

Agentic RAG Chatbot: In-depth Documentation

This document provides a comprehensive deep dive into the architecture, design choices, and implementation details of the Agentic RAG Chatbot. It aims to clarify how the system works, even for those without a deep technical background.

Page 1: Introduction to RAG and Agentic AI

What is Retrieval-Augmented Generation (RAG)?

Imagine you have a brilliant but forgetful assistant. This assistant (a Large Language Model, or LLM) knows a lot about general topics because they've read countless books and articles. However, they might not know the very latest information, or specific details from your private documents. Sometimes, they might even confidently make up answers (this is called "hallucination").

RAG is a technique to solve this. It's like giving your assistant an "open book" exam. When you ask a question:

1. **Retrieval:** The assistant first quickly looks up relevant information from a specific set of documents (your "open book").
2. **Augmentation:** It then takes your question *and* the relevant information it found, and uses both to formulate an answer.
3. **Generation:** Finally, it generates a coherent response based on this augmented knowledge.

This ensures the answers are accurate, up-to-date, and grounded in your specific data, reducing hallucinations.

What is Agentic AI?

Think of a complex task, like building a house. You wouldn't have one person do everything (design, plumbing, electrical, carpentry). Instead, you'd have specialized professionals (agents) who each handle a specific part of the job.

Agentic AI applies this concept to software. Instead of one large, monolithic program trying to do everything, an agentic system breaks down a complex problem into smaller, manageable tasks, each handled by a dedicated "agent."

- Each agent has a clear role and responsibility.
- Agents communicate with each other to pass information and coordinate tasks.
- This makes the system more modular, easier to understand, maintain, debug, and scale.

Why Combine RAG with Agentic AI?

Combining RAG with an agentic architecture offers significant advantages:

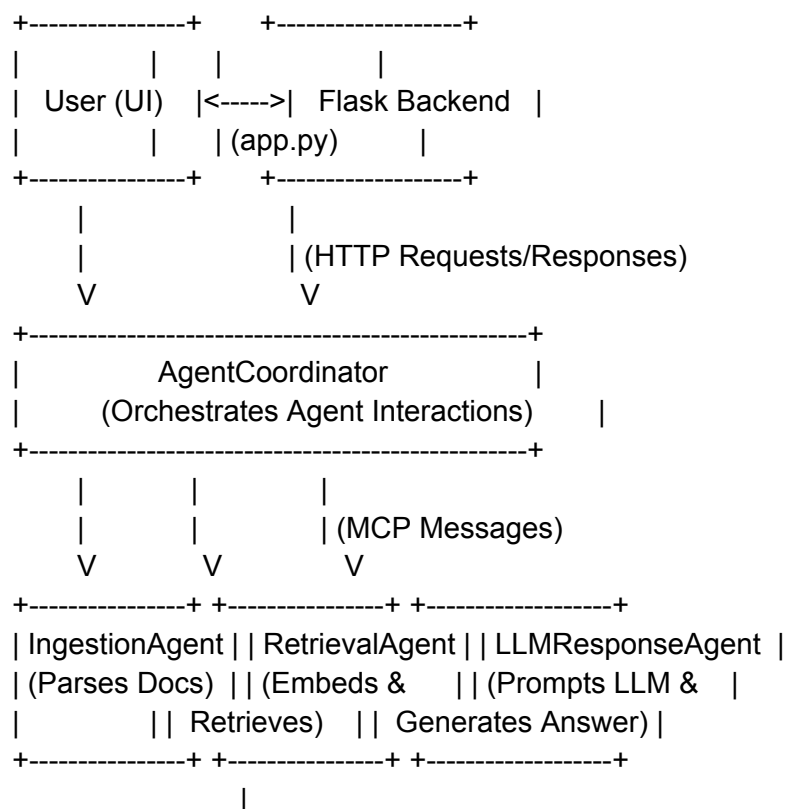
- **Clarity and Modularity:** Each step of the RAG process (document parsing, embedding, retrieval, LLM interaction) is handled by a distinct agent, making the system's logic very clear.
- **Maintainability:** If you want to change how documents are parsed (e.g., add support for a new format), you only need to modify the `IngestionAgent`, without touching other parts of the system.
- **Scalability:** Different agents could potentially run on different machines or services, allowing the system to handle more load by scaling individual components.
- **Robustness:** Issues in one agent are less likely to bring down the entire system.
- **Enhanced Control:** The agentic structure allows for fine-grained control over the flow of information and the decision-making at each stage.

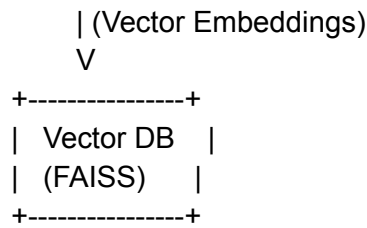
Project Goal

The goal of this project is to build a functional RAG chatbot that can answer questions based on user-uploaded documents, demonstrating a clear agent-based architecture and using a custom Model Context Protocol for inter-agent communication.

Page 2: System Architecture & Flow

High-Level System Diagram (Conceptual)





Agent Roles and Responsibilities

1. IngestionAgent:

- **Role:** The "Librarian" who organizes new books.
- **Responsibility:** Takes raw documents (PDF, DOCX, etc.), extracts all the text, and breaks that text down into smaller, manageable pieces called "chunks." These chunks are like individual paragraphs or sections, making them easier to search and process. It also adds metadata like the original source filename to each chunk.

2. RetrievalAgent:

- **Role:** The "Indexer" and "Search Engine" of the library.
- **Responsibility:**
 - **Indexing:** Takes the text chunks from the IngestionAgent and converts them into numerical representations called "embeddings." These embeddings are then stored in a special database called a "vector database" (using FAISS). Embeddings allow the computer to understand the meaning of text and find similar pieces of information very quickly.
 - **Retrieval:** When a user asks a question, it converts the question into an embedding and searches the vector database to find the chunks that are most semantically similar (i.e., most relevant in meaning) to the question.

3. LLMResponseAgent:

- **Role:** The "Smart Assistant" who answers questions.
- **Responsibility:** Takes the user's original question and the most relevant chunks retrieved by the RetrievalAgent. It then carefully crafts a "prompt" (a set of instructions and context) for a powerful Large Language Model (LLM). The LLM uses this prompt to generate a coherent, human-like answer that is directly based on the provided retrieved information. It also identifies the original source documents for transparency.

4. AgentCoordinator:

- **Role:** The "Project Manager" or "Orchestrator."
- **Responsibility:** This agent doesn't do the core work itself but manages the flow. It receives requests (like document uploads or chat queries) from the user interface and decides which other agents need to be involved, in what order. It sends messages to and receives responses from the other agents, ensuring the entire process runs smoothly.

System Flow Diagram (with Message Passing - MCP)

A. Document Upload Flow

1. **User (UI) - Flask Backend (app.py):** User uploads document.pdf.
2. **Flask Backend - AgentCoordinator:** app.py calls `coordinator.handle_document_upload(file_path)`.
3. **AgentCoordinator - IngestionAgent:**
 - **MCP Message (Type: INGEST_DOCUMENT_REQUEST)**
 - Payload: `{"file_path": "path/to/document.pdf"}`
 - IngestionAgent processes the document, extracts text, and creates chunks.
4. **IngestionAgent - AgentCoordinator:**
 - **MCP Message (Type: INGEST_DOCUMENT_RESPONSE)**
 - Payload: `{"chunks": [...], "file_path": "path/to/document.pdf"}`
5. **AgentCoordinator - RetrievalAgent:**
 - **MCP Message (Type: INDEX_CHUNKS_REQUEST)**
 - Payload: `{"chunks": [...]}`
 - RetrievalAgent generates embeddings and adds them to FAISS.
6. **RetrievalAgent - AgentCoordinator:**
 - **MCP Message (Type: INDEX_CHUNKS_RESPONSE)**
 - Payload: `{"status": "indexed", "num_chunks": N}`
7. **AgentCoordinator - Flask Backend:** Returns success/error status.
8. **Flask Backend - User (UI):** Displays upload status message.

B. Chat Query Flow

1. **User (UI) - Flask Backend (app.py):** User asks "What are the KPIs?".
2. **Flask Backend - AgentCoordinator:** app.py calls `coordinator.handle_chat_query(query)`.
3. **AgentCoordinator - RetrievalAgent:**
 - **MCP Message (Type: RETRIEVE_CONTEXT_REQUEST)**
 - Payload: `{"query": "What are the KPIs?"}`
 - RetrievalAgent searches FAISS for relevant chunks.
4. **RetrievalAgent - AgentCoordinator:**
 - **MCP Message (Type: RETRIEVE_CONTEXT_RESPONSE)**
 - Payload: `{"retrieved_context": [...], "query": "What are the KPIs?"}`
5. **AgentCoordinator - LLMResponseAgent:**
 - **MCP Message (Type: GENERATE_RESPONSE_REQUEST)**
 - Payload: `{"query": "What are the KPIs?", "retrieved_context": [...]}`
 - LLMResponseAgent formats the prompt and calls the LLM.
6. **LLMResponseAgent - AgentCoordinator:**

- **MCP Message (Type: GENERATE_RESPONSE_RESPONSE)**
 - Payload: {"answer": "...", "source_context": [...]}
7. **AgentCoordinator - Flask Backend:** Returns the LLM's answer and sources.
 8. **Flask Backend - User (UI):** Displays the chatbot's answer and sources.

Page 3: Technical Deep Dive

Model Context Protocol (MCP) - mcp/message_protocol.py

This file defines the MCPMessage class, which is a simple yet effective way to standardize communication between agents. Each message is a Python object (converted to a dictionary/JSON for conceptual passing) with the following structure:

```
{
  "sender": "AgentName",    // Name of the sending agent (e.g., "RetrievalAgent")
  "receiver": "AgentName",  // Name of the intended recipient (e.g., "LLMResponseAgent")
  "type": "MESSAGE_TYPE",   // Describes the purpose of the message (e.g.,
  "CONTEXT_RESPONSE")
  "trace_id": "unique-id-123", // A unique identifier to track a single request's journey
  "payload": {               // The actual data being transmitted
    "key1": "value1",
    "key2": ["item1", "item2"]
  }
}
```

This structured approach makes debugging easier and ensures all agents understand the data they are exchanging.

Ingestion Agent - agents/ingestion_agent.py

- **Purpose:** To turn unstructured documents into usable text chunks.
- **Key Libraries:**
 - PyPDF2: For extracting text from PDF files.
 - python-pptx: For reading text from PowerPoint presentations.
 - pandas: Used for reading and converting CSV data into a string format.
 - python-docx: For extracting text from Word documents.
 - markdown: To convert Markdown files into plain text.
- **Chunking Logic:** Documents are split into chunk_size (e.g., 500 words) pieces with chunk_overlap (e.g., 50 words). Overlapping chunks are crucial to ensure that important information isn't split across two chunks, losing context. Each chunk is stored with its content and the original source filename.

Retrieval Agent - agents/retrieval_agent.py

- **Purpose:** To create a searchable knowledge base and find relevant information.
- **Key Libraries:**
 1. **sentence-transformers:** This library provides pre-trained models (like all-MiniLM-L6-v2) that convert human-readable text into dense numerical vectors (embeddings). Texts with similar meanings will have embeddings that are numerically close to each other.
 2. **faiss:** (Facebook AI Similarity Search) A highly efficient library for performing similarity searches on large collections of vectors. It allows us to quickly find the "nearest neighbors" (most similar chunks) to a given query embedding.
- **Process:**
 1. **Embedding Generation:** Text chunks are passed through the SentenceTransformer model to get their embeddings.
 2. **FAISS Indexing:** These embeddings are added to a FAISS IndexFlatL2 index. This index is a data structure optimized for fast similarity lookups.
 3. **Retrieval:** When a query comes in, its embedding is generated, and FAISS is used to find the top_k (e.g., 3) most similar chunk embeddings. The original text content of these chunks, along with their source, is then returned.

LLM Response Agent - `agents/llm_response_agent.py`

- **Purpose:** To generate the final, intelligent answer using a Large Language Model.
- **Key Library:**
 - **transformers (Hugging Face):** This powerful library provides access to a vast collection of pre-trained LLMs.
- **LLM Integration (Highlight!):**
 - We specifically use the `pipeline("text2text-generation", model="google/flan-t5-small")`.
 - `google/flan-t5-small` is a smaller, instruction-tuned version of Google's T5 model. It's chosen for its balance of performance and relatively lower resource requirements compared to much larger models.
 - **Prompt Engineering:** The `_format_prompt` method is critical here. It constructs a carefully designed text input for the LLM. This prompt includes:
 - **System Instructions:** Directives to the LLM on how to behave (e.g., "You are a helpful AI assistant.").
 - **Groundedness Instruction:** Crucially, it tells the LLM to *only* use the provided "Retrieved Information" and to state if it cannot find the answer there. This is how we ensure the LLM stays "on topic" and doesn't hallucinate.
 - **Context:** The actual content of the relevant chunks retrieved by the `RetrievalAgent`.
 - **User Query:** The original question from the user.
 - The LLM then processes this entire prompt and generates a new, coherent answer.

Agent Coordinator - `agents/agent_coordinator.py`

- **Purpose:** The central nervous system of the application.
- **Functionality:** It instantiates the other three agents and acts as the intermediary. It receives high-level requests (like "upload this file" or "answer this question") and translates them into a sequence of calls to the appropriate agents, simulating the MCP message passing. It consolidates results from different agents and returns the final outcome to the Flask application. It also manages clearing all indexed data.

Flask Backend - `app.py`

- **Purpose:** Provides the web interface for the chatbot.
- **Key Features:**
 - Uses the Flask web framework to handle HTTP requests (GET for the main page, POST for uploads, chat, and clearing data).
 - Manages file uploads, saving them to the `documents/` directory.
 - Sets `MAX_CONTENT_LENGTH` to allow larger file uploads (increased to 200MB).
 - Acts as the bridge between the web UI and the `AgentCoordinator`, passing requests and responses.
 - Includes robust error handling and logging for server-side operations.

Frontend (HTML, CSS, JS) - `templates/index.html`, `static/style.css`, `static/script.js`

- **`index.html`:** The main web page structure. It includes a sidebar for controls and a main area for the chat conversation. It uses Tailwind CSS for rapid styling and a clean, responsive layout.
- **`style.css`:** Contains custom CSS rules to enhance the visual appeal, such as custom scrollbars, loader animations, and specific styling for chat messages and source contexts. It also includes media queries for basic mobile responsiveness.
- **`script.js`:** Handles all client-side interactivity:
 - Sending user messages and file uploads to the Flask backend using `fetch` API calls.
 - Dynamically updating the chat display with user and bot messages.
 - Displaying loading indicators and temporary notification messages.
 - **Crucially, it implements client-side request timeouts using `AbortController`** for both uploads and chat queries. This prevents the browser from hanging indefinitely if the server is slow or unresponsive, providing a better user experience by giving specific "timed out" feedback.

Page 4: Setup, Usage, Challenges, and Future Scope

Setup and Usage

For detailed instructions on setting up the project (cloning, virtual environment, installing dependencies) and how to use the chatbot interface (uploading documents, asking questions, clearing data), please refer to the README .md file in the project's root directory.

Challenges Faced

Developing this Agentic RAG Chatbot presented several interesting challenges:

1. **Document Parsing Robustness:** Handling diverse document formats (PDF, PPTX, DOCX, CSV, TXT/MD) reliably is complex. Each format has its quirks, and extracting clean, usable text can be tricky, especially with complex layouts or embedded objects. Initial attempts faced issues with certain file structures, leading to incomplete text extraction.
2. **Large File Processing & Timeouts:** When dealing with larger documents (e.g., multi-page PDFs), the process of reading, chunking, embedding, and indexing can be time-consuming. This led to client-side "Network Error" messages due to default browser fetch timeouts. We addressed this by:
 - Increasing the MAX_CONTENT_LENGTH in Flask to allow larger uploads.
 - Implementing explicit client-side fetch timeouts using AbortController in JavaScript, providing more informative "timed out" messages to the user instead of generic network errors.
3. **LLM Integration and Prompt Engineering:** Transitioning from a simulated LLM to a real one (google/flan-t5-small) required careful prompt engineering. Crafting instructions that effectively guide the LLM to use *only* the provided context and respond appropriately when information is missing is an iterative process. Balancing conciseness with clarity in the prompt is key to efficient and accurate responses.
4. **Dependency Management:** Managing Python dependencies, especially large ones like torch (required by transformers) and faiss-cpu, can sometimes lead to environment-specific installation issues. Ensuring a consistent and reproducible setup required clear requirements.txt and emphasis on virtual environments.

Future Scope & Improvements

The current implementation provides a strong foundation, but several areas can be enhanced for a more robust, scalable, and user-friendly production-ready system:

1. **Persistent Storage for Vector Store:**
 - **Current State:** The FAISS index and document metadata are stored in memory and are lost when the Flask server restarts.
 - **Improvement:** Implement a mechanism to save the FAISS index to disk (e.g., faiss.write_index) and load it on startup. For more advanced needs, integrate a dedicated persistent vector database (e.g., ChromaDB, Pinecone, Weaviate) that handles storage and scaling automatically.
2. **Enhanced Conversational Memory:**
 - **Current State:** The chatbot displays previous turns, but the LLM doesn't inherently use past conversation context to understand follow-up questions (e.g., "What about Q2?" after "What KPIs were tracked in Q1?").

- **Improvement:** Implement a "chat history" component that dynamically adds relevant previous turns to the LLM's prompt. This could also involve a "query rephrasing" sub-agent that uses an LLM to rephrase a short follow-up question into a self-contained query based on the conversation context.
- 3. **Advanced RAG Techniques:**
 - **Improvement:** Explore techniques beyond simple top-k retrieval:
 - **Re-ranking:** Use a separate, smaller model to re-rank the initial retrieved chunks to ensure the most relevant ones are prioritized.
 - **Hybrid Search:** Combine semantic search (FAISS) with keyword-based search for a more comprehensive retrieval.
 - **Multi-hop Reasoning:** For complex questions requiring information from multiple disparate chunks, implement a mechanism for the LLM to iteratively retrieve more information.
- 4. **Asynchronous Document Processing:**
 - **Current State:** Document processing (parsing, embedding, indexing) happens synchronously on the main Flask thread, which can block the UI for large files.
 - **Improvement:** Implement asynchronous task queues (e.g., using Celery with Redis or RabbitMQ) to offload document processing to background workers, allowing the UI to remain responsive.
- 5. **User Authentication and Multi-User Support:**
 - **Current State:** The application is single-user.
 - **Improvement:** Add user authentication (e.g., Flask-Login, OAuth) and partition the document storage and vector index per user, allowing multiple users to have their own private knowledge bases.
- 6. **Containerization:**
 - **Improvement:** Provide a Dockerfile and docker-compose.yml to encapsulate the application and its dependencies, making deployment and environment setup consistent across different machines.
- 7. **More Robust Error Handling & Logging:**
 - **Improvement:** Implement more granular error types, better user-facing error messages, and a centralized logging system for easier debugging in production.
- 8. **Frontend Enhancements:**
 - **Improvement:** Add a visual list of currently uploaded documents, progress bars during processing, and more interactive elements for a richer user experience.

These improvements would transform the current prototype into a more robust, scalable, and feature-rich RAG solution suitable for production environments.