

AI Medical Chatbot using RAG

(MediRAG)



PROJECT LAYOUT

Phase 1 –Setup Memory for LLM (Vector Database)

- Load raw PDF(s)
- Create Chunks
- Create Vector Embeddings
- Store embeddings in ChromaDB

Phase 2–Connect Memory with LLM

- Setup LLM (Mistral with HuggingFace)
- Connect LLM with ChromaDB
- Create chain

Phase 3–Setup UI for the Chatbot

- Chatbot with Streamlit
- Load Vector store (ChromaDB) in cache
- Retrieval Augmented Generation–RAG

Tools and Technologies

- Langchain (AI Framework for LLM applications)
- HuggingFace (ML/AI Hub)
- Mistral (LLM Model)
- ChromaDB (Vector Database)
- Streamlit (For Chatbot UI)
- Python (Programming Language)
- VS Code (IDE)

Project layout



Medical Chatbot RAG Pipeline

Setup Memory for LLM (Vector Database)



Load raw PDF(s)

Import medical documents from the data/ directory.



Create Chunks

Split documents into manageable text chunks using LangChain's RecursiveCharacterTextSplitter.



Create Vector Embeddings

Generate embeddings with HuggingFace's sentence-transformers/ or openai for gpt-4o-mini.



Store Embeddings in ChromaDB

Persist embeddings in a ChromaDB vector store for efficient retrieval.

Connect Memory with LLM & UI



**(Mistral-7B-Instruct
Groq-v-0.3)**



Groq

Initialize a Mistral-7B or Groq-hosted model.



Create Chain with ChromaDB

Build a RetrievalQA chain.



Setup UI for the Chatbot

Develop a professional UI with Streamlit.



Streamlit Cache



Use Streamlit's caching

Retrieval Augmented Generation

Retrieval Augmented Generation (RAG)

Implement RAG for context-aware responses

Improvement Potential/Next Steps

- Add authentication in the UI
- Make use of self-upload document functionality
- Add multiple documents and embed them together
- Add Unit testing of RAG applications

Summary

- **Modern AI Chatbot for Medical Knowledge** MediBot leverages Retrieval-Augmented Generation (RAG) to deliver accurate, context-aware responses from medical PDFs.
- **Modular 3-Phased Development** Structured into memory setup (ChromaDB), LLM integration (LangChain & Groq), and a professional Streamlit UI.
- **Key Technologies Highlighted**
 1. Streamlit for interactive UI
 2. LangChain & HuggingFace for RAG pipeline
 3. ChromaDB for vector storage
 4. End-to-End RAG workflow