

DOCUMENT QA SYSTEM

Description:

This project implements a document-based Question-Answering (QA) system that allows users to upload documents (PDF, TXT, or Excel) and ask questions about the content. The system processes the document, extracts relevant text, and enables users to query the information using an AI model. It also stores past queries in an SQLite database for persistent history.

The backend uses LangChain for document processing, Hugging Face transformers for embeddings and text generation, FAISS for vector search, and Streamlit for the frontend UI.

Key Features:

Multi-Format Document Processing: Supports PDFs, TXT files, and Excel sheets, extracting textual content for analysis.

AI-Powered Q&A System: Utilizes FAISS vector search and Hugging Face transformers to efficiently retrieve and generate responses.

Persistent Query History: Stores user queries in an SQLite database, allowing users to view and revisit past questions.

Data Visualization for Excel Files: If structured data is uploaded, users can visualize it within the Streamlit interface.

User-Friendly Web Interface: A clean and intuitive UI built using Streamlit for easy document uploads and question interactions.

Scalability & Deployment: Designed for local execution and potential cloud deployment for broader accessibility.

Technical Stack:

Frontend: Streamlit

Backend Processing: LangChain for document loading and QA

Embedding & Retrieval: FAISS + Hugging Face's MiniLM embeddings

LLM Model: Transformer-based text generation (e.g., T5, OPT, or a locally hosted model)

Database: SQLite for storing user query history

Deliverables:

Web Application – A Streamlit-based UI to upload documents, process them, and allow interactive Q&A.

Document Processing Pipeline – Supports PDFs, text files, and Excel sheets, converting them into structured formats for retrieval.

QA System – Uses FAISS-based vector search with embeddings and a Hugging Face transformer model for text generation.

Persistent History Feature – Stores users' queries in an SQLite database for future reference.

Deployment Guide – Documentation for setting up and running the application locally or on a server.

Project Plan with Milestones:

Step 1 - Research Setup & Document Processing – Week 3(02/03 – 02/10)

Identify necessary libraries and set up a basic Streamlit app. Implement support for PDF, text, and Excel files.

Step 2 - QA System Integration – Week 4,5 (02/11 – 02/24)

Develop and integrate the QA system using FAISS and Hugging Face transformers.

Step 3 - Query Storage & UI Enhancement & Visualization – Week 6(02/25 – 03/10)

Implement SQLite database to store and retrieve previous user queries. Improve UI and add optional data visualization for Excel files.

Step 4 - Testing & Optimization – Week 7,8(03/11 – 03/24)

Test performance, optimize retrieval accuracy, and handle edge cases.

Step 5 – ML- OPS Deployment – Week 9 – 13(03/31 – 04/28)

Step 6 – Deployment & Documentation – Week 14(04/29 – 05/05)

Finalize deployment strategy and write user documentation.

Team Members:

Dinesh Kothandaraman

Sri Harshetha Amaravadi

Damodar Reddy Chirapureddy

Veda Samohitha Chaganti