# RETRIEVAL-AUGMENTED GENERATION (R.A.G.) WITH OLLAMA

## AN INTERACTIVE CHATBOT FOR DOCUMENT ANALYSIS

# Project Overview

## Goals

The primary goal of this project is to create a versatile, AI-driven chatbot that leverages Retrieval-Augmented Generation (RAG) techniques to analyze, summarize, and interact with uploaded PDF documents. This chatbot not only facilitates document-based query resolution but also serves as a powerful tool to test and optimize RAG model performance by adjusting a variety of parameters, including language model (LLM) selection, embedding models, chunk size, and LLM-specific hyperparameters.

Additionally, the tool can operate as a standalone chatbot, enabling general interactions without document uploads, and making it adaptable for broader use cases.

## Data Sources

The core data source for this chatbot comprises uploaded PDF documents. These documents are parsed, segmented into manageable chunks, and stored in a vector database. The vector database allows for efficient retrieval of content relevant to the user's query. Key aspects of the data flow include:

- **User-Uploaded PDF Files:** The app supports a wide range of documents, parsed into text chunks for processing.

- **External API (Ollama):** Ollama's LLM and embedding models provide contextual embeddings and generate responses based on retrieved document segments.

- **ChromaDB:** This vector database handles the embeddings for document chunks, supporting semantic search within document content.

## Project Purpose and Additional Use Cases

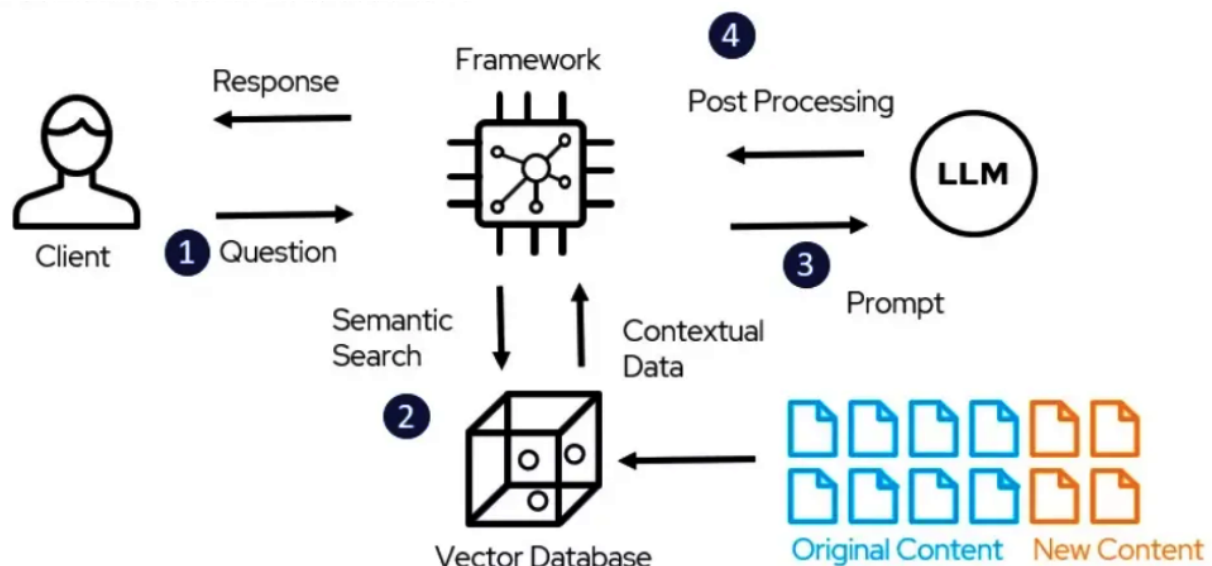In addition to enabling rich interaction with uploaded documents, the tool can:

- Serve as a **testing platform** to evaluate RAG performance across various settings by switching LLMs, embedding models, chunk sizes, and LLM parameters.

- Facilitate conversations without document uploads, allowing users to interact with the LLM directly for general conversational purposes.

# System Design

## Architecture and Workflow

The RAG application employs a carefully structured architecture to balance efficient data processing with a smooth user experience. The architecture includes:

- **Document Parsing and Chunking:** Uploaded documents are parsed into text, split based on user-configurable chunks and overlap sizes, and then processed for embeddings. Recursive splitting enables the chatbot to work effectively with large documents.

- **Dynamic Vector Database Creation:** A unique vector database is created when a document is uploaded or the embedding model is changed. This database stores chunks of text as vectorized embeddings, allowing for fast similarity searches on user queries.

- **Embeddings and Semantic Search:** Semantic embeddings from the selected model are used to compare the query with document segments, enabling contextual responses.

- **LLM Query Processing:** After retrieving the relevant document segments, the app passes these to the LLM. Various configurable parameters fine-tune response generation to meet user requirements.

# Innovative Features

- **Parameter Flexibility:** The user can adjust parameters, such as chunk size and model settings, allowing for deep experimentation and insights into how different models perform on specific tasks. (Figure 1)

- **Side-by-Side Comparison:** The app enables users to view parsed document content alongside the original file, ensuring transparency in how the LLM interprets the text. (Figure 2)

- **Standalone Chatbot Mode:** Beyond document analysis, the tool can also engage in open-ended conversation, demonstrating its flexibility. (Figure 3)

**Fig 1. -
Parameter Flexibility**



**Fig 2. -
Side-by-Side Comparison**



**Fig 3. -
Standalone Chatbot Mode**

# Technical Implementation

## AI Model Quality

- **Embedding Models:** The app uses Ollama's embedding models to vectorize document chunks, enabling the chatbot to retrieve contextually relevant information with high accuracy.

- **Language Model (LLM):** Various LLMs are available for query processing, allowing users to test and optimize for response quality. Each LLM can be fine-tuned with parameters such as temperature, top-k sampling, and penalties, refining response generation.

## Code Efficiency

- **Session State Management:** By leveraging Streamlit's `st.session_state`, the app optimizes performance by reusing existing databases and embedding models where possible. This approach reduces redundant processing and improves runtime efficiency.

- **Efficient Data Handling with ChromaDB:** ChromaDB manages vector embeddings and retrieval, ensuring fast and accurate responses even for large documents.

## Seamless Integration with APIs

The app's seamless interaction with Ollama and ChromaDB APIs highlights the system's adaptability. Document uploads, embeddings, LLM queries, and vector searches occur with minimal latency, ensuring a responsive user experience.

## Source code

The code is published on GitHub at,
https://github.com/khamsakamal48/Local-RAG-with-Ollama/

## Pre-requisites

The application expects that Ollama is installed on your environment and accessible on http://localhost:11434 for API connections. Please refer to the documentation of Ollama as mentioned on https://ollama.com/ to install it in your environment.

Load at least one Large Language Model and appropriate Embedding Model from https://ollama.com/library

# User Experience

## User-Friendliness and Interactivity

The app's layout and sidebar options make it easy to navigate:

- **Intuitive Sidebar Options:** Model selection, chunk size, overlap, and LLM parameters are available in the sidebar, enabling real-time adjustments.

- **Document Upload and Comparison:** Users can upload documents and view parsed content side-by-side with the original file, promoting trust in the AI's data handling.

- **Flexible LLM Configuration:** The extensive range of parameters allows users to customize LLM output. These parameters include temperature, top-k, and top-p, making it ideal for those testing response precision.

**Asking ChatBot:**
**"Why is the sky blue? Respond in one sentence."**



Used LlaMa 3.2 as LLM with temp: 0.1 and top-p: 0.9 & max_length: 50

# Clear and Contextual Responses

The app provides clear responses relevant to user queries, emphasizing content extracted from the document or conversational context. For example, users querying about specific sections of the document receive responses directly related to those sections. This clarity makes the chatbot highly usable for document analysis.

**Asking ChatBot:**
**"Summarise the <u>document</u> in 200 words with bullet points in a tone and language that even a high school student can understand."**



**LLM:** LlaMa 3.2, **Embedding Model:** Mixed Bread AI, **Chunk Size:** 7500
**Document:** https://www.cse.iitb.ac.in/~pb/papers/mts23-maml.pdf

# Challenges Faced

## Embedding Dimension Mismatch

Switching embedding models with differing vector dimensions initially resulted in dimension mismatch errors. Addressing this required dynamically resetting the vector database whenever embedding models were changed to ensure consistency.

## Efficient Session Management

Efficiently managing session states was crucial due to the app's flexible parameter settings. Implementing logic to reset or retain databases and session states depending on the changes made (e.g., document upload, embedding model switch, chunk size change) improved efficiency and resolved frequent database errors.

## Parameter Tuning for Optimal Response Generation

The variety of LLM parameters posed a challenge in maintaining response quality. Iterative testing and feedback loops helped optimize defaults that balance coherence, relevance, and user control in response generation.

# Conclusion

This RAG-based chatbot application leverages the power of Ollama's embedding models and LLMs to provide a robust document analysis and conversational tool. By enabling dynamic parameter adjustment, the app allows users to explore RAG's full potential and engage with uploaded documents or directly with the LLM for broader interactions. The carefully designed user interface, efficient session management, and seamless API integration make it a highly adaptable solution suitable for a wide range of applications in document analysis and interactive AI.

This tool demonstrates the versatility of RAG and its potential in real-world applications, making it an invaluable resource for both technical exploration and end-user engagement.