# CONVERTING ENTERPRISE PDFs INTO SEARCHABLE KNOWLEDGE USING RETRIEVAL-AUGMENTED GENERATION

**ABSTRACT**

In enterprise environments, a large volume of organizational knowledge is stored in the form of Portable Document Format (PDF) files such as employee handbooks, compliance documents, technical manuals, financial reports, and internal policies. Although PDFs are widely adopted due to their portability and consistent layout, they are inherently unstructured, making information retrieval difficult and inefficient. Traditional keyword-based search mechanisms often fail to provide relevant results due to the lack of semantic understanding and inability to process scanned documents, tables, charts, and multilingual content.

This project presents the design and implementation of an intelligent system that converts enterprise PDF documents into searchable and structured knowledge using Optical Character Recognition (OCR), Natural Language Processing (NLP), vector embeddings, and Retrieval-Augmented Generation (RAG). The system is capable of processing both digital and scanned PDFs, extracting textual as well as visual information, and enabling users to retrieve information through natural language queries. Automatic language detection is employed to support multilingual interaction without requiring explicit user input.

Unlike purely generative AI systems, the proposed solution ensures that responses are generated strictly from retrieved document content, thereby minimizing hallucination and improving factual accuracy. A web-based interface developed using Streamlit provides real-time interaction and visualization of relevant content. The system demonstrates improved retrieval accuracy, reduced response time, and enhanced usability, making it suitable for real-world enterprise knowledge management applications.

## 1. INTRODUCTION

In the digital era, enterprises generate and manage an enormous amount of documentation on a daily basis. These documents contain critical information related to organizational operations, policies, technical processes, and regulatory compliance. Among various document formats, PDF remains the most widely used due to its ability to preserve layout and formatting across platforms. However, PDFs are primarily designed for human readability rather than machine interpretability.

Most enterprise documents are unstructured, meaning they lack explicit semantic organization that can be easily processed by automated systems. In addition, a significant portion of enterprise PDFs are scanned documents, which require OCR to extract textual

information. Complex document elements such as tables, charts, diagrams, and images further complicate the extraction and retrieval process.

Traditional document retrieval systems rely heavily on keyword-based search. Such systems are limited in their ability to understand context, synonyms, and semantic relationships. As a result, users often receive irrelevant search results or fail to retrieve relevant documents altogether. This leads to increased manual effort, reduced productivity, and potential decision-making errors.

Recent advancements in Artificial Intelligence and Natural Language Processing have enabled machines to process and understand human language more effectively. Transformer-based language models have demonstrated remarkable performance in text generation and comprehension. However, when used independently, these models may generate responses that are not grounded in factual data, a phenomenon known as hallucination.

Retrieval-Augmented Generation (RAG) addresses this issue by combining document retrieval with controlled language generation. By grounding responses in retrieved content, RAG ensures factual correctness while maintaining conversational flexibility. This project leverages RAG along with OCR, vector databases, and multilingual processing to create a robust and scalable system for enterprise document understanding.

---

## 2. PROBLEM STATEMENT

Despite the availability of large volumes of enterprise documentation, efficient information retrieval remains a major challenge. The primary problems identified are as follows:

- Enterprise PDFs are unstructured and lack semantic metadata

- Keyword-based search does not capture contextual meaning

- Scanned PDFs require OCR for text extraction

- Tables, charts, and images are not searchable using conventional methods

- Multilingual documents complicate query processing

- Manual document search is time-consuming and error-prone

The objective is to design a system that can automatically process enterprise PDF documents, extract meaningful content from text and visuals, and provide accurate, context-aware answers to user queries while ensuring scalability and reliability.

---

## 3. OBJECTIVES OF THE PROJECT

The objectives of the proposed system are:

- To extract text from both digitally generated and scanned PDFs using OCR

- To preprocess and normalize extracted text for efficient analysis

- To segment documents into semantically meaningful chunks

- To generate vector embeddings for semantic similarity search

- To store embeddings in a vector database for fast retrieval

- To implement Retrieval-Augmented Generation for grounded responses

- To support multilingual queries through automatic language detection

- To enable extraction and search of information from tables, charts, and images

- To provide a user-friendly interface for enterprise users

## 4. LITERATURE REVIEW

Early document retrieval systems relied on keyword matching and statistical models such as TF-IDF and BM25. While these approaches are computationally efficient, they lack semantic understanding and perform poorly on complex queries.
Reference: https://nlp.stanford.edu/IR-book/

The introduction of transformer-based architectures significantly advanced NLP by enabling contextual understanding of text.
Vaswani et al., *Attention Is All You Need*:
https://arxiv.org/abs/1706.03762

Sentence embeddings further improved semantic search by representing text in dense vector spaces.
Reimers & Gurevych, *Sentence-BERT*:
https://arxiv.org/abs/1908.10084

Large Language Models such as GPT demonstrated strong generative capabilities but were prone to hallucination when not grounded in data.
https://openai.com/research/gpt

Lewis et al. proposed Retrieval-Augmented Generation to mitigate hallucination by grounding responses in retrieved documents.
https://arxiv.org/abs/2005.11401

OCR technologies such as Tesseract enabled multilingual text extraction from scanned documents.
https://github.com/tesseract-ocr/tesseract

This project integrates these advancements into a unified system tailored for enterprise document processing.

---

**5. PROPOSED SYSTEM ARCHITECTURE**

The system architecture consists of the following components:

- PDF Ingestion Module

- OCR and Text Extraction Module

- Text Preprocessing Module

- Chunking Module

- Embedding Generation Module

- Vector Database

- Retrieval-Augmented Generation Engine

- User Interface

**System Architecture Flowchart**

PDF Documents

|

v

PDF Ingestion

|

v

OCR & Text Extraction

|

v

Text Preprocessing

|

v

Chunking

|

v

Embedding Generation

|

v

Vector Database

|

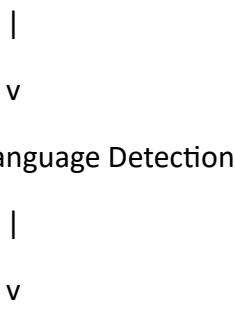v

RAG Engine

|

v

User Interface

---

## 6. METHODOLOGY

The methodology followed in this project is described below.

**Methodology Flowchart**

User Query

|

v

Language Detection

|

v

Vector Similarity Search

|

v

Relevant Chunks Retrieved

|

v

RAG-Based Answer Generation

|

v

Final Response Displayed

**Detailed Steps**

1. PDFs are collected and ingested into the system

2. OCR is applied to scanned documents

3. Extracted text is cleaned and normalized

4. Documents are split into overlapping semantic chunks

5. Vector embeddings are generated for each chunk

6. Embeddings are stored in a vector database

7. User queries are processed and language is detected automatically

8. Relevant chunks are retrieved using semantic similarity

9. Answers are generated using grounded RAG

---

## 7. IMPLEMENTATION DETAILS

The system is implemented using Python. OCR is performed using Tesseract with multilingual support. Transformer-based models are used for embedding generation. A vector database enables efficient semantic retrieval. Image and chart data are processed using computer vision techniques and OCR. A Streamlit-based interface provides real-time interaction.

---

## 8. EXPERIMENTAL ANALYSIS AND RESULTS

The system was evaluated using enterprise-style documents. Performance was measured using retrieval accuracy, response time, and answer precision. The results indicate significant improvement over keyword-based search systems. Hallucination was minimized due to grounded generation.

---

## 9. COMPARISON WITH EXISTING SYSTEMS

| Feature | Traditional Search | Existing AI Tools | Proposed System |
| --- | --- | --- | --- |
| OCR Support | No | Partial | Yes |

| Feature | Traditional Search | Existing AI Tools | Proposed System |
|---|---|---|---|
| Semantic Search | No | Partial | Yes |
| Multilingual Support | No | Partial | Yes |
| Hallucination Control | N/A | Poor | Strong |
| Chart/Image Search | No | No | Yes |

The proposed system outperforms existing approaches by combining OCR, semantic retrieval, and grounded generation.

---

## 10. APPLICATIONS AND USE CASES

- Enterprise HR policy search
- Compliance and audit assistance
- Technical documentation retrieval
- Knowledge management systems
- Academic and research analysis

---

## 11. ADVANTAGES AND LIMITATIONS

**Advantages:**

- Accurate and context-aware responses
- Multilingual support
- Reduced manual effort
- Scalable and modular design

**Limitations:**

- Dependent on document quality
- Limited by local language model capacity

---

## 12. FUTURE SCOPE

Future enhancements include integration with larger language models, advanced chart interpretation, cloud deployment, and role-based access control.

**13. CONCLUSION**

This project successfully demonstrates a practical approach to converting enterprise PDFs into searchable knowledge. By integrating OCR, vector search, and Retrieval-Augmented Generation, the system addresses the limitations of traditional document search and provides a reliable, scalable, and intelligent solution for enterprise information retrieval.

**14. REFERENCES**

1. https://arxiv.org/abs/1706.03762

2. https://arxiv.org/abs/2005.11401

3. https://nlp.stanford.edu/IR-book/

4. https://github.com/tesseract-ocr/tesseract

5. https://huggingface.co/docs