

# PROJECT REPORT

**Name : Pariksith G**

**Subject : Generative – AI**

**Class : AI & DS (B section)**

**Reg No : 710723243077**

**Roll No : 23AD077**

**Submission Date : 19.01.2026**

**GitHub Link: <https://github.com/pariksith/Document-Intelligence-System>**

## **DOCUMENT INTELLIGENCE SYSTEM**

### **1. INTRODUCTION**

In today's digital era, vast amounts of information are stored in the form of documents such as PDFs, reports, research papers, notes, and manuals across education, business, and online platforms. Extracting meaningful information from these documents and finding precise answers often requires manual reading, searching, and analysis, which can be time-consuming and inefficient. Students, researchers, and professionals frequently face difficulties when working with large or complex documents, leading to reduced productivity and delayed decision-making.

The Document Intelligence System is an AI-powered solution designed to automatically understand, analyze, and retrieve information from uploaded documents. The system enables users to upload PDF files and ask natural language questions, instantly receiving accurate and context-aware answers. By leveraging modern techniques such as Natural Language Processing (NLP), vector embeddings, and retrieval-augmented generation (RAG), the system identifies relevant document sections and generates human-like responses based on the content.

This project demonstrates the practical integration of Generative AI with web technologies to build an intelligent and user-friendly document analysis application. The system minimizes manual effort, enhances information accessibility, and accelerates knowledge discovery. It also highlights the real-world application of AI in document understanding, showcasing how intelligent automation can significantly improve efficiency, learning, and digital productivity across various domains..

### **2. PROBLEM STATEMENT**

Traditional document analysis processes require users to:

- Manually read and scan large volumes of documents
- Search for relevant information across multiple pages
- Interpret complex content to extract meaningful insights

This creates challenges such as:

- Difficulty in quickly locating precise answers
- Increased time and effort when working with lengthy PDFs or notes
- Reduced productivity for students, researchers, and professionals

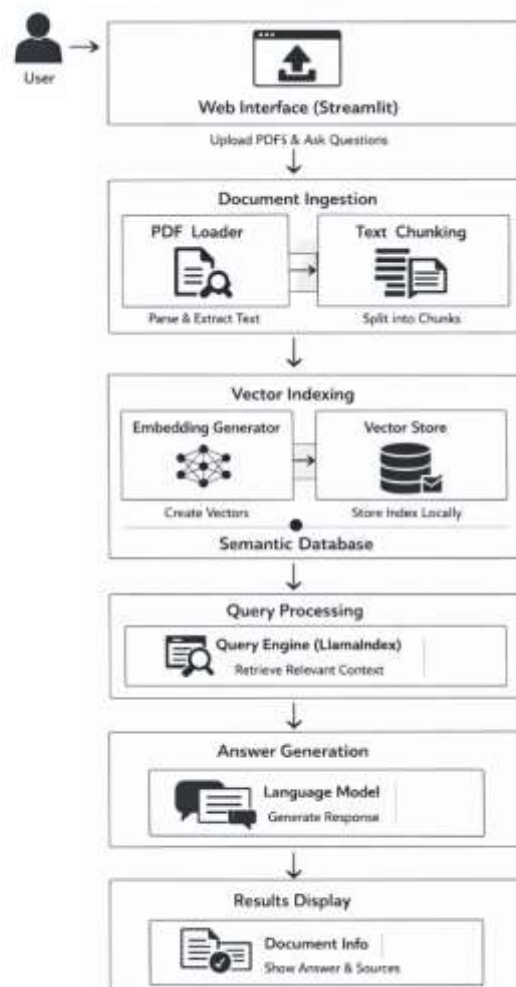
### 3. OBJECTIVES OF THE PROJECT

The main objectives of this project are:

- To design a simple and user-friendly web interface for uploading and managing documents such as PDFs and notes.
- To automatically understand and analyze document content using Artificial Intelligence techniques.
- To enable users to ask natural language questions and retrieve accurate, context-aware answers from documents.

### 4. SYSTEM ARCHITECTURE

Document Intelligence System Architecture



## **5. MODULE**

### **Module 1: User Interface Module (Document Interaction Interface)**

The User Interface module provides a clean, minimal, and interactive web interface through which users can interact with the Document Intelligence System. It allows users to upload PDF documents and enter natural language queries to retrieve relevant information from the uploaded content. The interface is designed to be intuitive and visually simple, enabling both technical and non-technical users to operate the system without prior training. This module plays a crucial role in user engagement by presenting upload options, query input fields, and answer displays in a structured manner. It also provides visual feedback such as indexing status, loading indicators, and result cards, helping users understand system operations clearly. Overall, this module ensures usability, accessibility, and a smooth interaction experience.

### **Module 2: Document Processing and Indexing Module**

This module is responsible for processing uploaded PDF documents and converting them into a structured format suitable for intelligent querying. It extracts textual content from PDF files, divides the content into smaller meaningful chunks, and generates vector embeddings for each chunk. These embeddings are stored in a local vector database, enabling efficient semantic search and retrieval. The module ensures that documents are indexed correctly and persistently stored for future queries. By organizing unstructured document data into searchable vector representations, this module forms the foundation of the document intelligence functionality and ensures accurate and context-aware information retrieval.

### **Module 3: Query Processing and Answer Generation Module**

The Query Processing module handles user questions and retrieves relevant information from the indexed documents. It converts user queries into vector form and performs semantic similarity search against the stored document embeddings to identify the most relevant sections. These retrieved contexts are then passed to a language model, which generates clear and meaningful answers based strictly on the document content. The module ensures that responses are accurate, concise, and contextually relevant. It also supports configurable result limits to control the number of retrieved snippets. This module plays a key role in transforming raw document data into useful insights, enabling intelligent question-answering over uploaded documents.

## **6. USER INTERFACE**

The user interface of the Document Intelligence System is developed using a simple and interactive web framework, providing an easy way for users to interact with the system. The interface offers the following features:

- Document upload option for users to upload PDF files
- Input field to enter natural language questions
- Submit button to initiate document analysis and query processing
- Display section for generated answers with relevant document

The UI is designed to be minimal, clean, and intuitive so that even non-technical users can use the system efficiently without prior knowledge..

## 7. TOOLS AND TECHNOLOGIES USED

### **Python:**

Python is used as the primary programming language for developing the core logic of the Document Intelligence System. It offers extensive libraries for natural language processing, document handling, and machine learning integration. Python enables efficient implementation of AI-driven features such as text extraction, embedding generation, and intelligent question answering.

### **Streamlit:**

Streamlit is used to build the web-based user interface. It allows rapid development of interactive applications and simplifies user interaction with document uploads and query inputs.

### **LlamaIndex:**

LlamaIndex is used for document indexing, chunking, and semantic retrieval. It helps transform unstructured document content into searchable vector representations.

### **HuggingFace Embeddings:**

HuggingFace models are used to generate embeddings that capture the semantic meaning of document text, enabling accurate similarity-based retrieval.

### **Vector Store (Local Storage):**

A local vector storage mechanism is used to store document embeddings, allowing fast and efficient retrieval without relying on external paid services.

## 8. OUTPUT SCREENSHOTS

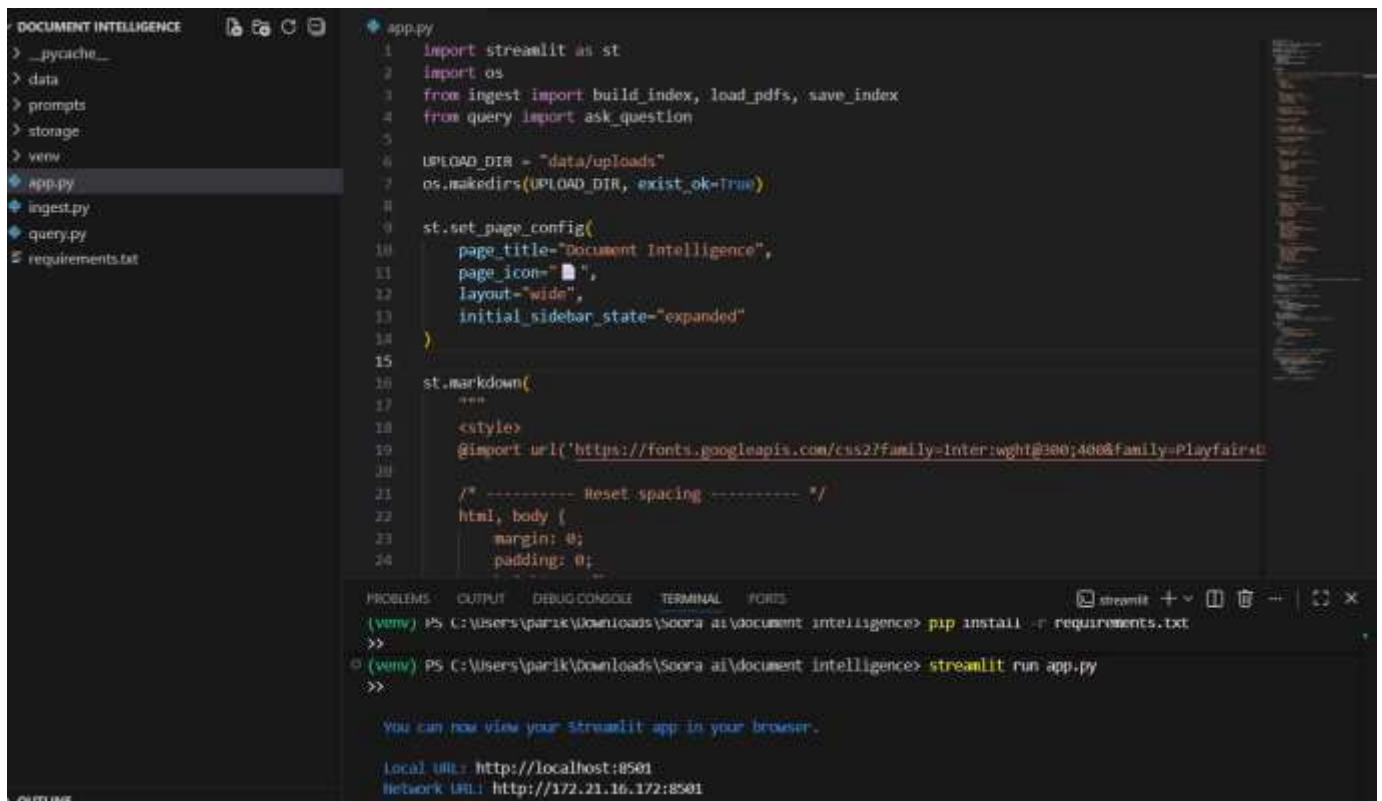


## Q1. Outline the Core Concepts of Data Mining and Predictive Analytics with Real-World Applications

1. Introduction In the era of big data and digital transformation, organizations generate enormous volumes of data from transactions, sensors, social media, IoT devices, and online platforms. Merely storing data is not sufficient; extracting meaningful insights from data has become essential for informed decision-making. Data Mining and Predictive Analytics play a crucial role in converting raw data into valuable knowledge. Data mining focuses on discovering hidden patterns, relationships, and trends in large datasets, while predictive analytics uses these patterns along with statistical and machine learning techniques to predict future outcomes. Together, they form the backbone of modern intelligent systems used in business, healthcare, finance, manufacturing, and government sectors.
2. Data Mining – Definition and Overview Definition Data Mining is the process of exploring large datasets to identify meaningful patterns, correlations, anomalies, and knowledge using statistical, machine learning, and database techniques. It is a key step in the Knowledge Discovery in Databases (KDD) process. Objectives of Data Mining
  - Discover hidden patterns and trends
  - Improve decision-making
  - Predict future behaviors
  - Detect fraud and anomalies
  - Support strategic planning
3. Knowledge Discovery in Databases (KDD) Process Data mining is part of the broader KDD process, which consists of multiple stages. Stages of KDD
4. Data Collection (File: SELF STUDY ANS.pdf, Page: 1)

Personalized learning

11. Advantages of Data Mining and Predictive Analytics
  - Improved decision-making
  - Cost reduction
  - Risk mitigation
  - Increased efficiency



The screenshot shows a Streamlit web application titled "Document Intelligence". The interface includes a sidebar with a file explorer showing folders like \_\_pycache\_\_, data, prompts, storage, venv, and files, and a list of files: app.py, ingest.py, query.py, and requirements.txt. The main area displays a code editor with the following Python code:

```
1 import streamlit as st
2 import os
3 from ingest import build_index, load_pdfs, save_index
4 from query import ask_question
5
6 UPLOAD_DIR = "data/uploads"
7 os.makedirs(UPLOAD_DIR, exist_ok=True)
8
9 st.set_page_config(
10     page_title="Document Intelligence",
11     page_icon="📄",
12     layout="wide",
13     initial_sidebar_state="expanded"
14 )
15
16 st.markdown(
17     """
18     <style>
19     @import url('https://fonts.googleapis.com/css2?family=Inter:wght@300;400&family=Playfair+T
20
21     /* ----- Reset spacing ----- */
22     html, body {
23     margin: 0;
24     padding: 0;
```

Below the code editor, the terminal shows the command to run the application:

```
(venv) PS C:\Users\parik\Downloads\Soora AI\document intelligence> pip install -r requirements.txt
>>
(venv) PS C:\Users\parik\Downloads\Soora AI\document intelligence> streamlit run app.py
>>
```

The terminal output indicates that the application is running and provides the local and network URLs:

```
You can now view your Streamlit app in your browser.

Local URL: http://localhost:8501
Network URL: http://172.21.16.172:8501
```

## **9. USE CASES**

- Automated question answering from academic notes, textbooks, and research papers
- Intelligent document analysis for students and educators
- Information retrieval from large PDF documents in corporate environments
- Smart document-based chat systems for knowledge management

## **10. ADVANTAGES**

- Eliminates the need for manual document searching and reading
- Saves time by providing instant, context-aware answers
- Simple and user-friendly interface for non-technical users
- Demonstrates practical implementation of Document Intelligence and AI
- Handles large volumes of documents efficiently

## **11. LIMITATIONS**

- Performance may decrease when processing very large or complex documents
- Answer quality depends on the accuracy of document text extraction and embeddings
- Limited to supported document formats such as PDF

## **12. FUTURE ENHANCEMENTS**

- Integration of more advanced language models for improved answer accuracy
- Addition of document summarization and keyword extraction features
- Deployment as a mobile or cloud-based application for wider accessibility
- Multi-language document understanding and question answering

### **13. CONCLUSION**

The Document Intelligence System demonstrates the effective application of artificial intelligence techniques to automate document understanding and information retrieval. By enabling users to upload documents and interact with them through intelligent querying, the system significantly reduces the manual effort required to search, read, and analyze large volumes of text. The integration of natural language processing and AI-based retrieval mechanisms allows the system to provide accurate, context-aware responses, improving efficiency and decision-making.