# RAG QA Bot Model Architecture and Process Overview

**Backend:-**

**1. Model Architecture**

The RAG (Retrieval-Augmented Generation) QA bot is designed to answer questions or summarize documents by combining information retrieval and generative response mechanisms. Here's how the architecture works:

- **Embedding Model**: The system uses the Hugging Face model all-MiniLM-L6-v2 for embedding text. This model transforms chunks of the document into vector representations, allowing efficient similarity searches.

- **Chroma Vector Store**: The vector embeddings are stored in a Chroma vector database. This structure enables quick access to relevant parts of the document when a query is made.

- **Cohere API**: The bot uses Cohere's generative language model (command-r-plus-04-2024) to generate natural language answers or summaries based on the context retrieved from the document.

**2. Approach to Retrieval**

The RAG bot retrieves relevant content from the document through the following steps:

- **Document Splitting**: When a PDF is uploaded, it is first split into smaller chunks (around 500 characters with some overlap). This is done to optimize processing and retrieval, making it easier to match questions with specific sections of the document.

- **Embedding**: Each document chunk is converted into a vector (numerical representation) using the embedding model. These vectors capture the semantic meaning of each chunk, which is then stored in the vector database (Chroma).

- **Similarity Search**: When a question is asked, it is embedded in the same way. The vector representation of the question is compared to the document chunks in the vector database, and the most relevant chunks (top 5, in this case) are retrieved for context.

**3. Generative Response Creation**

Once the relevant document chunks are retrieved, the generative process begins:

- **Context Preparation**: The retrieved document chunks are assembled into a prompt using a predefined template. This prompt includes both the retrieved context and the user's question, ensuring that the generative model has all the necessary information to formulate an answer.

- **Answer Generation**: The prompt is then passed to the Cohere model, which generates a natural language response. The model uses the context to create an answer or summary that closely aligns with the information found in the document.

- **Post-processing**: The generated response is cleaned up (removing extra spaces, ensuring readability) and then returned to the user.

**Summary of Workflow**

1. **PDF Upload**: The user uploads a PDF document.

2. **Document Splitting**: The document is divided into chunks, and embeddings are generated for each chunk.

3. **Question Input**: The user can ask a question or request a summary.

4. **Retrieval**: Relevant chunks of the document are retrieved using vector similarity search.

5. **Response Generation**: Cohere's language model generates an answer or summary based on the retrieved content.

6. **Output**: The user receives the answer to their question or the document's summary.

## Frontend :-

**Overview**

The **PDF QA Assistant** allows users to upload PDF documents, generate summaries, and ask questions about the content of the document. The assistant will read the PDF and respond with relevant information based on the user's questions. It also provides the ability to export conversation history.

---

**How to Use**

**1. Uploading a PDF**

- On the **sidebar** of the app, you'll find an option titled **"📁 Document Upload"**.

- Click the **"Browse"** button and select a PDF file from your local device. Make sure the file is in **.pdf** format.

- Once the file is uploaded, the assistant will be ready to generate a summary and answer any questions you may have about the document.

**2. Generating a Summary**

- After uploading a PDF, click the **"📊 Generate Summary"** button in the sidebar.

- The assistant will analyze the document and generate a summary. This might take a few seconds depending on the length of the PDF.

- Once the summary is ready, it will appear in the **"Document Summary"** section in the main content area.

**3. Asking Questions**

- After uploading your PDF, you can ask questions about its content.

- Type your question into the chat input box at the bottom of the screen labeled **"Ask a question about your document"**.

- Press **Enter** or click the **send button** to submit your question.

- The assistant will process your query and respond in the chat with the answer.

### 4. Viewing Responses

- The conversation between you and the assistant will be displayed in a chat format in the main content area.

- You can scroll up to view the full conversation history.

- For each question you ask, the assistant will respond with relevant information from the PDF.

### 5. Exporting Conversation History

- You can export your conversation with the assistant to a PDF or text file.

- In the sidebar, click the **" 📤 Export Conversation"** button.

- If the conversation is available, you can download it as a **PDF** or **text** file.

### 6. Starting a New Session

- If you'd like to upload a new PDF and clear your conversation history, click the **" 🔄 New Session"** button in the sidebar.

- This will reset the app, and you'll be able to start fresh with a new PDF file.

---

**Features Summary**

- **Upload PDF**: Easily upload any PDF file for analysis.

- **Generate Summary**: Get a concise summary of the uploaded document.

- **Ask Questions**: Query the document content for quick answers.

- **Export Conversation**: Save the conversation with the assistant in a PDF or text format.

- **New Session**: Reset the app and start a new session with a different document.

---

**Notes**

- Ensure that you upload valid PDF files.

- The assistant may take a few seconds to generate responses depending on the complexity and size of the document.

- Always reset the session if you want to upload a new document or restart the analysis.

**Decisions, Challenges, and Solutions in the RAG QA Bot Development**

**Key Decisions**

1. **Model Selection**: We chose the Hugging Face model all-MiniLM-L6-v2 for embeddings due to its balance of speed and accuracy. The Cohere model (command-r-plus-04-2024) was used for generating responses because it offers coherent, fluent answers.

2. **Vector Store**: Chroma was selected for storing embeddings as it allows fast and efficient similarity searches.

3. **Chunking Strategy**: Documents were split into chunks of 500 characters with overlap to ensure context wasn't lost and to improve the accuracy of retrieval.

**Challenges Faced and Solutions**

1. **Session Refresh Issue**:

   o **Problem**: After uploading a new PDF, the system sometimes showed answers based on the previously processed PDF.

   o **Solution**: We ensured that all variables storing previous document data were cleared when a new PDF was uploaded to avoid stale data issues.

2. **Performance Delays**:

   o **Problem**: The initial processing, including embedding generation and similarity search, took time, especially for large PDFs.

   o **Solution**: We optimized by adjusting the chunk size and overlap. For future improvements, integrating a caching mechanism or batch processing could further speed up performance.

3. **Embedding Accuracy**:

   o **Problem**: Embeddings sometimes led to slightly irrelevant context retrieval, affecting the quality of responses.

   o **Solution**: We refined the similarity search to retrieve the top 5 most relevant chunks, balancing performance and precision.

4. **PDF Format Issues**:

   o **Problem**: Some PDFs contained embedded images or non-standard text encoding, causing errors during processing.

   o **Solution**: We added a check to validate that the uploaded file is a valid PDF and provided clear error messages for unsupported formats.

**Conclusion**

By addressing these challenges and making decisions aimed at balancing speed, accuracy, and ease of use, we developed a robust system that efficiently processes PDFs, retrieves relevant information, and generates coherent responses. Future optimizations can further enhance performance and usability.