

Medical Chatbot RAG Implementation

LLM: OpenAI

Embeddings: OpenAI Embeddings

VectorDB: FAISS db

Frontend: streamlit

project utilizes Retrieval-Augmented Generation (RAG) to create a highly efficient medical chatbot. Here's a step-by-step breakdown of the process:

Upload PDF: The first step involves uploading a PDF document containing the necessary medical information. This document serves as the primary source of data for the chatbot.

PDF Chunking: Once the PDF is uploaded, the document undergoes a chunking process. This involves splitting the PDF into smaller, manageable chunks of text to facilitate easier processing and analysis.

Embedding Generation: After chunking, each text chunk is converted into embeddings using OpenAI Embeddings. These embeddings represent the semantic meaning of the text in a numerical format, which is essential for efficient retrieval and processing.

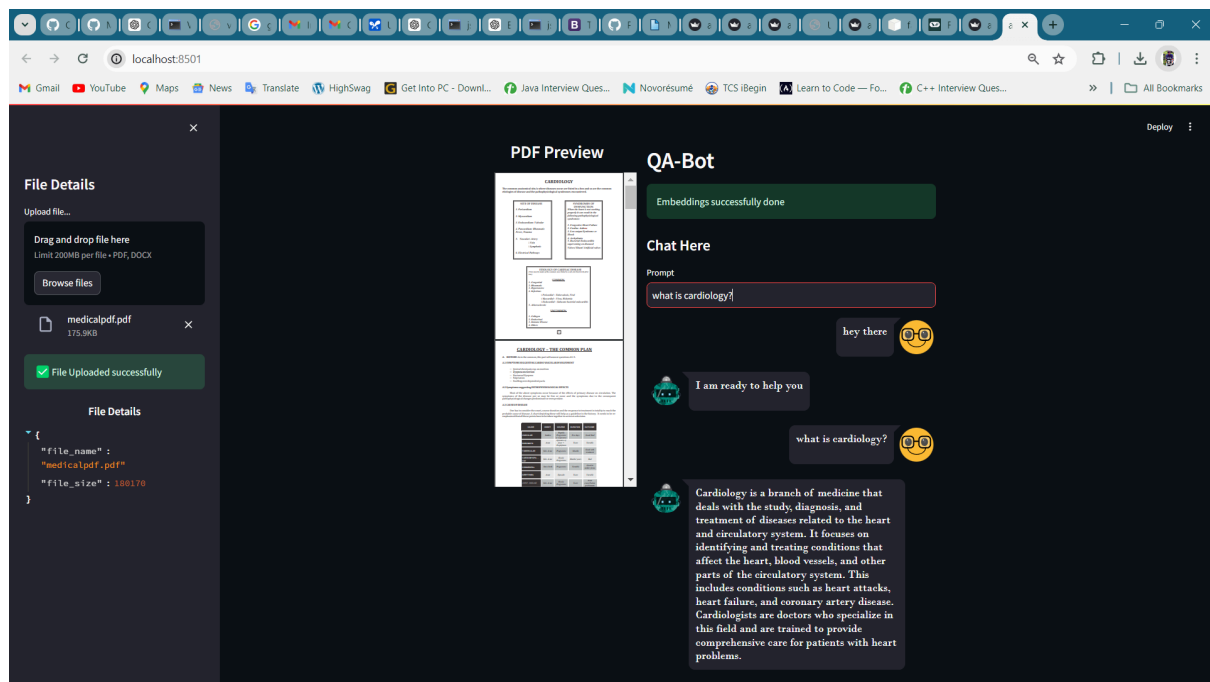
Storing Embeddings in VectorDB: The generated embeddings are stored in a FAISS (Facebook AI Similarity Search) database, a highly optimized vector database designed for similarity search and clustering of dense vectors.

RetrievalQA Chain for Query Processing: When a user asks a question, the query is first converted into a vector format. The RetrievalQA chain then retrieves the most relevant data from the FAISS vector database using similarity search.

Response Generation: The retrieved information is used to generate a precise and relevant response to the user's query, leveraging the capabilities of OpenAI's language model (LLM).

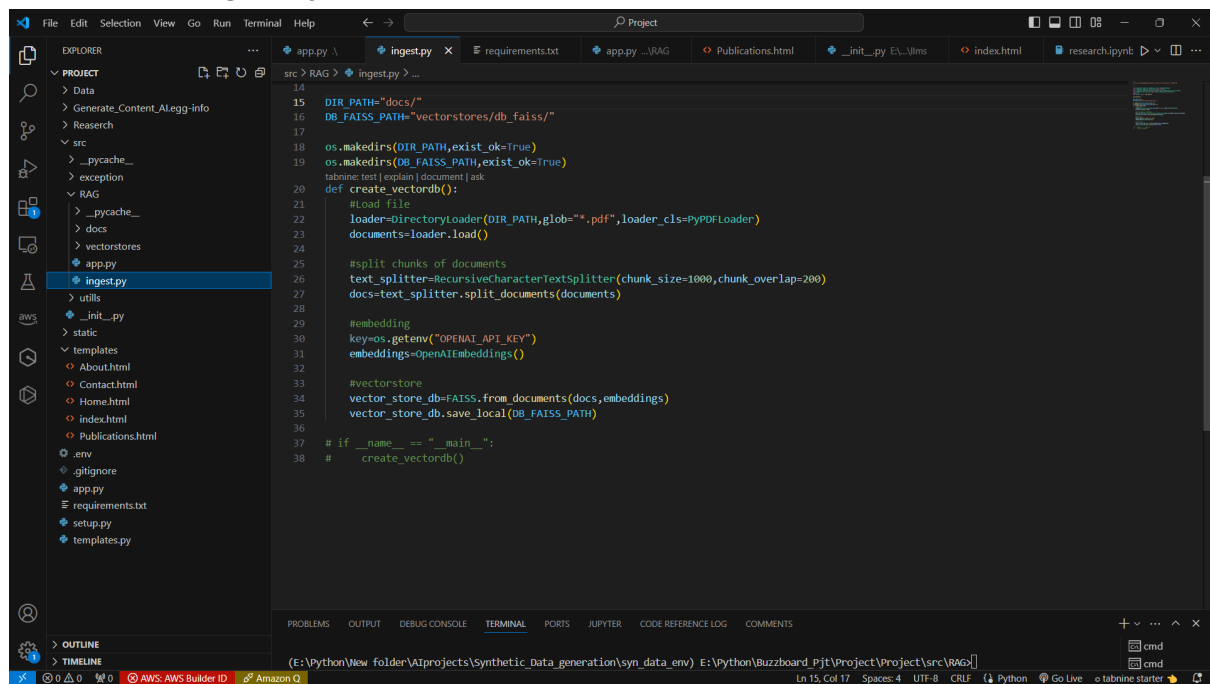
This RAG implementation ensures that our medical chatbot provides accurate, contextually relevant answers by effectively combining retrieval and generation techniques

Output:



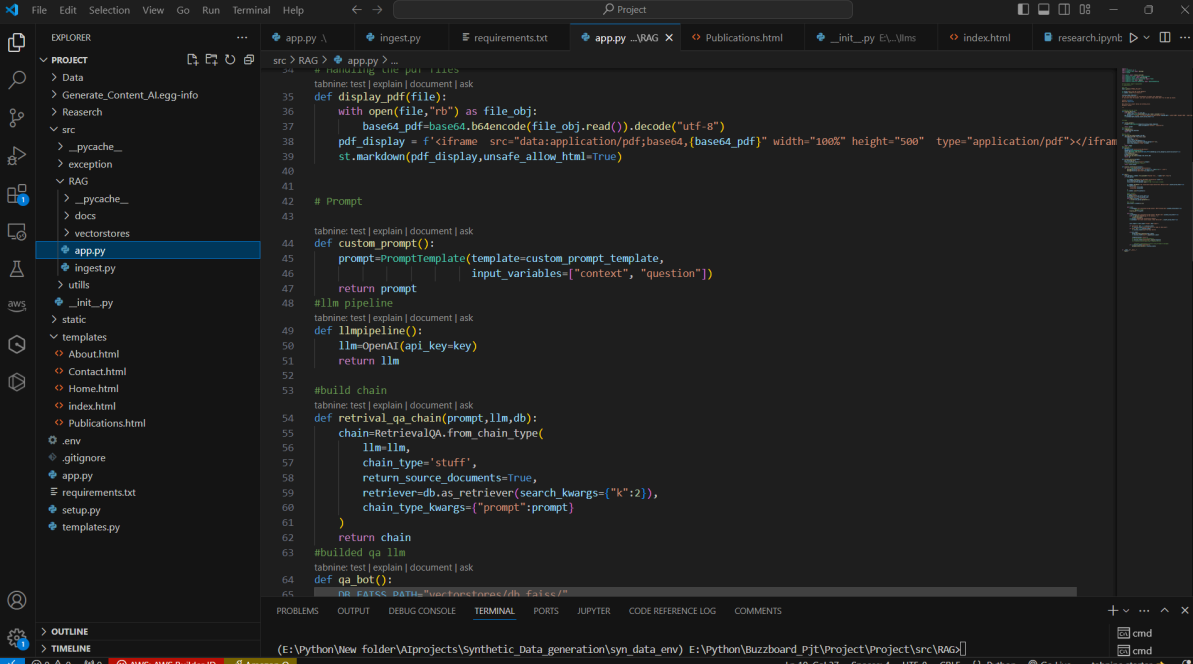
Below Code handling the DB creation :

path `src/RAG/ingest.py`



This code handling Retrieval data from vector db

src/RAG/app.py



```
src > RAG > app.py >
# Handling the RAG pipeline
tabnine: test | explain | document | ask
def display_pdf(file):
    with open(file, "rb") as file_obj:
        base64_pdf = base64.b64encode(file_obj.read()).decode("utf-8")
        pdf_display = f'<iframe src="data:application/pdf;base64,{base64_pdf}" width="100%" height="500" type="application/pdf"></iframe>'
        st.markdown(pdf_display, unsafe_allow_html=True)

# Prompt
tabnine: test | explain | document | ask
def custom_prompt():
    prompt = PromptTemplate(template=custom_prompt_template,
                             input_variables=["context", "question"])
    return prompt

# llm pipeline
tabnine: test | explain | document | ask
def llm_pipeline():
    llm = OpenAI(api_key=key)
    return llm

# build chain
tabnine: test | explain | document | ask
def retrieval_qa_chain(prompt, llm, db):
    chain = RetrievalQA.from_chain_type(
        llm=llm,
        chain_type='stuff',
        return_source_documents=True,
        retriever=db.as_retriever(search_kwargs={"k": 2}),
        chain_type_kwargs={"prompt": prompt}
    )
    return chain

# build qa llm
tabnine: test | explain | document | ask
def qa_bot():
    DR_FATSS_PATH = "usr/fatss/dh-fatss/"
```

(E:\Python\New folder\AIprojects\Synthetic_Data_generation\syn_data_env) E:\Python\Buzzboard_Pjt\Project\Project\src\RAG\