# Project Report:

# Conversio

"Transforming PDFs into Interactive Knowledge Bases"

By:

El Kajam Hamza

El Amrani Ranya

El Rhirhayi Taha

2024-2025

1. **Introduction**

The **Chat with PDF** application is an AI-powered tool designed to revolutionize how users interact with PDF documents. By leveraging advanced natural language processing (NLP) and vector database technologies, the app allows users to upload PDFs, ask questions in natural language, and receive precise, context-aware answers instantly. This project combines cutting-edge technologies like **Google Gemini**, **Pinecone**, **FAISS**, and **Streamlit** to create a seamless and intuitive user experience.

## 2. Problem Statement

### The Challenge

- **Inefficient Document Interaction**: Extracting specific information from large PDF documents is time-consuming and often requires manual searching.

- **Limited Search Capabilities**: Traditional search methods (e.g., Ctrl+F) are ineffective for complex queries or context-based questions.

- **Lack of Scalability**: Existing solutions often struggle with large datasets or require significant computational resources.

### The Need

There is a growing demand for tools that can:

1. Automate the process of extracting insights from documents.

2. Provide accurate, context-aware answers to user queries.

3. Scale efficiently for both small and large datasets.

## 3. Solution Overview

### What It Does

The **Chat with PDF** application addresses these challenges by:

1. Allowing users to upload PDF documents.

2. Extracting and processing text from the documents.

3. Using AI to generate embeddings and store them in a vector database.

4. Enabling users to ask questions and receive instant, accurate answers.

### Key Features

- **PDF Upload and Processing**: Supports multiple PDFs and extracts text efficiently.

- **AI-Powered Question Answering**: Uses Google Gemini for embeddings and context-aware answers.

- **Hybrid Vector Storage**: Combines Pinecone (cloud) and FAISS (local) for flexibility and reliability.

- **User-Friendly Interface**: Built with Streamlit for a clean, interactive experience.

## 4. Architecture

The application is built on a modular architecture with the following components:

### Frontend

- **Streamlit**: Provides a simple, interactive user interface for uploading PDFs and asking questions.

- **Features**:
  - File uploader for PDFs.
  - Chat-like interface for questions and answers.
  - System status display (e.g., Pinecone connection).

### Backend

- **Python**: Handles the core logic, including text extraction, embeddings, and question answering.

- **Libraries Used**:
  - **PyPDF2**: Extracts text from PDFs.
  - **LangChain**: Manages text splitting, embeddings, and conversational chains.
  - **GoogleGenerativeAI**: Generates embeddings and powers the question-answering system.

### Large Language Model (LLM)

- **Google Gemini**: Used for:
  - Generating embeddings from text chunks.
  - Answering user questions based on the document context.

- **Advantages**:
  - State-of-the-art performance.
  - Scalable and reliable for production use.

## Database

- **Pinecone**: A cloud-based vector database for storing and querying embeddings.
  - Enables fast similarity search for large datasets.
  - Scalable and managed service.
- **FAISS**: A local vector store used as a fallback when Pinecone is unavailable.
  - Lightweight and efficient for small to medium datasets.
  - Ensures reliability and offline functionality.

## 5. Implementation

## Workflow

1. **PDF Upload**: Users upload PDFs through the Streamlit interface.
2. **Text Extraction**: The app extracts text using PyPDF2.
3. **Text Splitting**: The text is split into smaller chunks using LangChain's RecursiveCharacterTextSplitter.
4. **Embeddings**: Google Gemini generates embeddings for each text chunk.
5. **Vector Storage**: Embeddings are stored in Pinecone (or FAISS as a fallback).
6. **Question Answering**: When a user asks a question:
   - The app retrieves relevant chunks from the vector store.
   - Google Gemini generates an answer based on the retrieved context.
7. **Display Answer**: The answer is displayed in the chat interface.

## Key Functions

- **get_pdf_text**: Extracts text from uploaded PDFs.
- **get_text_chunks**: Splits text into smaller chunks for processing.
- **get_hybrid_vector_store**: Stores embeddings in Pinecone and FAISS.

- **get_conversational_chain**: Generates answers using Google Gemini.
- **user_input**: Handles user questions and retrieves answers.

## 6. Results

**Outcomes**

- **Efficient Document Interaction**: Users can quickly find information in PDFs without manual searching.
- **Accurate Answers**: The app provides context-aware answers using Google Gemini.
- **Scalable Solution**: The hybrid vector storage (Pinecone + FAISS) ensures reliability and scalability.

**User Experience**

- The Streamlit interface is intuitive and easy to use.
- Real-time feedback (e.g., processing time, vector store used) enhances transparency.

## 7. Challenges and Solutions

**Challenges**

1. **Dependency Conflicts**: Issues with Pinecone and FAISS versions.
   - **Solution**: Pinned compatible versions and added fallback mechanisms.
2. **Embedding Storage**: Embeddings were not being stored correctly in Pinecone.
   - **Solution**: Debugged and fixed Pinecone initialization and storage logic.
3. **Handling Large PDFs**: Processing large documents was slow.
   - **Solution**: Optimized text splitting and chunking.

**Lessons Learned**

- Proper version management is crucial for avoiding dependency conflicts.
- Hybrid storage solutions (Pinecone + FAISS) provide flexibility and reliability.
- User feedback is invaluable for improving the app's functionality.

**8. Future Enhancements**

**Planned Features**

1. **Support for More File Formats**: Add support for Word, Excel, and PowerPoint files.

2. **Multi-Language Support**: Enable the app to work with documents in multiple languages.

3. **Cloud Integration**: Integrate with Google Drive, Dropbox, and OneDrive for seamless document access.

4. **Advanced Query Handling**: Add support for multi-part questions and contextual follow-ups.

5. **User Authentication**: Implement user accounts and document management features.

6. **Analytics and Insights**: Provide usage statistics and document insights.

**Scalability**

- Deploy the app on cloud platforms (e.g., AWS, GCP) for scalability.

- Optimize the backend for handling larger datasets and more concurrent users.

**9. Conclusion**

The **Chat with PDF** application demonstrates the power of combining AI, vector databases, and modern web frameworks to solve real-world problems. By automating the process of extracting insights from PDFs, the app saves time and improves productivity for users. With its modular architecture and scalable design, the app is well-positioned for future enhancements and widespread adoption.

[GitHub Repo of the project.](#)

**10. References**

- **Streamlit:** https://streamlit.io/

- **Google Gemini**: https://ai.google/

- **Pinecone**: https://www.pinecone.io/

- **LangChain**: https://www.langchain.com/

- **FAISS**: https://github.com/facebookresearch/faiss