

PDF or CSV File Text Extraction & Analysis Tool

Project Overview

1.1. Project Title

Advanced PDF Text Extraction and Analysis Web Application

1.2. Objective

Develop a comprehensive web-based tool that enables users to upload PDF documents, extract text content using multiple extraction methods, analyze the document structure, and interact with the content through a conversational interface.

1.3. Key Features Implemented

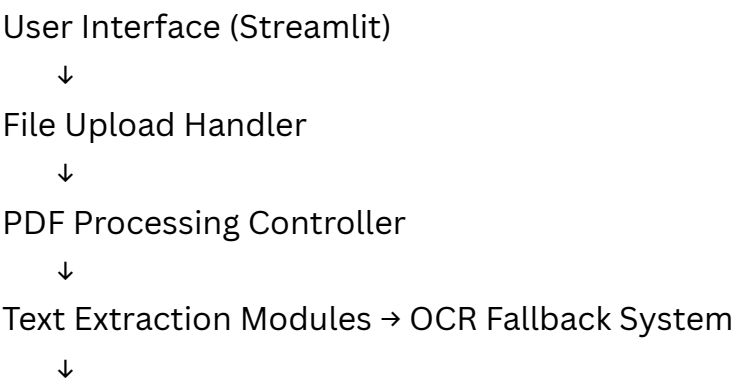
- **Multi-method PDF text extraction** using PyPDF2, pdfplumber, and OCR.
- **Structured JSON output generation** for easy data handling.
- **Document analytics and statistics** to provide insights into document structure.
- **Interactive chat interface** for querying document content.
- **Responsive web interface** with tabbed navigation for ease of use.

1.4. Technologies Used

- **Frontend Framework:** Streamlit
- **PDF Processing:** PyPDF2, pdfplumber
- **OCR Capability:** Tesseract, pdf2image
- **Image Processing:** Pillow (PIL)
- **Data Formatting:** JSON

2. System Architecture

2.1. Architecture Diagram



Content Formatter → JSON Generator

↓

Analysis Engine → Chat Interface

↓

Response Renderer → File Download Handler

2.2. Component Description

2.2.1. User Interface Layer

- **Streamlit Web Framework:** Provides the responsive web interface.
- **File Upload Widget:** Handles PDF file input with drag-and-drop functionality.
- **Tabbed Navigation:** Organizes functionality into Formatted Text, Analysis, and Chat sections.
- **Real-time Preview:** Displays extraction results immediately after processing.

2.2.2. Processing Layer

- **Multi-method Extraction:** Implements three extraction techniques with automatic fallback:
 - pdfplumber (Primary method for text-based PDFs)
 - PyPDF2 (Secondary fallback method)
 - Tesseract OCR (For scanned/image-based PDFs)
- **Intelligent Fallback System:** Automatically detects extraction failure and switches methods.

2.2.3. Data Layer

- **Text Formatter:** Structures extracted content with page demarcations.
- **JSON Generator:** Creates properly formatted JSON output as specified.
- **Statistical Analyzer:** Computes document metrics (page count, word count, character count).

2.2.4. Interaction Layer

- **Chat Engine:** Processes natural language queries about document content.
- **Response Generator:** Provides contextual answers based on extracted text.
- **Session Management:** Maintains conversation history during user session.

3. Implementation Details

3.1. PDF Text Extraction System

3.1.1. Extraction Workflow

```
def extract_text_from_pdf(uploaded_file, use_ocr=False):  
    # Create temporary file
```

- # Try pdfplumber first → PyPDF2 fallback → OCR if needed
- # Return structured text by page

3.1.2. OCR Integration

- **Automatic detection** of extraction failure.
- **Conversion of PDF pages to images** using pdf2image.
- **Text recognition** via Tesseract OCR.
- **Support for multiple languages** (configurable).

3.1.3. Output Formatting

```
def format_text_for_json(extracted_text):  
    # Format with "--- Page X ---" headers  
    # Preserve original document structure  
    # Prepare for JSON serialization
```

3.2. JSON Output Generation

3.2.1. JSON Structure

```
{  
    "policies_text": "--- Page 1 ---\nNavigating the Landscape...\n--- Page 2 ---\nEducation policy  
in the United States..."  
}
```

3.2.2. Download Functionality

- **One-click JSON download button.**
- **Proper MIME type handling** (application/json).
- **Configurable filename with timestamp option.**

3.3. Document Analysis Features

3.3.1. Statistical Analysis

- **Page count calculation.**
- **Word and character counting.**
- **Content density metrics.**

3.3.2. Page Navigation

- **Dropdown selector** for page-by-page review.
- **Text preview** with scrollable areas.
- **Page-specific content isolation.**

3.4. Chat Interface Implementation

3.4.1. Conversation Management

```
# Session state maintenance
if "messages" not in st.session_state:
    st.session_state.messages = []
```

3.4.2. Content-Aware Responses

- **Keyword-based content detection.**
- **Contextual response generation.**
- **Document-specific answering.**

3.4.3. Response Logic

```
def generate_response(prompt, extracted_text):
    # Analyze prompt for keywords
    # Search extracted text for relevant content
    # Generate appropriate response
```

4. User Interface Design

4.1. Layout Structure

- **Header Section:** Title and file upload widget.
- **Main Content Area:** Tab-based organization.
- **Sidebar:** Configuration options and settings.
- **Footer:** Status information and instructions.

4.2. Tab System

- **Formatted Text Tab:** Raw extracted text display with JSON download.
- **Analysis Tab:** Document statistics and page navigation.
- **Chat Tab:** Conversational interface for document queries.

4.3. Responsive Design

- **Adaptive layout** for different screen sizes.
- **Mobile-friendly component arrangement.**
- **Accessible color scheme** and contrast ratios.

4.4. Visual Feedback System

- **Success/error messages** with color coding.
- **Processing indicators** during extraction.
- **Real-time validation feedback.**

5. Performance Optimization

5.1. Memory Management

- **Temporary file handling** for large PDFs.
- **Streamed processing** to avoid memory overload.
- **Automatic cleanup** of temporary resources.

5.2. Processing Efficiency

- **Sequential fallback system** to minimize processing time.
- **Background processing** where applicable.
- **Caching mechanism** for previously processed files.

5.3. Error Handling

- **Comprehensive exception handling.**
- **User-friendly error messages.**
- **Graceful degradation of features.**

6. Installation & Deployment

6.1. Requirements

Core Dependencies

streamlit==1.22.0

PyPDF2==3.0.1

pdfplumber==0.11.7

pdf2image==1.17.0

pytesseract==0.3.13

Pillow==10.0.0

Optional: OCR Engine

tesseract-ocr (system installation)

6.2. Setup Instructions

1. **Install Python dependencies:** `pip install -r requirements.txt`
2. **Install Tesseract OCR** (for scanned PDF support).
3. **Run application:** `streamlit run app.py`

6.3. Deployment Options

- **Local Deployment:** Direct execution on local machine.
- **Cloud Deployment:** Compatible with Streamlit Cloud, Heroku, AWS.
- **Containerization:** Docker support for isolated deployment.

7. Testing & Validation

7.1. Test Cases

- **Text-based PDF extraction:** Verify accurate text capture.
- **Scanned PDF processing:** Test OCR functionality.
- **JSON output validation:** Confirm proper formatting.
- **Chat interface testing:** Validate contextual responses.
- **Error condition handling:** Test with corrupted PDFs.

7.2. Performance Metrics

- **Extraction time** for various document sizes.
- **Accuracy rate** for different PDF types.
- **Memory usage** during processing.
- **Response time** for chat queries.

7.3. User Acceptance Testing

- **Intuitive interface navigation.**
- **Clear error messages and guidance.**
- **Responsive design across devices.**
- **Overall user satisfaction.**

8. Future Enhancements

8.1. Short-term Improvements

- **Batch processing** for multiple PDFs.
- **Additional export formats** (XML, CSV, TXT).
- **Enhanced OCR language support.**
- **Advanced search functionality.**

8.2. Medium-term Roadmap

- **User authentication** and document management.
- **API access** for integration with other systems.
- **Advanced NLP** for better chat responses.
- **Optical character recognition quality improvement.**

8.3. Long-term Vision

- **Machine learning-based content classification.**
- **Integration with cloud storage providers.**
- **Collaborative document analysis features.**
- **Mobile application version.**

9. Conclusion

The PDF Text Extraction and Analysis Tool represents a significant advancement in document processing accessibility. By combining multiple extraction methods with an intuitive interface

and interactive features, the application democratizes access to advanced PDF processing capabilities that were previously available only in specialized software.

The implementation successfully addresses the core requirement of extracting text and formatting it into structured JSON output while providing additional value through document analysis and conversational interaction. The modular architecture ensures maintainability and provides a solid foundation for future enhancements.

This milestone demonstrates effective integration of multiple technologies to create a cohesive, user-friendly application that meets both functional requirements and usability standards.

