# TASK-4: Architecture & Evaluation Report

## ➔ About the Project

This project implements an **Agentic Document Question Answering (QA) System** that enables users to ask natural language questions over uploaded documents and receive accurate, context-grounded responses. The system is designed using an agent-based Retrieval-Augmented Generation (RAG) architecture and emphasizes **local LLM inference**, **modular design**, and **containerized deployment**.

The project was developed as part of an **Internship Technical Task** and demonstrates practical implementation of modern AI workflows using LangGraph, vector databases, and FastAPI.
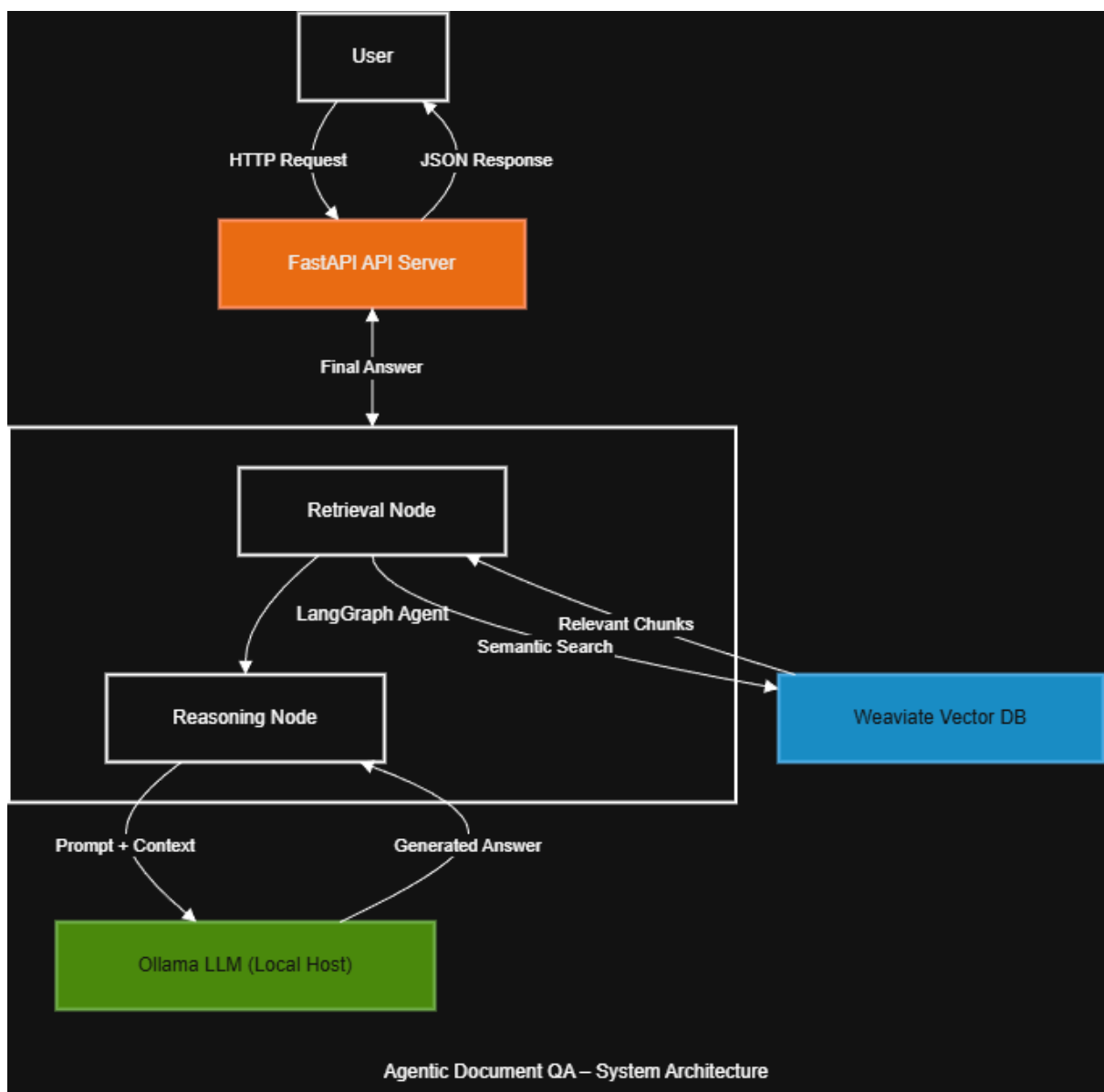
## ➔ Project Description

The system ingests PDF documents, splits them into semantically meaningful chunks, and stores vector embeddings in a Weaviate vector database. When a user submits a query, a LangGraph-based agent orchestrates the workflow by retrieving relevant document chunks and passing them to a local Large Language Model (LLM) running via Ollama. The model generates responses strictly grounded in the retrieved content, ensuring accuracy and minimizing hallucinations.

The application exposes its functionality through a RESTful API built with FastAPI and is fully containerized using Docker and Docker Compose for reproducible deployment.

## ➔ Technologies Used

| Component | Technology |
|---|---|
| Agent Framework | LangGraph |
| LLM Framework | LangChain |
| Vector Database | Weaviate |
| Embedding Model | Sentence-Transformers |
| LLM Inference | Ollama (Local) |
| API Framework | FastAPI |
| Containerization | Docker, Docker Compose |
| Programming Language | Python |

# ➔   <u>Architecture Diagram</u>



Agentic Document QA – System Architecture

# ➔     <u>Data Flow Summary</u>

The system receives user queries through a FastAPI REST endpoint. The query is forwarded to a LangGraph-based agent, where a retrieval node performs semantic similarity search on document embeddings stored in a Weaviate vector database. The retrieved document chunks are then passed to a reasoning node, which invokes a locally running Ollama large language model to generate a response strictly grounded in the retrieved context. The final answer is returned to the user in structured JSON format via the API.

# ➔     <u>Evaluation Observations</u>

The system successfully retrieves relevant document segments and produces accurate, context-aware responses for factual queries related to the ingested documents. The agentic workflow enforces strict grounding using retrieved context, effectively reducing hallucinations. Local inference using Ollama demonstrated low latency and stable performance for lightweight models. Semantic retrieval via Weaviate was efficient for document-level search. The Dockerized deployment ensured reproducibility and simplified system evaluation.

Screenshots

## ➔ <u>Strengths</u>

- Fully local inference (privacy-preserving)
- Clear separation of retrieval and reasoning
- Agentic design improves response grounding

- Production-ready API with FastAPI
- Scalable vector store architecture

# ➔ <u>Limitations & Future Improvements</u>

- Embedding generation depends on network availability during first load
- Evaluation performed on a single document corpus
- Can be extended with multi-document ingestion
- Can support streaming responses in future versions

# ➔ <u>Conclusion</u>

This project successfully demonstrates the design and implementation of an agentic document question answering system using modern AI frameworks. By combining semantic retrieval, agent-based reasoning, and local LLM inference, the system delivers accurate and privacy-preserving responses. The modular architecture and containerized deployment make the application extensible, reproducible, and suitable for real-world usage.