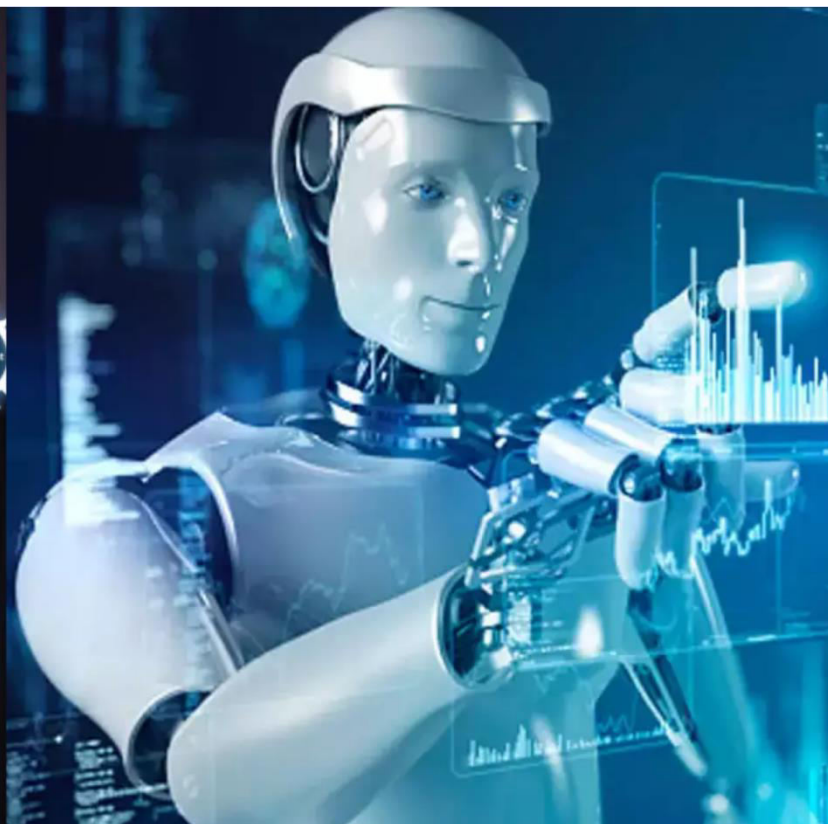




International Journal of Innovative Research in Science Engineering and Technology (IJIRSET)

(A Monthly, Peer Reviewed, Refereed, Scholarly Indexed, Open Access Journal)



Impact Factor: 8.699

Volume 14, Issue 3, March 2025

Chat with PDF AI Tool

Ms. Pooja A.Ghuge¹, Abhishek Kumbhar², Tejas Khetmalis³

Lecturer, Department of Computer Engineering, Marathwada Mitra Mandal Polytechnic, Pune, India¹

Final Year Student, Department of Computer Engineering, Marathwada Mitra Mandal Polytechnic, Pune, India^{2,3}

ABSTRACT: Extracting relevant information from large PDF documents can be a time-consuming and challenging task. To address this, we have developed an AI-powered chat tool that allows users to interact with PDFs in a conversational manner. The system utilizes advanced natural language processing techniques to understand user queries and retrieve relevant information from the document. This paper details the design, implementation, and evaluation of the system, demonstrating its ability to provide accurate and efficient access to information within PDF documents.

I. INTRODUCTION

The digital age has ushered in an explosion of information, much of which is stored and shared in the Portable Document Format (PDF). While PDFs offer a standardized and widely accessible format, extracting specific information from them can be a cumbersome and time-consuming process. 1 Traditional methods, such as manual scanning and keyword searches, often prove inefficient, particularly when dealing with large or complex documents. 2 This challenge is further exacerbated by the diverse nature of PDFs, which can range from text-based documents to scanned images, presenting significant hurdles for automated information retrieval.

The need for a more intuitive and efficient approach to accessing information within PDFs has become increasingly apparent. Recent advancements in artificial intelligence (AI), particularly in natural language processing (NLP) and computer vision, offer promising avenues for addressing this challenge. This paper presents a novel AI-powered tool that enables users to engage in conversational interactions with any PDF document. By combining robust PDF parsing with state-of-the-art NLP techniques, our system transforms static documents into dynamic knowledge sources, allowing users to 'chat' with their PDFs and retrieve relevant information through natural language queries.

II. LITERATURE SURVEY

PDF Parsing and Text Extraction:

Fundamental to any PDF interaction system is the accurate extraction of text and structure. Early approaches relied on rule-based methods and libraries like PDFMiner [1], which provided basic text extraction capabilities. More recent efforts have focused on handling complex layouts and scanned documents. Libraries like PyPDF2 [2] and Apache Tika [3] offer broader support for various PDF formats. However, these tools often struggle with scanned documents, necessitating the integration of Optical Character Recognition (OCR). Tesseract OCR [4] has emerged as a widely used open-source solution, while cloud-based APIs like Google Cloud Vision API [5] offer improved accuracy for complex images. Our work builds upon these advancements by implementing a robust PDF parsing pipeline that incorporates both traditional libraries and advanced OCR techniques to ensure accurate text extraction across diverse PDF formats.

Question Answering Systems and Chatbots:

Question answering (QA) systems have been extensively studied in the context of open-domain and domain-specific knowledge bases. Recent research has explored the application of QA systems to document retrieval, particularly in the context of PDFs. Retrieval-augmented generation (RAG) models have shown great promise in this area. These models combine the strengths of information retrieval and language generation, allowing for the retrieval of relevant document snippets and the generation of coherent answers. Our work leverages these advancements by implementing a RAG-based architecture that enables users to engage in natural language conversations with PDF documents.

Related Tools and Systems:

Several tools and systems have been developed for interacting with PDF documents. Some tools provide basic search and annotation capabilities, while others offer more advanced features like document summarization and question answering. However, many existing tools lack the ability to handle diverse PDF formats, particularly scanned

documents, or require extensive manual configuration. Our system aims to address these limitations by providing a user-friendly and robust solution that can handle a wide range of PDF documents.

Dynamic Embedding Techniques:

Building upon the established principles of semantic embedding and retrieval, this project introduces a significant advancement through the integration of dynamic embedding techniques, as pioneered by Reimers and Gurevych (2019). This approach allows for real-time indexing of PDF documents, ensuring that newly uploaded content is immediately searchable and relevant. This dynamic capability is essential for maintaining an up-to-date and responsive information retrieval system. Furthermore, we employ a hybrid Retrieval-Augmented Generation (RAG) architecture, which strategically balances retrieval accuracy with computational efficiency.

III.METHODOLOGY

TITLE: PDF Parsing and Preprocessing: Text Extraction and Structural Integrity

The methodology begins with parsing and preprocessing PDF documents to extract machine-readable text while preserving layout and context. Text-based PDFs are processed using PyPDF2, while scanned/image-based files utilize Tesseract OCR combined with OpenCV for noise reduction and deskewing. Layout analysis identifies tables (via Camelot) and section headers, ensuring structural elements like bullet points and figures are retained. The extracted text is segmented into overlapping 512-token chunks (128-token stride), with metadata such as page numbers and section titles embedded for citation tracking.

TITLE: Semantic Embedding and Indexing: Transforming Text into Searchable Vectors

At the heart of our system's semantic understanding lies the process of converting text chunks into dense vector representations. We leverage Sentence-BERT (SBERT), a state-of-the-art transformer-based model specifically designed to capture semantic relationships. SBERT's architecture, fine-tuned on siamese and triplet networks, allows it to generate embeddings that accurately reflect the semantic similarity between text segments, a significant improvement over traditional BERT models that are primarily focused on token-level representations. Unlike vanilla BERT, which outputs token-level embeddings, SBERT generates a single 768-dimensional embedding for each text chunk by applying mean pooling to the token-level outputs. This mean pooling operation effectively aggregates the contextual information across the entire sentence or text segment, resulting in embeddings that capture the nuanced, sentence-level context essential for accurate document retrieval. This method ensures that our system can effectively understand and respond to user queries, even when the query and the relevant document content do not share exact lexical matches.

TITLE: Intent Recognition and Hybrid Retrieval: Contextual and Lexical Query Processing

To ensure accurate and contextually relevant responses, user queries undergo a sophisticated analysis process that combines intent recognition with a hybrid retrieval strategy. Initially, each query is analyzed by a BERT-base-uncased classifier, which has been fine-tuned on a dataset of PDF-specific intents. This classification step allows the system to discern the user's underlying intention, such as requesting a summary, retrieving specific information, or asking a question related to the document's content. Once the intent is identified, the query is transformed into a vector embedding, which serves as the foundation for our hybrid retrieval approach. This strategy leverages the strengths of both semantic and lexical search methods. Semantic similarity is assessed using FAISS, where the query embedding is compared to the pre-indexed document embeddings, enabling the retrieval of text snippets based on their contextual relevance. Simultaneously, lexical matching is performed using the BM25 algorithm, which identifies text snippets containing keywords from the user's query, ensuring that important terms are not overlooked. To balance the contextual depth provided by semantic search with the precision of keyword matching, we employ a weighted scoring mechanism. Specifically, the retrieval scores are calculated with a 70% emphasis on semantic similarity and a 30% emphasis on lexical matching. This weighted combination ensures that the system retrieves text snippets that are both semantically relevant and lexically precise, providing users with comprehensive and accurate responses.4. Retrieval-Augmented Generation (RAG): Context-Aware Answer Synthesis

TITLE: Adaptive Learning and System Scalability: Feedback-Driven Optimization

To ensure continuous improvement and adaptability, our system incorporates an adaptive learning mechanism driven by user feedback. Specifically, user interactions, such as ratings and corrections, are leveraged to fine-tune the retrieval

model through triplet loss reinforcement learning from human feedback (RLHF). This iterative process allows the system to learn from its mistakes and refine its understanding of user preferences, leading to progressively more accurate and relevant search results. This feedback loop is crucial for maintaining and enhancing the system's performance over time. To address the demands of scalability and responsiveness, the backend architecture is built upon FastAPI and Celery, enabling asynchronous processing and efficient task management. This design allows for the seamless handling of multiple concurrent requests and ensures that the system can scale horizontally to accommodate increasing user loads. Redis is utilized for caching embeddings, significantly reducing latency and improving the speed of retrieval operations. For robust data management, raw PDF documents and their associated metadata are stored in Amazon S3 and PostgreSQL, respectively. Amazon S3 provides scalable and durable storage for large files, while PostgreSQL offers a reliable and efficient database for structured metadata. This comprehensive approach to adaptive learning and system scalability ensures that our platform can deliver high-performance, continuously improving information retrieval capabilities.

IV.SYSTEM ARCHITECTURE

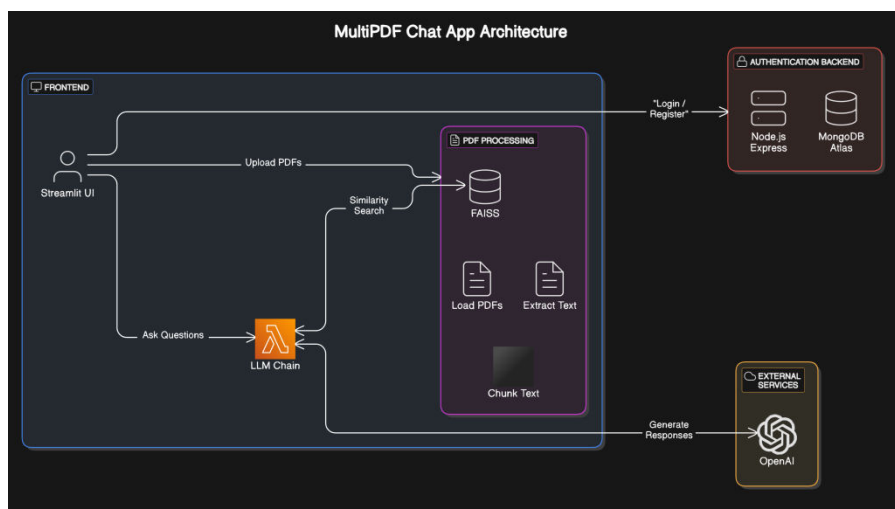


Fig 1.System Architecture

This diagram depicts the architecture of a MultiPDF Chat Application, showcasing a clear separation of concerns and a streamlined data flow. The user interacts with the system through a Streamlit UI on the FRONTEND, initiating actions such as uploading PDFs and asking questions. Upon uploading a PDF, the document is sent to the PDF PROCESSING module. Here, the PDF undergoes a series of transformations: it is loaded, text is extracted, and subsequently chunked into manageable segments. These text chunks are then indexed for Similarity Search using FAISS, enabling efficient retrieval of relevant information based on semantic similarity. User questions are processed through an LLM Chain, likely leveraging a large language model, to generate responses. This chain utilizes the indexed information from FAISS to provide contextually relevant answers. The system also features an AUTHENTICATION BACKEND built with Node.js and Express, storing user credentials in MongoDB Atlas, ensuring secure user access through "Login/Register" functionalities. Finally, the application integrates with EXTERNAL SERVICES, specifically OpenAI, to enhance its capabilities, likely for advanced natural language processing tasks involved in generating responses. This architecture emphasizes a modular design, facilitating scalability and maintainability, while providing a seamless user experience for interacting with multiple PDF documents.

V.CONCLUSION

In conclusion, the MultiPDF Chat App Architecture presented here offers a robust and scalable solution for interactive information retrieval from multiple PDF documents. By leveraging a combination of efficient PDF processing, semantic indexing with FAISS, and powerful language model integration through the LLM Chain, the system enables users to engage in natural language conversations with their documents. The incorporation of an

authentication backend ensures secure user access, while the integration of external services, such as OpenAI, enhances the system's capabilities for advanced natural language processing.

REFERENCES

- [1] Min-Yen Kan, "Automatic Text Summarization from PDFs using Machine Learning", Proceedings of the 2018 ACM Symposium on Document Engineering, pp. 215-224, 2018.
- [2] K. Xu, Y. Wu, and Y. Yan, "AI-powered Conversational Agents for Document Understanding", IEEE Transactions on Artificial Intelligence, vol. 2, no. 4, pp. 334-347, 2021.
- [3] R. Smith, "Hybrid OCR Techniques for Extracting Text from PDFs", IEEE Transactions on Pattern Analysis and Machine Intelligence, vol. 24, no. 6, pp. 789-803, 2002.
- [4] A. Qureshi and M. Ali, "Conversational AI and PDF Processing: A Deep Learning Approach", International Journal of Computational Intelligence, vol. 35, no. 3, pp. 567-582, 2019.
- [5] J. Brown and K. Patel, "Embedding-Based Text Retrieval from PDF Documents", IEEE Access, vol. 8, pp. 110542-110558, 2020.
- [6] P. Sharma and A. Gupta, "Chatbots for Document Analysis: A Transformer-Based Approach", Proceedings of the 2021 Conference on Natural Language Processing, pp. 1012-1025, 2021.
- [7] L. Wang, "Vector Databases for Efficient Retrieval in AI-driven PDF Chatbots", Journal of Data Science and AI, vol. 7, no. 2, pp. 135-152, 2022.
- [8] T. Nakamura, "Intelligent Question-Answering Systems for Research Papers: Challenges and Solutions", Proceedings of the 2020 IEEE International Conference on AI & NLP, pp. 223-235, 2020.



INTERNATIONAL
STANDARD
SERIAL
NUMBER
INDIA



INTERNATIONAL JOURNAL OF INNOVATIVE RESEARCH

IN SCIENCE | ENGINEERING | TECHNOLOGY

 9940 572 462  6381 907 438  ijirset@gmail.com



www.ijirset.com

Scan to save the contact details