
2. Attention compares every token with every other token to find relevance.
3. Memory is not a single vector the model can access the entire sequence directly.
4. Long-range dependencies are easy because any token can look anywhere.
5. Training is fast due to full parallelization and efficient gradient flow.
6. With retrieval augmentation (RAG), Transformers gain external memory far beyond the model size.

RNNs remember by carrying a compressed hidden state through time, while Transformers remember by directly looking back at any token using attention.