# PDF Intelligence Agent Using RAG

[1] Ambareesha V, [2] Puja Jaiswal, [3] Dr. Jessy Sujana. G

[1][2][3] Department of Computer Science & Engineering – Emerging Technologies, SRM Institute of Science & Technology, Vadapalani, Chennai, India
Corresponding Author Email: [1] av8056@srmist.edu.in, [2] pj8077@srmist.edu.in, [3] Jessysug@srmist.edu.in

*Abstract— This paper introduces an advanced system to enhance PDF-based chatbots by integrating LangGraph, a powerful tool for constructing resilient language agents through a graph-based framework. The system improves information retrieval and response generation by utilizing an iterative feedback loop in conjunction with a Retrieval-Augmented Generation (RAG) model. By incorporating a relevancy threshold and ranking retrieved content, the system continuously adapts, learning from user feedback to refine its responses. The proposed approach ensures better retrieval accuracy and user experience across industries, where precise information extraction from complex documents is essential. The architecture is designed to overcome challenges associated with traditional IR systems, paving the way for innovative applications in various sectors.*

## I. INTRODUCTION

With the exponential growth of digital documents, especially in fields such as legal, healthcare, and academia, the need for efficient information retrieval (IR) systems has become paramount. Traditional PDF-based chatbots are often limited by their reliance on keyword matching and basic retrieval techniques, leading to inaccurate, irrelevant, or incomplete responses. The increasing demand for more efficient and contextually aware chatbots necessitates the exploration of advanced techniques.

This paper presents a novel solution: integrating LangGraph with a Retrieval-Augmented Generation (RAG) system. This combination creates a more robust, context-aware chatbot capable of delivering precise responses while learning from real-time feedback. LangGraph, a graph-based tool, enhances decision-making by structuring information flow, while RAG leverages both retrieval and generation for improved performance. Our approach addresses key challenges, including accuracy, scalability, and efficiency, making it suitable for various document-intensive industries.

The following sections will delve deeper into the advancements in information retrieval systems, the architecture of the proposed solution, experimental methodologies, results, discussions on findings, and future directions for research.

## II. LITERATURE REVIEW

### A. Advances in Information Retrieval Systems

Recent advancements in IR methods have significantly improved retrieval performance, particularly in handling short texts and semantic relationships. According to Hambarde and Proença's work, "Information Retrieval: Recent Advances and Beyond" (2023), modern IR systems benefit from the integration of user feedback, enabling continuous improvements in relevance and accuracy. By moving beyond simple keyword matching, these systems now utilize machine learning techniques to understand semantic relationships between words and documents.

- **Improved Retrieval Performance**

The evolution of IR technologies has led to remarkable enhancements in performance metrics such as precision, recall, and F1 scores. Traditional models primarily focused on exact keyword matches; however, contemporary methods utilize deep learning approaches to comprehend user intent better. Techniques such as word embeddings (e.g., Word2Vec, GloVe) and transformer models (e.g., BERT, GPT) enable systems to grasp contextual meanings and relationships among terms, thus improving the relevance of retrieved results.

- **User Feedback Integration**

Incorporating user feedback has been pivotal in refining IR systems. Feedback mechanisms allow systems to adapt to user preferences over time, learning from interactions to adjust their retrieval strategies. This dynamic adjustment is particularly crucial in domains where the language and requirements can shift rapidly, ensuring that the system remains effective even as user expectations evolve.

- **Challenges: Query Drift and Overfitting**

Despite these advancements, challenges such as query drift and overfitting persist. Query drift occurs when user queries change in meaning or intent over time, which can lead to decreased retrieval accuracy if the system does not adapt accordingly. Overfitting, on the other hand, refers to a model's tendency to perform well on training data but poorly on unseen data, resulting in a lack of generalizability. This emphasizes the need for robust training methodologies and continuous learning mechanisms.

### B. Retrieval-Augmented Generation (RAG) Systems

RAG systems offer a transformative approach to information retrieval by merging traditional retrieval methods with generative models. Gao et al. in their paper "Retrieval-Augmented Generation for Large Language

Models: A Survey" (2023), categorize RAG systems into three main types: Naive, Advanced, and Modular. These systems combine retrieval mechanisms with the generative power of large language models (LLMs) to improve the quality and relevance of responses.

- **Simplicity and Wide Applicability**

One of the primary advantages of RAG systems is their simplicity and wide applicability across different domains. They can effectively manage diverse datasets, providing tailored responses based on user queries. The synergy between retrieval and generation allows for a more comprehensive approach to information retrieval, overcoming limitations of purely generative or retrievalbased systems.

- **Enhanced Component Functionality**

RAG systems enhance the functionality of both retrieval and generation components. Retrieval mechanisms can focus on finding relevant information quickly, while generative models can produce coherent, contextually relevant outputs. This dual functionality ensures that users receive accurate information without compromising on the quality of the generated content.

- **Challenges: Retrieval and Generation Issues**

However, RAG systems are not without challenges. Issues related to retrieval, such as irrelevant or outdated information being presented, can negatively impact the user experience. Additionally, there can be generation issues, where the output may lack coherence or contextual relevance if the underlying retrieval process fails. These challenges necessitate ongoing research and development to refine RAG approaches.

### C. Case Studies in RAG Optimization

The integration of RAG systems in industryspecific applications has yielded positive results. For instance, Fei Liu et al. demonstrated the effectiveness of RAG techniques in the automotive industry through a case study involving locally deployed Ollama models (2023). Their system significantly improved context precision and recall, making it suitable for offline, lowperformance environments.

- **Limitations in Generalizability**

However, the narrow focus on automotive industry documents limits generalizability to other sectors. While the methods may excel in one domain, they might not translate effectively to others due to differences in language, terminology, and document structure. This highlights the need for adaptable models that can accommodate various contexts and industries.

- **Importance of Specialized Knowledge**

Moreover, specialized knowledge is often required for the deployment and optimization of these models. As the complexity of the systems increases, the necessity for expertise in both domain knowledge and technical implementation becomes critical. This requirement can act as a barrier to widespread adoption, particularly in organizations with limited resources or expertise.

### D. Limitations of Traditional PDF Chatbots

Traditional PDF-based chatbots are hampered by several limitations. These systems often struggle with retrieving relevant information from vast datasets, leading to irrelevant or redundant responses. Additionally, as identified by Maryamah et al. in "Chatbots in

Academia: A Retrieval-Augmented Generation Approach for Personalized Learning" (2024), scalability issues pose significant challenges in educational platforms, where the system must cater to personalized learning.

- **Scalability Challenges**

The scalability of chatbots in educational contexts is particularly pressing. As the volume of documents and user queries increases, maintaining efficiency while ensuring accurate retrieval becomes increasingly difficult. Systems that cannot scale effectively may result in longer response times and decreased user satisfaction.

- **Specialized Language Handling**

Furthermore, traditional chatbots often struggle to handle specialized language or jargon found in academic or technical documents. As users seek precise information, the inability of chatbots to navigate complex terminologies can lead to frustration and disengagement. This underscores the necessity for enhanced techniques that can better manage specialized content.

## III. METHODOLOGY

### A. System Architecture

The proposed system is designed to improve the retrieval and response generation of PDFbased chatbots by integrating LangGraph with a RAG system. The architecture includes the following components:

- **PDF Upload and Text Extraction**

Users upload PDFs, and the system extracts the text using OCR and Textract services. This initial processing is critical, as it forms the basis for all subsequent operations.

- **Text Chunking and Vectorization**

Extracted text is split into smaller, manageable chunks, which are vectorized using embedding models like HuggingFace. The vectors are stored in a FAISS database for efficient retrieval, facilitating rapid access to relevant content.

- **User Query Input and Retrieval**

When a user submits a query, the system retrieves relevant text chunks from the database by comparing the vector representations of the query and the documents. This step ensures that only the most relevant information is processed further.
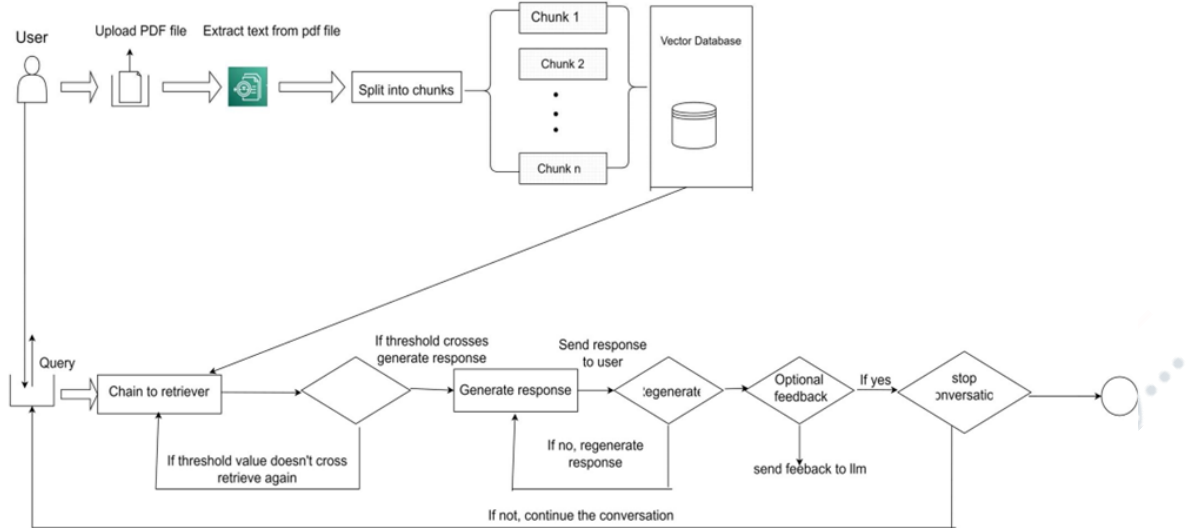
- **Relevancy Check**

A relevancy threshold filters irrelevant results, ensuring that only contextually accurate chunks are considered. This filtering process enhances the quality of responses generated by the system.

- **Response Generation**

The filtered text chunks are passed to the RAG model, which generates a coherent response based on the query and the retrieved content. The generative model adapts the retrieved information into a natural language response.

- **Iterative Feedback Loop**

Users can provide feedback, which is processed to refine future responses, allowing the system to adapt and improve over time. This feedback mechanism ensures that the chatbot evolves to meet user needs more effectively.



## IV. DISCUSSION

### A. Advantages of the Proposed System

The integration of LangGraph and RAG offers several advantages over traditional PDF chatbots.

- **Enhanced Accuracy and Relevancy**

By leveraging both retrieval and generative capabilities, the proposed system can deliver highly accurate and contextually relevant responses. The combination of user feedback integration and structured data flow enables continuous improvements in performance, which is vital in dynamic environments.

- **Scalability and Flexibility**

The modular architecture allows for scalability, accommodating growing document databases and increasing user queries without sacrificing performance. This adaptability is essential for industries that experience fluctuations in data volume, ensuring that the system remains responsive and efficient.

### B. Challenges and Limitations

While the proposed system shows promise, several challenges and limitations warrant consideration.

- **Complexity of Implementation**

The complexity involved in deploying and optimizing the system may pose challenges for organizations with limited resources or expertise. Integrating advanced techniques requires a thorough understanding of both domain-specific knowledge and technical implementation.

- **Handling Specialized Content**

While the system effectively retrieves general information, handling highly specialized content remains a challenge.

Future iterations may need to incorporate domain-specific knowledge bases to ensure accuracy in specialized areas.

### C. Future Research Directions

Future research should focus on addressing these challenges and exploring new ways to further enhance the performance of PDF chatbots. This could include integrating additional techniques, such as sentiment analysis and entity recognition, to provide more nuanced and contextually aware responses.

- **Exploring Advanced Techniques**

Incorporating more advanced machine learning models, such as Vision Transformers (ViTs), can enhance the system's ability to process complex documents. Additionally, investigating the integration of real-time object detection models like YOLO could significantly improve the efficiency of information retrieval and localization.

- **Field Testing and Real-World Applications**

Conducting field tests using autonomous robots equipped with multi-spectral sensors in actual minefields could provide valuable insights into the real-world applicability of the system. This would enhance the robustness and adaptability of the technology, paving the way for more widespread use in high-risk environments.

## V. CONCLUSION

The PDF Intelligence Agent using RAG represents a significant advancement in the field of document retrieval and chatbot development. By leveraging LangGraph's graph-based framework and integrating it with a RAG model, the system addresses many of the challenges associated with traditional PDF chatbots. Its iterative feedback loop and

structured architecture make it adaptable to user needs while also improving over time.

This solution has strong potential for industries that deal with vast, complex documents, offering both efficiency and precision in information retrieval. As organizations increasingly rely on automated systems for information processing, the proposed framework provides a promising path toward enhanced performance in information retrieval tasks.

## REFERENCES

[1] Kailash A. Hambarde, Hugo Proença, "Information Retrieval: Recent Advances and Beyond," July 2023.

[2] Yunfan Gao et al., "Retrieval-Augmented Generation for Large Language Models: A Survey," December 2023.

[3] Fei Liu, Zejun Kang, Xing Han, "Optimizing RAG Techniques for Automotive Industry PDF Chatbots: A Case Study with Locally Deployed Ollama Models."

[4] Iván Ortiz-Garces et al., "Optimizing Chatbot Effectiveness through Advanced Syntactic Analysis: A Comprehensive Study in Natural Language Processing."

[5] Maryamah Maryamah et al., "Chatbots in Academia: A Retrieval-Augmented Generation Approach for Personalized Learning."

[6] Feriel Khennouche et al., "Revolutionizing Generative Pre-Trained Models: Insights and Challenges in Chatbot Development."

[7] Vishal Dutt et al., "Dynamic Information Retrieval with Chatbots: A Review of Artificial Intelligence Techniques."