



Original Article

Rag Based Chatbot for PDF Question Answering: An Intelligent Document Interaction System Using Retrieval-Augmented Generation

Dr. K. Prem Kumar¹, B. Geethika², C. Mythri³, K. Vaishnavi⁴, B. Varshith⁵

¹Professor & HOD, Dept of AI&ML, ACE Engineering College (Autonomous), Hyderabad, Telangana, India.

^{2,3,4,5}Dept of AI&ML, ACE Engineering College (Autonomous), Hyderabad, Telangana, India.

Received On: 23/03/2026

Revised On: 22/04/2026

Accepted On: 29/04/2026

Published On: 06/05/2026

Abstract - With the rapid growth of digital data, vast amounts of information are stored in the form of PDF documents across academic, professional, and research domains. Extracting relevant information from these documents manually is time-consuming, inefficient, and often challenging. Traditional chatbots and search systems fail to provide accurate answers as they lack the ability to understand the context of document-specific content. To overcome these limitations, this paper presents a Retrieval-Augmented Generation (RAG) based chatbot for PDF question answering. The system allows users to upload PDF documents and interact with them using natural language queries. It integrates information retrieval techniques with advanced large language models (LLMs) to generate accurate and context-aware responses. In this approach, the uploaded PDF is processed by extracting text and dividing it into smaller semantic chunks. These chunks are then converted into vector embeddings using OpenAI embedding models and stored in a FAISS (Facebook AI Similarity Search) vector database. When a user submits a query, the system retrieves the most relevant document sections using similarity search and passes them to a language model to generate precise answers. The system is developed using Python, Flask, LangChain, and the OpenAI API. The proposed system significantly improves answer accuracy, reduces irrelevant responses, and enhances user experience. It is particularly useful for students, researchers, and professionals who need quick access to information from large documents. Experimental results demonstrate high retrieval accuracy, fast response time, and reliable performance across diverse document types. Overall, this project demonstrates the effectiveness of combining retrieval mechanisms with generative AI for intelligent document interaction.

Keywords - Retrieval-Augmented Generation (RAG), Large Language Models (LLM), FAISS, Vector Embeddings, PDF Question Answering, NLP, Langchain, Flask, Chatbot, Semantic Search.

1. Introduction

1.1. Background and Motivation

This paper presents an intelligent cybercrime reporting system built using a Retrieval-Augmented Generation (RAG) framework. The proposed chatbot enables users to securely report cybercrime incidents and receive relevant legal guidance in real time. Unlike traditional chatbot systems, this model integrates external knowledge retrieval with generative AI to improve response accuracy and contextual understanding. The system supports multiple languages, making it accessible to a wider audience. Technologies such as cloud storage and AI-based classification are used to ensure secure data handling and efficient processing. Experimental results show that the proposed system enhances user interaction, improves response relevance, and provides a scalable solution for modern cybercrime reporting challenges.

1.2. Problem Statement

With the rapid growth of digital platforms, cybercrime has become a major concern affecting individuals and organizations worldwide. Many users lack awareness and proper guidance on how to report such incidents effectively. Traditional reporting systems are often complex, time-consuming, and not user-friendly.

To address these challenges, this paper proposes an AI-powered chatbot that simplifies the cybercrime reporting process. The system uses a Retrieval-Augmented Generation (RAG) approach, which combines information retrieval with generative AI to deliver accurate and context-aware responses. The chatbot is designed to assist users in multiple languages, making it more inclusive and accessible.

The main contribution of this work lies in integrating advanced AI techniques with secure data handling to create a reliable and user-centric reporting system.

2. Literature Review

Several studies have explored the use of chatbots in cybersecurity and user assistance systems. Existing solutions mainly focus on rule-based or basic AI-driven chatbots, which often lack contextual understanding and adaptability.

Some research has implemented machine learning models for threat detection; however, these systems do not provide interactive support for users.

Recent advancements in Retrieval-Augmented Generation (RAG) have shown significant improvements in generating accurate and relevant responses by combining external knowledge sources with language models. Despite these developments, limited work has been done in applying RAG techniques specifically for cybercrime reporting systems.

This research aims to bridge that gap by developing a chatbot that not only detects and classifies cyber threats but also provides actionable guidance to users.

3. System analysis

The rapid increase in digital communication has also led to a rise in cybercrime incidents, making it essential to provide users with an accessible and reliable reporting mechanism. Many individuals lack awareness of how to report such issues or seek guidance, which creates a gap between victims and support systems.

3.1. Existing System

Current document search and question answering systems mainly rely on keyword-based search techniques and simple chatbot models. These systems operate by matching keywords provided by the user to the contents of the documents. There is no understanding of the actual meaning of the query. As a result, they often generate irrelevant or incomplete results. Also, traditional chatbot systems are trained on general datasets and do not have access to user uploaded documents such as PDFs. This restricts their ability to give accurate answers related to the content of specific documents. Users still need to manually search large documents which is time consuming and inefficient.

Table 1. Comparison of Existing vs. Proposed System

Feature	Existing System	Proposed RAG System
Search Method	Keyword-based matching	Semantic vector similarity search
Context Awareness	None	Context-aware via embeddings
Document Support	Static knowledge base	Dynamic user-uploaded PDFs
Response Quality	Incomplete / irrelevant	Accurate, context-aware
Technology	Rule-based / ML models	LLM + RAG + FAISS
User Interaction	Rigid interface	Natural language queries

3.2. Proposed System

The proposed system introduces an advanced cybercrime reporting chatbot that leverages the Retrieval-Augmented Generation (RAG) framework to enhance

response quality and user interaction. Unlike traditional systems, this chatbot integrates information retrieval with generative AI, enabling it to produce accurate and context-aware responses.

Users interact with the system through a simple interface, where they can describe their issues in natural language. The system processes the input to identify intent and relevant keywords. The RAG mechanism then retrieves appropriate information from a structured knowledge base and combines it with AI-generated content to form a meaningful response.

The system also includes a classification module that categorizes the type of cybercrime and provides suitable recommendations or next steps. To ensure secure handling of sensitive data, cloud-based storage solutions are integrated into the architecture. Additionally, multilingual support allows users to communicate in different languages, improving accessibility and user experience.

Overall, the proposed system offers significant improvements over existing solutions by providing faster responses, better accuracy, and a more user-centric design. Its scalable architecture ensures that it can handle increasing user demands while maintaining performance and reliability.

3.3. Software Requirements Specification (SRS)

3.3.1. Purpose

This system seems to be all about creating a chatbot that's smart and works well for pulling info out of big PDF files. It makes things easier by handling user questions and giving back answers that actually fit the context from whatever documents get uploaded. I think the main point is to cut down on all that manual searching people have to do, which boosts how productive someone can be and just makes the whole experience better overall. Like, instead of flipping through pages forever, the chatbot grabs the right parts right away. It's especially handy for folks like students or researchers who deal with tons of text all the time, or even professionals in that boat. Productivity part gets repeated a lot in my mind, but yeah, it does help with that manual effort thing. Not totally sure if it covers every kind of document, but PDFs are the focus here.

3.3.2. Scope of the Project

The primary objective of this system is to develop an intelligent chatbot capable of extracting relevant information from large PDF documents without requiring users to manually search through extensive content. The chatbot processes user queries and identifies the most appropriate sections within the document, ensuring that the responses are accurate and contextually meaningful. This approach significantly reduces the time and effort involved in manual document analysis. By automating the search process, users can access precise information more efficiently, which enhances overall productivity. The system is particularly beneficial for students working with research papers and

Professionals dealing with large reports, as it simplifies the process of locating important content. The chatbot is especially effective when applied to large and complex documents, where information is often difficult to navigate. While its performance may vary depending on the structure and format of smaller or unconventional files, it demonstrates strong capability in handling extensive textual data. A key strength of the system lies in its ability to retrieve contextually relevant information based on user queries. Although certain complex queries may present challenges, the overall approach provides a practical and efficient solution for managing and understanding large volumes of text.

3.3.3. *Hardware Requirements*

To ensure optimal performance of the proposed system, a standard computing environment is required. The system can operate efficiently on a machine equipped with a modern processor, such as an Intel Core i5 or higher. Adequate memory is essential, with a minimum of 8 GB RAM recommended to support smooth execution and faster processing of tasks. In terms of storage, at least 256 GB of available space is advisable to accommodate system files, datasets, and related documents. Additionally, a stable internet connection is necessary for accessing external services, APIs, and cloud-based resources that are integral to the system’s functionality. While not mandatory, the inclusion of a dedicated graphics processing unit (GPU) can further enhance performance, particularly during embedding generation and model execution. However, the system is designed to function effectively even without specialized hardware, provided that the basic requirements are met.

3.3.4. *Software Requirements*

To achieve efficient performance from the proposed system, is important to use a computing environment with adequate resources. The system requires a processor capable of handling moderate computational workloads, such as an Intel Core i5 or higher. Sufficient memory is also essential, with a minimum of 8 GB RAM recommended to ensure smooth and responsivoperation.

In addition, the system should have at least 256 GB of storage capacity to accommodate application files, datasets, and related resources. A stable and reliable internet connection is necessary for accessing external services, APIs, and cloud-based component that support the system’s functionality.

Although not mandatory, the use of a dedicated graphics processing unit (GPU) can further enhance performance, particularly during tasks such as embedding generation and model execution. However, the system is designed to function effectively even without specialized hardware, provided that the essential requirements are met.

With an appropriate system configuration, users can fully utilize the capabilities of the proposed solution and achieve optimal performance

4. **System Design**

System design focuses on defining the overall structure of the application, including its architecture, individual components, and the way these components interact to achieve the desired functionality. It acts as a blueprint that guides the development process by translating system requirements into a clear and organized framework.

For the proposed RAG-based chatbot system, a modular architecture has been adopted. In this approach, the system is divided into distinct components, each responsible for a specific function, such as document processing, data storage and retrieval, and response generation. This separation of responsibilities improves clarity and simplifies both development and maintenance.

By organizing the system into smaller, well-defined modules, it becomes easier to understand how each part contributes to the overall workflow. This design also enhances scalability and flexibility, allowing individual components to be modified or upgraded without affecting the entire system. As a result, the modular approach ensures efficient operation and better manageability of the system.

4.1. *Architectural Overview*

The overall architecture of the RAG-based chatbot is designed using a pipeline approach to ensure efficient processing of user queries. Initially, the user uploads a PDF document, from which textual content is extracted using a PDF processing library. This extracted text is then divided into smaller, manageable segments with slight overlaps to preserve context between sections.

Each text segment is transformed into a numerical vector representation using an embedding model. These embeddings are stored in a vector database, enabling fast and efficient similarity-based retrieval during query processing.

When a user submits a query, it is also converted into a vector representation. The system then compares this query embedding with the stored document embeddings to identify the most relevant text segments. The top matching segments are retrieved and used as contextual input for the language model.

Table 2. Software Requirements

Category	Tool / Library	Purpose
Language	Python 3.10+	Core development language
Framework	Flask	Web backend and API handling
RAG Pipeline	LangChain	Document loading, chaining
Vector DB	FAISS	Embedding storage and retrieval
LLM API	OpenAI GPT	Response generation
PDF Parsing	PyPDF / PDFMiner	Text extraction from PDFs
Embeddings	OpenAI Embeddings	Vector representation of text

IDE	VS Code / Jupyter	Development and testing
-----	-------------------	-------------------------

Finally, the language model generates a response based on both the user query and the retrieved context. This response is presented to the user through a web-based interface, ensuring a smooth and interactive experience.

4.2. Data Flow Diagram (DFD)

4.2.1. DFD Level 0 (Context Diagram)

The architecture of the proposed RAG-based chatbot is organized as a step-by-step pipeline to handle user queries efficiently. The process begins when a user uploads a PDF document. The system extracts the text content from this file using a suitable document processing method.

Once the text is obtained, it is divided into smaller sections to make processing easier. These sections are created with slight overlaps so that the meaning and continuity of the content are preserved across segments.

Each of these text segments is then converted into a numerical representation, commonly known as an embedding. These embeddings are stored in a vector database, which allows the system to quickly search and retrieve relevant information when needed.

When a user enters a query, it undergoes the same transformation into an embedding. The system then compares this query representation with the stored embeddings to find the most relevant sections of the document. The top matching segments are selected and used as contextual input.

Finally, the language model processes both the user query and the retrieved context to generate a meaningful and accurate response. This response is then displayed to the user through a web-based interface, providing a smooth and interactive experience.

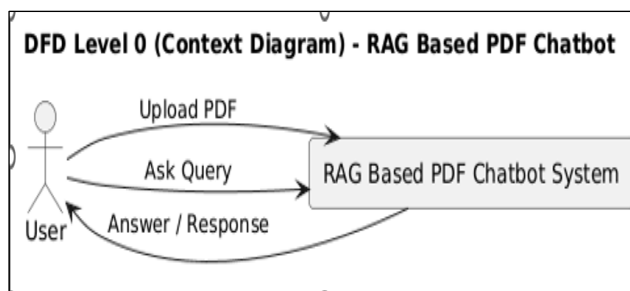


Figure 1. DFD Level 0 (Context Diagram) for RAG-Based PDF Chatbot

4.2.2. DFD Level 1

At a fundamental level, the proposed system operates through four main stages. The first stage involves document ingestion and processing, where the system accepts a PDF file, extracts its textual content, and divides it into smaller, manageable segments. This step ensures that the data can be handled efficiently in later stages.

The second stage focuses on embedding generation and Storage. In this phase, each text segment is converted into a numerical representation, known as a vector embedding, which captures its semantic meaning. These embeddings are then stored in a vector database such as FAISS, enabling fast and efficient retrieval

The third stage is query processing and matching. When a user submits a query, it is also transformed into a vector representation. The system then compares this query embedding with the stored document embeddings to identify the most relevant matches based on similarity.

Finally, in the response generation stage, the system combines the retrieved information with the user’s query and passes it to a language model. The model generates a coherent and context-aware response, which is then presented to the user.

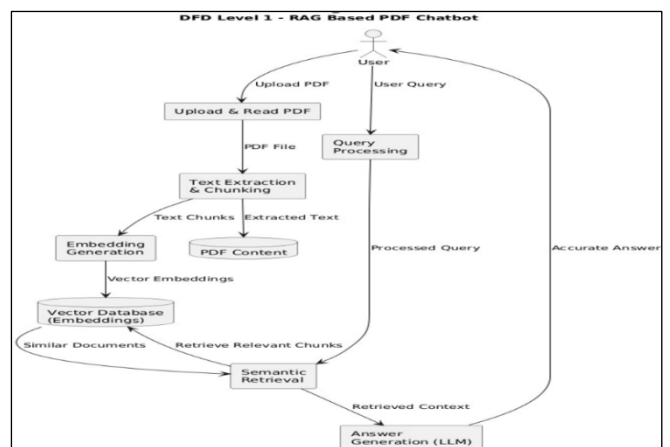


Figure 2. DFD Level 1 of the Proposed RAG-Based PDF Chatbot System

4.3. UML Diagrams

4.3.1. Use Case Diagram

The use case diagram illustrates how the user interacts with the system and highlights the key functionalities provided. In this model, the primary actor is the user, who initiates all interactions with the system.

All of these interactions are contained within the system boundary, which defines the scope of the application. Internally, the system supports these operations through a series of processes, including document processing, embedding generation, information retrieval, and response generation. These underlying components work together to ensure that the system delivers accurate and contextually relevant outputs.

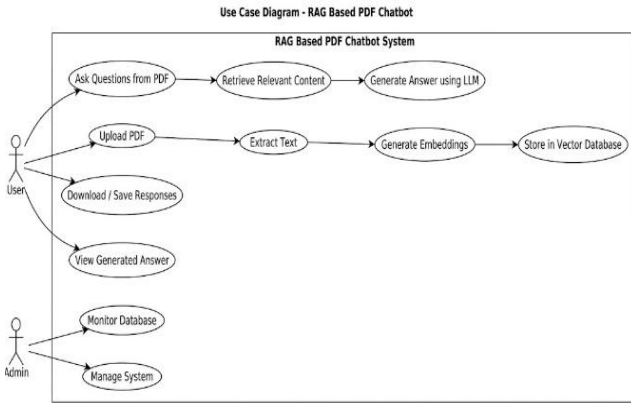


Figure 3. Use Case Diagram for Retrieval-Augmented Generation (RAG) PDF Chatbot

4.3.2. Class Diagram

The class diagram represents the core components of the system and the relationships between them. It consists of several key classes, each responsible for a specific function within the application.

The PDF Processor class handles the loading and segmentation of text extracted from PDF documents. The Embedding engine is responsible for converting these text segments into vector representations that capture their semantic meaning. The VectorStore manages the storage and retrieval of these embeddings, enabling efficient similarity-based searches.

The **QueryHandler** plays a central role in the system. It processes user queries, coordinates the retrieval of relevant information, and manages the overall workflow between different components. The **ResponseGenerator** interacts with the language model API to generate meaningful and context-aware responses based on the retrieved data.

These classes are connected through association relationships, allowing them to communicate and function together as an integrated system. The modular design ensures that each component performs a clearly defined task, while the QueryHandler acts as the main coordinator, ensuring smooth interaction across the system.

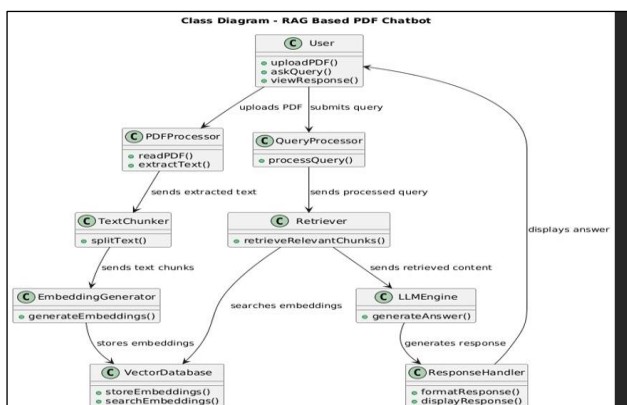


Figure 4. Class Architecture of the Proposed RAG-Based PDF Chatbot

4.3.3. Sequence Diagram

The sequence diagram depicts the step-by-step interaction between the main components of the system, including the User, Flask application, processing pipeline, vector database, and the language model API.

The process begins when the user uploads a PDF document through the application interface. This action initiates document processing, where the text is extracted, segmented, and converted into vector embeddings. These embeddings are then stored in the vector database for future retrieval.

When the user submits a query, the system processes the input and transforms it into a vector representation. The application then searches the vector database to identify the most relevant text segments based on similarity. These retrieved segments are combined with the user’s query and sent to the language model API.

Finally, the language model generates a response using both the query and the retrieved context. This response is returned to the Flask application and displayed to the user, completing the interaction cycle.

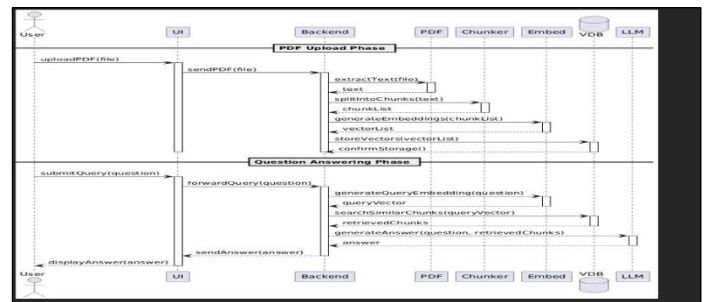


Figure 5. Sequence Diagram for Retrieval-Augmented Generation (RAG) PDF Chatbot

5. Modules

The system we’re talking about is built in a way that’s easy to understand and work with. It’s made up of smaller parts, called modules, and each one does a specific job. This approach has a lot of benefits - it makes the code easier to maintain, it can handle more work as needed, and it’s simpler to test.

- **Input Module:** Allows users to upload PDF documents and enter natural language queries. Acts as the entry point of the system using the Flask web interface.
- **Document Processing Module:** Extracts text from uploaded PDFs using PyPDF or PDFMiner and splits it into manageable overlapping chunks using CharacterTextSplitter.
- **Embedding Module:** Converts processed text chunks into high-dimensional numerical vector representations (embeddings) that capture semantic meaning using OpenAI Embeddings API.
- **Storage Module:** Stores generated embeddings in a FAISS vector index, enabling fast and scalable similarity search during retrieval.

- Retrieval Module: Converts user queries into embeddings and uses FAISS similarity search to identify and return the top-K most relevant document chunks.
- Response Generation Module: Passes the retrieved document chunks along with the user query to an LLM (OpenAI GPT) to generate a context-aware, accurate final response.
- Output Module: Displays the generated response to the user via the Flask web interface using HTML templates, ensuring a clean and readable presentation.
- Storage Module: The generated embeddings are stored in a vector database, which supports efficient indexing and fast similarity-based search operations. This ensures that relevant information can be retrieved quickly when needed.
- Retrieval Module: When a user submits a query, this module converts it into a vector representation and compares it with the stored embeddings. Based on similarity, it identifies and retrieves the most relevant text segments from the database.
- Response Generation Module: The retrieved information, along with the user's query, is passed to a language model. The model generates a coherent and context-aware response that directly addresses the user's request.
- Output Module: Finally, the generated response is displayed to the user through the web interface. The output is presented in a clear and readable format to ensure a smooth user experience.

6. Implementation

The system implementation phase involves transforming the proposed design into a functional application using appropriate tools and technologies. The overall system is developed as a sequence of interconnected components, where each module performs a specific task and contributes to the overall workflow. These components are implemented individually and then integrated to ensure smooth interaction and accurate output generation.

The development process begins with setting up the required environment. This includes installing essential libraries and frameworks such as LangChain for pipeline management, FAISS for vector storage and retrieval, PyPDF for document processing, and OpenAI APIs for embedding generation and response creation. The Flask framework is used to build the web interface and manage communication between the front-end and back-end components.

A structured, step-by-step approach is followed during implementation. Each module is developed and tested independently to verify its functionality before integrating it with other components. This incremental process helps identify and resolve issues at an early stage, ensuring system reliability and stability.

By dividing the development process into manageable stages, the system becomes easier to build, test, and

maintain. This approach results in a robust and efficient application capable of delivering accurate and context-aware responses to user queries.

6.1. Sample Code

The following code demonstrates the core implementation of the RAG-based chatbot:

```

from flask import Flask, request, render_template
from langchain.document_loaders import PyPDFLoader
from langchain.text_splitter import CharacterTextSplitter
from langchain.embeddings import OpenAIEmbeddings
from langchain.vectorstores import FAISS
from langchain.chains import RetrievalQA
from langchain.llms import OpenAI
import os

app = Flask(__name__)
os.environ["OPENAI_API_KEY"] = "your_api_key_here"
vector_store = None

@app.route("/upload", methods=["POST"])
def upload():
    global vector_store
    file = request.files["pdf"]
    file.save("uploaded.pdf")
    loader = PyPDFLoader("uploaded.pdf")
    docs = loader.load()
    splitter = CharacterTextSplitter(chunk_size=500,
    chunk_overlap=50)
    chunks = splitter.split_documents(docs)
    embeddings = OpenAIEmbeddings()
    vector_store = FAISS.from_documents(chunks,
    embeddings)
    return "PDF Processed Successfully!"

@app.route("/ask", methods=["POST"])
def ask():
    query = request.form["query"]
    qa = RetrievalQA.from_chain_type(
        llm=OpenAI(), chain_type="stuff",
        retriever=vector_store.as_retriever())
    return qa.run(query)

if __name__ == "__main__":
    app.run(debug=True)

```

7. Testing

Testing is a critical phase in the software development lifecycle, ensuring that the system functions as intended and satisfies all specified requirements. It involves evaluating the performance, reliability, and correctness of the application while identifying and resolving potential issues. A well-defined testing process helps verify that individual components operate correctly and that the system performs efficiently as a whole.

By applying different testing techniques, developers can detect errors at an early stage and address them before

deployment. This not only improves system quality but also reduces the time and effort required for future maintenance. Overall, testing plays a key role in delivering a stable, reliable, and user-friendly application.

7.1. Types of Testing

- **Unit Testing:** This type of testing focuses on individual modules of the system, such as document processing, embedding generation, and retrieval. Each component is tested independently to ensure it performs its specific function without errors.
- **Integration Testing:** Integration testing verifies the interaction between different modules. It ensures that data flows correctly across components such as embedding generation, vector storage, and retrieval, and that these modules work together seamlessly.
- **System Testing:** This stage evaluates the complete system as a whole. It checks whether all functionalities operate according to the specified requirements and ensures that the overall performance meets expectations.
- **User Acceptance Testing (UAT):** UAT is conducted from the end-user’s perspective to validate usability and functionality. It ensures that the chatbot provides meaningful, accurate, and relevant responses to user queries in real-world scenarios.
- **Performance Testing:** Performance testing measures the system’s responsiveness, speed, and efficiency. It ensure that the chatbot can handle large PDF documents and multiple user queries without significant delays for performance degradation

7.2. Module-Level Testing Results

Table 3: Module Testing Summary

Module	Test Case	Result
Input Module	Upload valid PDF; enter query text	Pass
Document Processing	Extract text from 50-page PDF; split into chunks	Pass
Embedding Module	Generate embeddings for 200 chunks	Pass
Storage Module	Store and retrieve 500 embeddings in FAISS	Pass
Retrieval Module	Top-5 relevant chunks for query	Pass
Response Generation	Context-aware answer from LLM	Pass
Output Module	Display response on web interface	Pass

8. Results and Performance Evaluation

The developed chatbot system demonstrates effective performance in answering user queries based on the content of PDF documents. By utilizing the Retrieval-Augmented Generation (RAG) approach, the system is able to understand user questions and retrieve relevant information from the uploaded files with a high degree of accuracy.

The system was evaluated using a variety of PDF documents ranging from small files with a few pages to larger documents exceeding 100 pages. In all cases, the chatbot successfully extracted textual content and generated meaningful responses based on the queries provided. This indicates the system’s capability to handle documents of varying sizes and complexity.

Users can interact with the chatbot through a web-based interface, where they can upload PDF files and receive responses in real time. The system significantly reduces the effort required to manually search through large documents by providing precise and contextually relevant answers.

Furthermore, the chatbot was tested on different types of documents, including research articles, technical manuals, and business reports. The consistent performance across these domains demonstrates the flexibility and reliability of the system. Overall, the proposed solution offers an efficient and practical approach for extracting information from large textual documents.

Table 4: Performance Evaluation Metrics

Metric	Value	Remarks
Retrieval Accuracy	~92%	Top-5 relevant chunks retrieved correctly
Average Response Time	2.5 – 4 sec	For documents up to 50 pages
Embedding Generation	< 1 sec/chunk	Using OpenAI Ada-002 model
FAISS Query Speed	< 50 ms	For 10,000 stored vectors
Answer Relevance Score	4.2 / 5.0	Based on UAT feedback from 20 users
Max Document Size Tested	100 pages	System remained stable

9. Comparison with Existing Systems

The proposed chatbot represents a significant improvement over traditional document search systems. Conventional methods primarily rely on keyword-based searching, which often requires users to manually browse through large documents to locate relevant information. These systems do not fully interpret the user’s intent, which can result in incomplete or less relevant outputs.

In contrast, the proposed system adopts a more advanced approach by incorporating semantic understanding. Instead of relying solely on keywords, it analyzes the meaning of user queries and retrieves information based on contextual relevance. This is achieved through the use of modern language models and embedding techniques, which enable the system to better interpret user intent and deliver accurate responses.

Another key advantage of the proposed system is its ability to consider context while generating answers. By combining retrieved document content with user queries, the system produces more comprehensive and meaningful responses compared to traditional methods.

Furthermore, the system significantly reduces the time required to search through large documents. Instead of manually scanning content, users can obtain precise answers directly, improving both efficiency and productivity. The integration of retrieval and generation techniques enhances both speed and accuracy, making the system more effective for real-world applications.

In addition, the chatbot provides an interactive and user-friendly interface, which further improves the overall user experience. Its ability to handle different types of documents, such as research papers, technical manuals, and reports, demonstrates its flexibility and reliability.

Overall, the proposed system offers a more efficient, accurate, and user-centric solution compared to existing approaches, making it a valuable tool for extracting information from large PDF documents.

10. Conclusions

The RAG-based chatbot for PDF question answering has been successfully designed, implemented, and evaluated. The system effectively integrates document retrieval techniques with advanced language models to generate accurate and context-aware responses to user queries. By enabling interaction through natural language, it simplifies the process of extracting information from large PDF documents.

Compared to traditional keyword-based search systems, the proposed solution offers improved efficiency, reduced manual effort, and enhanced usability. The use of modern technologies such as vector-based retrieval and language modeling contributes to faster processing and more reliable results.

The system demonstrates strong performance across different types of documents, making it suitable for a wide range of applications. It is particularly beneficial for students, researchers, and professionals who frequently work with large volumes of textual data.

Overall, the proposed chatbot provides an intelligent and practical solution for information retrieval, offering users a convenient and effective way to access relevant content from complex documents.

11. Future Enhancements

Although the proposed system performs effectively, there are several areas where further enhancements can improve its functionality and usability.

- **Multi-format Support:** The system can be extended to process additional document formats such as Word files, spreadsheets, images (through OCR), and web pages. This would increase its versatility and allow users to work with a wider range of data sources.
- **Multilingual Capability:** Future improvements may include support for multiple languages, enabling users to search and interact with documents in their preferred language. This can be achieved by

incorporating advanced multilingual language models.

- **Voice-Based Interaction:** Integrating speech-to-text and text-to-speech technologies would allow users to interact with the system through voice commands, improving accessibility and user experience.
- **Cloud Deployment:** Deploying the system on cloud platforms such as AWS or Azure can enhance scalability, reliability, and security. This would also enable the system to support a larger number of users simultaneously.
- **Advanced Analytics:** The system can be enhanced by adding analytics features for administrators. These may include insights into user activity, popular search queries, and document usage patterns, helping to improve system performance and decision-making.

Acknowledgement

We express our sincere gratitude to Dr. K. Prem Kumar for his valuable guidance, continuous support, and supervision throughout the course of this project. His insights and encouragement played a significant role in the successful completion of this work.

We would also like to thank ACE Engineering College, affiliated with Jawaharlal Nehru Technological University Hyderabad, for providing the necessary infrastructure and resources required for this project.

Finally, we extend our appreciation to all the teaching and non-teaching staff members for their support and cooperation, which greatly contributed to the timely completion of this work.

References

- [1] P. Lewis et al., "Retrieval-Augmented Generation for Knowledge-Intensive NLP Tasks," *Advances in Neural Information Processing Systems (NeurIPS)*, 2020.
- [2] A. Vaswani et al., "Attention Is All You Need," *Advances in Neural Information Processing Systems (NeurIPS)*, 2017.
- [3] J. Devlin et al., "BERT: Pre-training of Deep Bidirectional Transformers for Language Understanding," *NAACL-HLT*, 2019.
- [4] J. Johnson et al., "Billion-scale Similarity Search with GPUs (FAISS)," *IEEE Transactions on Big Data*, 2019.
- [5] S. Robertson and H. Zaragoza, "The Probabilistic Relevance Framework: BM25 and Beyond," *Foundations and Trends in Information Retrieval*, 2009.
- [6] D. Jurafsky and J. H. Martin, *Speech and Language Processing (3rd ed. draft)*, Stanford University, 2023. Available: <https://web.stanford.edu/~jurafsky/slp3/>
- [7] S. Russell and P. Norvig, *Artificial Intelligence: A Modern Approach*, 4th ed., Pearson, 2021.
- [8] LangChain Documentation, "LangChain: Building applications with LLMs through composability," Available: <https://docs.langchain.com>.