

**DESIGN AND IMPLEMENTATION OF A PDF TEXT EXTRACTION SYSTEM WITH  
PYTHON**

**BY**

**BOLAJI MATTHEW**

**MATRIC NO.: 2021006064**

A PROJECT REPORT SUBMITTED TO THE DEPARTMENT OF COMPUTER SCIENCE,  
LADOKE AKINTOLA UNIVERSITY OF TECHNOLOGY, OGBOMOSHO, IN PARTIAL  
FULFILLMENT OF THE REQUIREMENTS FOR THE AWARD OF THE DEGREE OF  
BACHELOR OF TECHNOLOGY (B.TECH) IN COMPUTER SCIENCE

**SUPERVISOR: PROF. W. O. ISMAILA**

**AUGUST 2025**

## **CERTIFICATION**

This is to certify that the project titled **DESIGN AND IMPLEMENTATION OF A PDF TEXT EXTRACTION SYSTEM WITH PYTHON**, submitted by **Bolaji Matthew (Matric No: 2021006064)**, was carried out under my supervision in the **Department of Computer Science, Ladoke Akintola University of Technology (LAUTECH), Ogbomosho, Oyo State, Nigeria**. This project is hereby approved as meeting part of the requirements for the award of the Bachelor of Technology (B.Tech.) degree in Computer Science.

---

**Supervisor**

---

**Date**

Prof. W. O. Ismaila, B.Tech., M.Tech., Ph.D.

Professor of Computer Science

Department of Computer Science

Ladoke Akintola University of Technology, Ogbomosho, Nigeria

---

**Head of Department**

---

**Date**

Prof. W. O. Ismaila, B.Tech., M.Tech., Ph.D.

Professor of Computer Science

Department of Computer Science

Ladoke Akintola University of Technology, Ogbomosho, Nigeria

## DEDICATION

This project is wholeheartedly dedicated to **God Almighty**, whose grace and mercy sustained me throughout the course of this work. I also dedicate it to my beloved parents, **Mr. and Mrs. Bolaji**, for their unwavering support, encouragement, and unconditional love, which have been a constant source of strength and motivation.

## ACKNOWLEDGEMENT

I express my profound gratitude to **God Almighty**, the giver of knowledge and wisdom, for granting me the strength, guidance, and perseverance to successfully complete this project.

My deepest appreciation goes to my supervisor, **Prof. W. O. Ismaila (B.Tech., M.Tech., Ph.D.)**, for his invaluable guidance, constructive feedback, and continuous encouragement throughout the duration of this research. His mentorship greatly enriched the quality of this work.

I remain indebted to my parents, **Mr. and Mrs. Bolaji**, for their financial, moral, and emotional support, as well as for providing me with the invaluable gift of education. I am eternally grateful for their sacrifices and love.

Finally, I acknowledge my colleagues, friends, and everyone who contributed directly or indirectly to the success of this project and my academic journey at LAUTECH. Your support and encouragement will always be remembered with gratitude.

## **ABSTRACT**

The proliferation of digital documents in Portable Document Format (PDF) poses significant challenges for automated data processing. Manual text extraction from diverse PDF formats, including image-based scans and complex layouts, is inefficient, labor-intensive, and error-prone. This project designs, implements, and evaluates a robust PDF Text Extraction System using Python.

The system adopts a modular pipeline architecture that distinguishes between text-based and image-based PDFs, applying appropriate extraction methods. For text-based documents, direct parsing is performed using the pdfplumber library. For scanned documents, an Optical Character Recognition (OCR) pipeline integrates pdf2image for page conversion, OpenCV for image preprocessing (employing algorithms such as Gaussian Blur for noise reduction and adaptive thresholding for binarization), and Tesseract (leveraging LSTM-based neural networks for character recognition) for text extraction.

Encapsulated in a secure, multi-user web application using Streamlit, the system features user authentication, batch processing, and export options in formats including .txt, .csv, .docx, .xlsx, .json, and .mp3. Evaluation on a diverse PDF dataset yields an F1-Score of 98.5% for text-based documents and 95% for high-quality image-based ones, with additional metrics like Precision, Recall, and Character Error Rate (CER) confirming high accuracy and efficiency. This scalable tool enhances text extraction for data mining and natural language processing applications.

## TABLE OF CONTENTS

CONTENTS	PAGES
<b>CERTIFICATION</b>	<b>I</b>
<b>DEDICATION</b>	<b>II</b>
<b>ACKNOWLEDGEMENT</b>	<b>III</b>
<b>ABSTRACT</b>	<b>IV</b>
<b>CHAPTER 1: INTRODUCTION</b>	<b>1</b>
1.1 BACKGROUND OF THE STUDY	1
1.2 PROBLEM STATEMENT	1
1.3 AIM AND OBJECTIVES OF THE STUDY	2
1.3.1 Aim of the Study	2
1.3.2 Objectives of the Study	2
1.4 SCOPE OF THE STUDY	3
1.5 SIGNIFICANCE	3
<b>CHAPTER 2: LITERATURE REVIEW</b>	<b>5</b>
2.1 INTRODUCTION	5
2.2 PDF FILE FORMAT	5
2.2.1 The Structure and Components of PDF Files	5
2.2.2 Document Structure	7
2.2.3 The Challenges of Parsing PDF Files due to their Complexity and Potential for Variations.	8
2.3 TEXT EXTRACTION TECHNIQUES	10
2.3.1 Optical Character Recognition (OCR)	10
2.3.2 The process of Image Preprocessing	13
2.3.3 PDF Parsing Libraries:	17
2.3.4 Natural Language Processing (NLP) Techniques:	22
2.3.5 The potential for improving accuracy and extracting meaningful information.	26

2.4 RELATED WORK	29
2.4.1 Review of Existing Projects and Research Related to PDF Text Extraction Using Python	30
2.4.2 Identifying Gaps in Existing Research and Projects	34
2.4.3 How This Project Will Contribute to the Field	35
<b>CHAPTER 3: SYSTEM DESIGN AND METHODOLOGY</b>	<b>36</b>
3.1 INTRODUCTION	36
3.1.1 Modular Text Extraction Pipeline Methodology for PDF Document Processing	36
3.1.2 Overview of Methodology Steps	38
3.2 SYSTEM REQUIREMENT SPECIFICATION	39
3.2.1 System Requirement for the PDF Text Extraction System	40
3.3 SYSTEM ARCHITECTURE	42
3.3.1 System Workflow (Textual Diagram)	43
3.3.2 Core Components	45
3.3.3 Technologies Used	46
3.4 DATA FLOW DIAGRAM (DFD)	47
3.4.1 Context Diagram (Level 0)	49
3.4.2 Structured Diagram (Level 1)	51
3.5 UML DIAGRAMS	52
3.5.1 Class Diagram	52
3.5.2 Use Case	54
3.5.3 Sequence Diagram	55
3.5.4 Data Dictionary	56
3.6 FLOWCHART	58
3.7 DATABASE MODEL	60
3.7.1 ER-Model	60
3.7.2 Relational Model	64
3.7.3 Sample Queries	64
3.8 KEY COMPONENTS OF THE SYSTEM	66
3.8.1 Data Preprocessing	66

3.8.2 Text Extraction:	70
3.8.3 Post-Processing	74
3.8.4 Evaluation Metrics	76
3.9 DEVELOPMENT ENVIRONMENT	80
3.9.1 Programming Language, Libraries, and Development Tools	80
3.9.2 Hardware and Software Requirements	82
<b>CHAPTER 4: IMPLEMENTATION</b>	<b>84</b>
4.1 DEVELOPMENT ENVIRONMENT SETUP	84
4.2 SYSTEM IMPLEMENTATION: THE MODULAR PIPELINE	85
4.2.1 Module 1: Input Handling and PDF Type Detection	85
4.2.2 Module 2: Preprocessing for Image-Based PDFs	85
4.2.3 Module 3: Text Extraction	85
4.2.4 Module 4: Post-Processing and Text Refinement	85
4.2.5 Module 5: Evaluation and Output	86
4.3 USER INTERFACE (UI) IMPLEMENTATION	86
4.3.1 Login and Registration Pages	86
4.3.2 Main Application Dashboard and File Uploader	89
4.3.3 Processing and Downloadable Results	91
4.4 CHAPTER SUMMARY	94
<b>CHAPTER 5: SYSTEM TESTING AND EVALUATION</b>	<b>95</b>
5.1 INTRODUCTION	95
5.2 TESTING METHODOLOGY	95
5.3 TEST ENVIRONMENT AND DATASET	96
5.4 EVALUATION METRICS	96
5.5 TEST RESULTS AND ANALYSIS	97
5.5.1 Classifier Performance	97
5.5.2 Text Extraction Performance	99
5.5.3 Functional Testing Results	101
5.5.4 Qualitative Evaluation	101
5.6 DISCUSSION OF LIMITATIONS	102



5.7 CHAPTER SUMMARY	102
<b>CHAPTER 6: SUMMARY, CONCLUSION, AND RECOMMENDATIONS</b>	<b>103</b>
6.1 INTRODUCTION	103
6.2 SUMMARY OF WORK	103
6.3 CONCLUSION	104
6.4 RECOMMENDATIONS FOR FUTURE WORK	104
<b>REFERENCES</b>	<b>106</b>
<b>APPENDIX A: FULL SOURCE CODE</b>	<b>111</b>
<b>APPENDIX B: USER MANUAL</b>	<b>124</b>
<b>APPENDIX C: LIST OF ACRONYMS</b>	<b>127</b>

## LIST OF FIGURES

FIGURES	PAGES
FIGURE 3.3.1: TEXTUAL DIAGRAM FOR PDF TEXT EXTRACTION SYSTEM	44
FIGURE 3.4.1: VISUALIZATION OF THE FLOW OF DATA	48
FIGURE 3.4.2: CONTEXT DIAGRAM OF THE PDF TEXT EXTRACTION SYSTEM	49
FIGURE 3.4.3: STRUCTURED DIAGRAM OF THE PDF TEXT EXTRACTION SYSTEM	51
FIGURE 3.5.1: CLASS DIAGRAM	53
FIGURE 3.5.2: USE CASE DIAGRAM	54
FIGURE 3.5.3: SEQUENCE DIAGRAM	55
FIGURE 3.6.1: FLOW CHART	59
FIGURE 3.7.1: ER-MODEL	63
FIGURE 4.3.1: REGISTRATION PAGE	87
FIGURE 4.3.2: LOGIN PAGE	88
FIGURE 4.3.3: MAIN DASHBOARD AND FILE UPLOADER	90
FIGURE 4.3.4: COMPRESSED DOWNLOADABLE RESULT	92
FIGURE 4.3.5: RESULT PREVIEW AND DOWNLOAD	93

## LIST OF TABLES

<b>TABLES</b>	<b>PAGES</b>
TABLE 2.3.1: PDF TEXT EXTRACTION LIBRARY COMPARISON	21
TABLE 3.3.1: CORE COMPONENTS OF THE PDF TEXT EXTRACTION SYSTEM	45
TABLE 3.3.2: TECHNOLOGIES USED FOR THE PDF TEXT EXTRACTION SYSTEM	46
TABLE 3.5.1: DATA DICTIONARY	56
TABLE 5.5.1: CLASSIFIER CONFUSION MATRIX	98
TABLE 5.5.2: TEXT EXTRACTION PERFORMANCE RESULTS	100
TABLE 6.4.1: LIST OF ACRONYMS	127

# CHAPTER 1: INTRODUCTION

## 1.1 BACKGROUND OF THE STUDY

The Portable Document Format (PDF) has become one of the most widely used standards for digital documents due to its platform independence, consistent formatting, and reliability in information sharing across domains such as education, healthcare, finance, and law (Adhikari & Agarwal, 2024). Despite its advantages, extracting meaningful text from PDF files remains a persistent challenge. The complexity of internal structures, varying layouts, and the presence of non-text elements such as tables, images, and scanned pages make text extraction prone to errors and inefficiencies (Compdf, 2025).

Manual extraction of text from PDFs is not only time-consuming but also error-prone, particularly when dealing with large datasets or documents with complex formatting. Studies indicate that manual approaches contribute to data quality issues and significantly reduce productivity (Auer et al., 2022). These limitations highlight the importance of automated approaches to ensure accuracy, scalability, and efficiency in processing large volumes of documents.

Recent advances in automated text extraction have shown promise. Python-based libraries such as PyMuPDF and pdfplumber enable structured and unstructured text extraction, while Optical Character Recognition (OCR) tools like Tesseract extend functionality to image-based or scanned PDFs (Adhikari & Agarwal, 2024). However, inconsistencies in spacing, garbled characters, and formatting loss continue to pose challenges (Compdf, 2025). To address these issues, robust solutions must integrate layout-aware parsing, OCR technologies, and scalable architectures. For instance, cloud-based pipelines have demonstrated the ability to process over one million PDF pages per hour by combining layout analysis with machine learning-based OCR (Auer et al., 2022).

In this context, the development of an automated, accurate, and scalable PDF text extraction system is necessary. Such systems can significantly enhance data accessibility and reliability, enabling timely decision-making in fields where information processing is critical.

## 1.2 PROBLEM STATEMENT

The Portable Document Format (PDF) is widely used for preserving document structure and ensuring cross-platform compatibility, extracting accurate and meaningful text from complex PDF files remains a significant challenge. The unstructured and heterogeneous nature of PDFs, ranging from multi-column layouts to embedded elements, makes manual and even automated extraction error-prone and inefficient (Adhikari & Agarwal, 2024).

Key difficulties arise from the presence of diverse layouts, such as headers, footers, sidebars, and multi-column text, which often result in mixed or disordered content during extraction. Similarly, tables and charts typically lose their structural integrity, producing disorganized or misaligned outputs when converted into plain text (Compdf, 2025). Interactive elements, including forms, hyperlinks, and annotations, add further complexity, as they are often overlooked or inconsistently processed. Additionally, scanned PDFs and image-based documents require Optical Character Recognition (OCR) to be converted into machine-readable text, a process that introduces further risk of inaccuracies (Auer et al., 2022).

The variability in document structures across different sources compounds the problem, as conventional extraction tools struggle to adapt to non-standardized layouts while maintaining contextual accuracy. These limitations highlight the inadequacy of existing manual and semi-automated methods and demonstrate the urgent need for robust, scalable, and accurate automated PDF text extraction systems.

### **1.3 AIM AND OBJECTIVES OF THE STUDY**

#### **1.3.1 AIM OF THE STUDY**

The aim of this study is to develop an automated and scalable PDF text extraction system using Python, incorporating Optical Character Recognition (OCR) and advanced parsing techniques to accurately process diverse document structures while ensuring efficiency, reliability, and usability.

#### **1.3.2 OBJECTIVES OF THE STUDY**

1. To design a scalable and modular system architecture for automated PDF text extraction, incorporating preprocessing, Optical Character Recognition (OCR), and layout-aware parsing components to handle diverse document structures such as multi-column text, tables, and scanned image-based files.
2. To collect and curate a representative dataset of PDF documents from diverse domains (e.g., academic, legal, financial, and scanned documents), ensuring variability in layouts and structures for robust testing and training of the extraction system.
3. To implement and integrate Python-based extraction libraries (e.g., PyMuPDF, pdfplumber, pdfminer) with OCR engines such as Tesseract, and to apply machine learning or natural language processing techniques where necessary (e.g., for layout detection, table recognition, and error correction).

4. To evaluate the accuracy, efficiency, and robustness of the system using appropriate performance metrics, including precision, recall, F1-score, processing time, and structural similarity indices, while optimizing the system for batch processing, speed, and resource utilization.

## 1.4 SCOPE OF THE STUDY

This study focuses on the design and implementation of an automated system for text extraction from PDF documents using Python-based tools and machine learning–driven approaches. The system is developed to handle both text-based and image-based PDFs by integrating Optical Character Recognition (OCR) and layout-aware parsing techniques. The scope specifically covers the extraction of text from documents containing varied structures such as multi-column layouts, tables, and charts, while ensuring the preservation of document integrity and minimizing common extraction errors.

The dataset for this study will consist of diverse PDF documents drawn from academic, legal, financial, and scanned sources, thereby reflecting real-world variability in formatting and content structures. The system will be evaluated using standard performance metrics, including precision, recall, F1-score, structural similarity, and processing time, in order to determine its accuracy, efficiency, and scalability.

The study is limited to text extraction and does not address advanced downstream tasks such as semantic analysis, document summarization, or natural language understanding. Furthermore, the system is designed for experimental and research purposes; therefore, aspects such as full-scale enterprise deployment, integration with large-scale cloud infrastructures, or multilingual OCR beyond English may be considered outside the scope of this project.

## 1.5 SIGNIFICANCE

This study is significant as it addresses a critical challenge faced by organizations, researchers, and individuals who rely heavily on Portable Document Format (PDF) documents for information management. By developing an automated text extraction system, the project contributes to improving productivity, accuracy, and efficiency across multiple domains.

First, the system supports **data-driven decision-making** by enabling the transformation of unstructured PDF documents into structured, machine-readable text. This facilitates downstream applications such as text mining, natural language processing (NLP), and information retrieval, thereby accelerating research, business intelligence, and legal analysis (Adhikari & Agarwal, 2024). Second, the integration of Optical Character Recognition (OCR) enhances **document searchability**, particularly for scanned and image-

based PDFs, making archival materials and regulatory records more accessible for compliance and auditing purposes (Auer et al., 2022).

Furthermore, the system contributes to **workflow automation and operational efficiency** by reducing reliance on manual data entry, which is often time-consuming and error-prone. Automating these processes not only minimizes human error but also reduces operational costs, supporting digital transformation efforts in data-intensive industries such as healthcare, finance, and education (Compdf, 2025). Finally, the study has implications for **machine learning applications**, as the extracted text can serve as training data for predictive models and NLP tasks such as document classification, summarization, and named entity recognition, thereby advancing broader applications of artificial intelligence.

By addressing these needs, the study provides both practical and academic contributions: practically, by offering a scalable solution to a real-world information management problem; and academically, by advancing knowledge on the integration of OCR, Python-based parsing libraries, and evaluation metrics for complex document processing.

# CHAPTER 2: LITERATURE REVIEW

## 2.1 INTRODUCTION

PDF (Portable Document Format) files are structured in a way that enables them to preserve the appearance of documents across different platforms and devices. The structure of PDF files is based on several key components, such as objects, streams, and page descriptions. Understanding these components is essential when working with text extraction and document manipulation, as each part plays a specific role in how the content is stored and rendered.

## 2.2 PDF FILE FORMAT

### 2.2.1 THE STRUCTURE AND COMPONENTS OF PDF FILES

**1. Objects:** PDFs are composed of **objects**, which are the basic building blocks of a PDF file.

Objects are classified into several types, each serving a specific function. These include:

- a. Boolean:** Represents a true or false value.
- b. Numeric:** Represents numbers, either integers or floating-point numbers.
- c. String:** Represents sequences of characters, often used for textual data or metadata.
- d. Name:** A sequence of characters that serve as identifiers or keys within a PDF.
- e. Array:** An ordered list of objects, which can include any type of object.
- f. Dictionary:** A collection of key-value pairs, where keys are names and values can be any type of object.
- g. Stream:** A sequence of data, often used to store large amounts of binary information like images, fonts, and page content.
- h. Null:** Represents a null value or absence of an object.

These objects are organized hierarchically within a PDF file, allowing for efficient rendering and interaction.

**2. Streams:** are a critical component of PDFs used to store large pieces of data, such as images, fonts, or the actual content of pages (text, graphics). Streams are always paired with a dictionary that defines attributes of the stream, such as the length of the stream data and how the data is encoded.



Key characteristics of streams:

- a. **Content Streams:** These contain the actual instructions for rendering a page, including the text to be displayed, the fonts used, the positioning of objects, and any graphical elements (lines, shapes, etc.).
- b. **Compressed Data:** To optimize file size, streams can be compressed using various algorithms (e.g., Flate, LZW). This compression complicates text extraction because the data must be decoded before being analyzed.

**3. Page Descriptions:** Each page in a PDF file is described using a combination of **page objects** and **content streams**. The page object contains metadata and references to resources such as fonts and images used on the page. The actual layout of the page, what is displayed and where, is defined in content streams, using a markup-like language called **PDF operators**.

Key components in a page description include:

- a. **Page Tree:** A hierarchical structure that organizes the pages in a PDF file. The root node of the page tree contains references to each page in the document, and each page node refers to resources (fonts, images) and the content stream that defines what is rendered on the page.
- b. **Text Objects:** Text in PDFs is defined within **text objects**, which are blocks that specify the positioning and styling of the text (e.g., font, size, color). Text is not stored in a straightforward or linear fashion, which is why extracting text while maintaining its original layout can be challenging.
- c. **Graphics Objects:** Graphics such as lines, shapes, and images are also defined within the content streams using graphics operators. This includes instructions for drawing vector graphics, positioning images, and applying transformations (scaling, rotation).

**4. Page Content and Operators:** PDFs describe the content of a page using a set of **operators** within content streams. These operators define how text, images, and graphics are positioned and rendered. Some common types of operators include:

- a. **Text Operators:** Operators that control how text is placed on the page, such as:
  - BT (Begin Text) and ET (End Text): These define the beginning and end of a text object.
  - Tj and TJ: These are used to show text strings.
  - Tm: Sets the text matrix, which determines the position, rotation, and scale of the text.
- b. **Graphics Operators:** These operators handle drawing shapes and images, such as:

m (MoveTo) and l (LineTo): These define paths for drawing lines and shapes.

f and S: These fill and stroke the paths created with m and l.

Do: Invokes an image or another external resource.

**5. Fonts and Encoding:** PDF files contain embedded **fonts** and **font descriptors** that are used to render text. Fonts can be either standard fonts (like Helvetica) or embedded fonts (custom fonts included in the PDF). Each font in a PDF has an associated encoding, which maps character codes to the actual glyphs that are rendered.

- a. **Font Types:** PDF supports several types of fonts, including TrueType, Type 1, and Type 3 fonts. The choice of font type affects how text is stored and rendered.
- b. **Character Encoding:** PDFs may use custom encodings that map character codes to specific glyphs in a font. This complicates text extraction, as the extracted character codes may need to be mapped back to readable text using the appropriate encoding.

**6. Resources:** Each page in a PDF can have associated **resources** that are used to render its content, such as fonts, images, and color spaces. The page object contains references to these resources, ensuring that the correct visual elements are available when the page is rendered.

- a. **Font Resources:** Defines the fonts available for use on a page.
- b. **Image Resources:** Images can be embedded directly into a page or referenced from an external file.
- c. **Color Spaces:** Defines the color model used for the content, which can include grayscale, RGB, CMYK, or spot colors.

## 2.2.2 DOCUMENT STRUCTURE

The overall structure of a PDF is made up of four main sections:

- 1. Header:** This contains basic information about the PDF file, such as the version of the PDF format being used.
- 2. Body:** The main part of the PDF, which contains all the objects, streams, and page descriptions.
- 3. Cross-Reference Table:** This table allows quick access to objects within the PDF by providing the byte offset of each object in the file. This structure enables efficient searching and rendering

of content.

4. **Trailer:** The trailer marks the end of the file and includes references to the cross-reference table and the root object (which points to the document's catalog and page tree).

### 2.2.3 THE CHALLENGES OF PARSING PDF FILES DUE TO THEIR COMPLEXITY AND POTENTIAL FOR VARIATIONS.

Parsing PDF files presents significant challenges due to the inherent complexity and potential variations in how PDF documents are structured. Unlike simpler file formats, PDFs are designed to preserve the exact appearance of a document across different devices and platforms, leading to a non-linear and highly flexible structure that complicates text extraction and other parsing tasks. Below are some of the key challenges, supported by references from the literature:

1. **Non-Linear Document Structure:** PDF files do not store text in a continuous, linear stream, as is the case with plain text files. Instead, text and other content (such as images and tables) are often stored in disconnected segments within content streams. These segments are organized based on the layout of the page rather than the logical flow of the document. This makes it difficult to extract text in the correct reading order, especially for multi-column layouts, footnotes, headers, and sidebars.
2. **Complex Page Layouts:** PDF files can have highly complex page layouts, including multiple columns, nested tables, text boxes, images, and floating elements. These elements are positioned precisely on the page using operators within content streams, and there is no standard way to define the logical structure of the document, such as chapters, sections, or paragraphs. Extracting text from PDFs while preserving its layout (e.g., maintaining table structures or detecting multi-column text) can be difficult because the layout is often intertwined with graphical elements and positioning instructions.
3. **Text as Graphics:** In some PDFs, especially those generated from scanned documents or certain design tools, text may be stored as graphical elements rather than as machine-readable text. For example, scanned PDFs contain images of text rather than the text itself. In these cases, Optical Character Recognition (OCR) is required to extract text, but OCR is prone to errors, especially when the quality of the scan is poor or when the document uses unusual fonts or handwriting.

- 4. Inconsistent Encoding and Fonts:** Text in PDFs is often stored using different character encodings and font types. Custom fonts and non-standard character encodings are frequently embedded in the PDF, which makes it difficult to map extracted character codes back to readable text. This can result in garbled or incorrect text during extraction, especially when non-Latin characters, special symbols, or mathematical notations are involved.
- 5. Lack of Semantic Structure:** Unlike HTML or XML documents, which include explicit tags to define headings, paragraphs, lists, and other semantic elements, PDF files typically do not include such semantic information. PDFs are designed to preserve the visual appearance of a document, but they often lack explicit markers that describe the logical structure of the content. As a result, text extraction tools must rely on heuristic approaches, such as analyzing font size, text positioning, and indentation, to infer structure (e.g., determining whether a line of text is a heading or part of a paragraph).
- 6. Variations in PDF Generators:** PDFs can be generated using a wide variety of software tools, each of which may create PDFs with different internal structures. Some PDFs are well-organized, with text stored in logical blocks, while others may store each individual character or word as a separate object, making it difficult to extract coherent sentences or paragraphs. Moreover, PDFs generated from different tools may store content in different ways, making it hard to develop a one-size-fits-all parser.
- 7. Tables and Structured Data:** Extracting structured data such as tables is particularly difficult because PDFs often store tables as a collection of individual text elements and graphical lines rather than as structured tables. There is no standard representation for tables in PDFs, which makes it challenging to detect table boundaries, interpret the relationships between cells, and ensure that the extracted data maintains its original structure.
- 8. Image-Based PDFs:** PDF files created by scanning paper documents often contain no machine-readable text at all. These PDFs store images of the document pages, making it impossible to directly extract text. OCR techniques must be applied to convert the images to text, but OCR accuracy depends on the quality of the scan, the clarity of the text, and the complexity of the document (e.g., handwritten notes, low-resolution scans).

## 2.3 TEXT EXTRACTION TECHNIQUES

### 2.3.1 OPTICAL CHARACTER RECOGNITION (OCR)

Optical Character Recognition (OCR) engines are essential tools for converting image-based PDFs or scanned documents into machine-readable text. Several OCR engines are available, each with unique strengths and weaknesses in terms of accuracy, speed, ease of use, and ability to handle different languages and document formats. This discussion focuses on two popular OCR engines, **Tesseract** and its Python wrapper, **Pytesseract**, along with other notable OCR engines, highlighting their key features, strengths, and weaknesses.

#### 2.3.1.1 DIFFERENT OCR ENGINES

1. **Tesseract OCR:** Tesseract is one of the most widely used and powerful open-source OCR engines. Originally developed by Hewlett-Packard in the 1980s, it is now maintained by Google. Tesseract can process images and PDFs, extract text from them, and handle a wide variety of languages and scripts.

#### Strengths:

- a. **Open-Source and Free:** One of the biggest advantages of Tesseract is that it is open-source and freely available. It can be easily integrated into applications and modified according to specific needs.
- b. **Multilingual Support:** Tesseract supports over 100 languages, including complex scripts like Arabic, Chinese, and Hindi. It can also be trained to recognize new languages or custom fonts using language models.
- c. **Accuracy:** Tesseract generally performs well on high-quality images and printed text. It has an option for page segmentation modes that allow users to customize the approach based on the layout of the document.
- d. **Custom Training:** Tesseract allows users to train the model on specific fonts or datasets, improving recognition accuracy for specialized documents or rare fonts.
- e. **Integration with PDF and Image Files:** Tesseract can work directly on PDFs or images and convert the extracted text into a structured format.

## Weaknesses:

- a. **Accuracy on Low-Quality Images:** Tesseract's accuracy decreases significantly when processing low-quality images, noisy backgrounds, or documents with non-standard fonts. It may struggle with documents that have poor resolution, complex layouts, or significant noise.
  - b. **Complex Layouts:** Tesseract does not handle complex layouts well, such as multi-column formats, tables, or documents with many images and text interspersed. Additional processing steps or layout analysis is required for such cases.
  - c. **Limited Layout Retention:** While Tesseract can extract text, it does not preserve the layout of the original document (e.g., tables or columns), making it less suitable for tasks that require structured output.
  - d. **Slow Performance:** Compared to commercial OCR engines, Tesseract may be slower in processing large documents, especially when handling high-resolution images or PDFs.
2. **Pytesseract** is a Python wrapper for the Tesseract OCR engine, making it easier for Python developers to integrate Tesseract into their applications. It simplifies the process of passing images to Tesseract and retrieving the extracted text, making it a popular choice for developers working in Python.

## Strengths:

- a. **Python Integration:** Pytesseract provides seamless integration with Python, making it easy to use in Python-based workflows for image and document processing. It works well with other libraries like Pillow and OpenCV, allowing for pre-processing images (resizing, noise reduction, thresholding) before sending them to Tesseract.
- b. **Simplified API:** Pytesseract simplifies interaction with Tesseract by providing easy-to-use functions for passing image files and retrieving text. This reduces the complexity involved in working directly with the Tesseract command line interface.
- c. **Flexible Input:** Pytesseract can handle various input types, including file paths, image objects, or binary image data. This flexibility is useful for developers working with different data sources (e.g., files, cameras, or screenshots).

### Weaknesses:

- a. **Dependent on Tesseract:** Since Pytesseract is a wrapper around Tesseract, it inherits all the strengths and weaknesses of the underlying Tesseract engine. This includes its difficulty with complex layouts and low-quality images.
- b. **Performance:** Pytesseract may not be the fastest solution for real-time OCR applications, as it relies on Tesseract's relatively slower processing times for large documents or images.

**3. ABBYY FineReader** is a commercial OCR engine known for its high accuracy and excellent layout retention. It is widely used in industries that require high-quality document digitization, such as legal and financial sectors.

### Strengths:

- a. **Highly Accurate:** ABBYY FineReader is known for its high accuracy in recognizing text, even in low-quality images or scanned documents. It performs exceptionally well with complex layouts and multi-column documents.
- b. **Layout Retention:** FineReader is excellent at preserving the original layout of documents, including tables, columns, and embedded images, making it a popular choice for users who need structured output (e.g., financial statements or legal contracts).
- c. **Wide Language Support:** Like Tesseract, ABBYY supports a wide range of languages and can recognize characters from non-Latin scripts with high precision.

### Weaknesses:

- a. **Commercial License:** FineReader is a paid software with a commercial license, which can be costly for individuals or organizations seeking a free solution.
- b. **Limited Customization:** FineReader is not open-source, so users cannot customize or extend it as easily as they can with open-source solutions like Tesseract.

**4. Google Cloud Vision OCR** is a cloud-based OCR service that offers high accuracy and scalability. It is part of Google Cloud's suite of machine learning tools.

### Strengths:

- a. **High Accuracy and Scalability:** Google Cloud Vision OCR provides high accuracy, especially for documents with printed text and complex layouts. It can handle large volumes of documents, making it ideal for organizations with big data needs.

- b. **Integration with Cloud Services:** Being a cloud-based service, it integrates well with other Google Cloud products, providing scalability and powerful post-processing options for developers.
- c. **Multilingual Support:** Google Cloud Vision supports many languages and scripts, making it suitable for global use cases.

#### **Weaknesses:**

- a. **Cloud Dependency:** The reliance on cloud services can be a drawback for users concerned about privacy or needing an offline solution. There are also ongoing costs associated with using the service.
- b. **Limited Customization:** Like ABBYY FineReader, Google Cloud Vision OCR is not customizable in the way that open-source solutions are.

### **2.3.2 THE PROCESS OF IMAGE PREPROCESSING**

Image preprocessing is a critical step in improving the accuracy of Optical Character Recognition (OCR). The process aims to enhance the quality of input images by removing noise, correcting distortions, and making the text clearer for OCR engines. Various preprocessing techniques such as **binarization**, **deskewing**, **noise removal**, and **thresholding** are commonly employed. Below is an explanation of these techniques, supported by references from the literature.

#### **2.3.2.1 BINARIZATION**

Binarization is the process of converting a grayscale or color image into a binary image where each pixel is either black (foreground) or white (background). This simplification helps the OCR engine focus on text by removing irrelevant information like colors or gradients.

**Thresholding:** Thresholding is the most common method used for binarization. It involves setting a threshold intensity value, where pixels with values below the threshold are converted to black (representing text), and those above are converted to white (representing the background). Thresholding can be global or adaptive:

1. **Global Thresholding:** A single threshold value is applied across the entire image. It works well for uniformly illuminated images.
2. **Adaptive Thresholding:** Different threshold values are applied based on local regions of the image, which is useful for images with varying lighting conditions or shadows.



#### 2.3.2.1.1 Importance of Binarization:

1. Binarization simplifies the image for OCR by reducing noise and focusing on text.
2. It improves contrast between the text and the background, which helps OCR engines distinguish characters more effectively.

#### 2.3.2.1.2 Challenges of Binarization:

In cases of non-uniform illumination, shadows, or textured backgrounds, global thresholding may fail. Adaptive methods like **Otsu's thresholding** or **adaptive Gaussian thresholding** are preferred for handling these cases.

#### 2.3.2.2 DESKEWING

Deskewing refers to the process of correcting the orientation of an image when it is not aligned properly. Documents scanned at an angle can cause skewed text, making it difficult for OCR engines to accurately recognize characters.

##### 2.3.2.2.1 Steps:

1. **Skew Angle Detection:** The first step in deskewing is detecting the angle by which the image is skewed. This can be done by analyzing the alignment of text lines using techniques like the Hough Transform or by identifying the dominant orientation of horizontal text lines.
2. **Rotation:** After the skew angle is detected, the image is rotated to align the text horizontally. This step ensures that the text is correctly positioned for OCR recognition.

##### 2.3.2.2.2 Importance:

Deskewing ensures that the text is properly aligned, which improves OCR accuracy, particularly for lines of text and characters that might otherwise appear distorted.

##### 2.3.2.2.3 Challenges:

In highly complex or noisy images, deskewing can be difficult, especially if the text is arranged in multiple orientations or contains non-text elements (such as images or diagrams).

### 2.3.2.3 NOISE REMOVAL

Noise in an image refers to unwanted artifacts like random black or white spots, lines, or distortions that can interfere with OCR. These artifacts can arise due to scanning issues, low resolution, or poor image quality.

#### 2.3.2.3.1 Steps:

- 1. Median Filtering:** This method replaces each pixel's value with the median value of the surrounding pixels, which helps reduce noise without blurring edges.
- 2. Gaussian Smoothing:** A Gaussian filter blurs the image slightly, reducing high frequency noise while preserving the edges of text.
- 3. Morphological Operations:** Techniques such as **erosion** and **dilation** are used to remove small noise components while preserving the shape of the text. Erosion removes small white noise, while dilation fills small black gaps in the text.

#### 2.3.2.3.2 Importance:

Reducing noise improves OCR accuracy by preventing the engine from misinterpreting noise artifacts as text characters.

#### 2.3.2.3.3 Challenges:

Over-filtering can lead to the loss of important details, especially in small or thin text characters, leading to inaccuracies in OCR.

### 2.3.2.4 THRESHOLDING AND CONTRAST ENHANCEMENT

Thresholding is closely related to binarization and refers to setting a threshold value to enhance the contrast between the text and the background. Contrast enhancement increases the distinction between characters and the background, making them more distinguishable.

#### 2.3.2.4.1 Steps:

- 1. Contrast Stretching:** This method involves adjusting the intensity levels of the image to span the full range of available pixel values (0 to 255 in an 8-bit image), making the text more visible.

**2. Histogram Equalization:** This technique redistributes the intensity levels in the image to improve contrast. It can be particularly useful for images with low contrast or uneven illumination.

#### **2.3.2.4.2 Importance:**

Enhancing contrast ensures that OCR engines can better distinguish text from the background, especially in faint or low-contrast images.

### **2.3.2.5 DILATION AND EROSION (MORPHOLOGICAL OPERATIONS)**

Morphological operations like **dilation** and **erosion** are used to clean up small imperfections in the text or background of the image. These operations are applied primarily in binary images.

#### **2.3.2.5.1 Steps:**

- 1. Dilation:** Expands the boundaries of objects in the image (i.e., the black text). This can help close small gaps or breaks in characters.
- 2. Erosion:** Shrinks the boundaries of objects in the image, helping to eliminate small white noise spots.

#### **2.3.2.5.2 Importance:**

These operations can help restore broken or thin text and remove small imperfections that might otherwise confuse the OCR engine.

#### **2.3.2.5.3 Challenges:**

Overuse of dilation can lead to merging of adjacent characters, and excessive erosion can cause characters to become too thin or break apart.

### **2.3.2.6 RESIZING AND RESOLUTION ADJUSTMENT**

OCR engines perform best when the input image has sufficient resolution. Low-resolution images often lead to errors in character recognition.

#### **2.3.2.6.1 Steps:**

- 1. Upsampling:** Low-resolution images can be resized to a higher resolution to improve text clarity, though this can introduce artifacts if done excessively.

**2. Optimal DPI:** For most OCR tasks, an image resolution of around 300 DPI (dots per inch) is recommended, as it provides a balance between file size and text clarity.

#### **2.3.2.6.2 Importance:**

Correcting resolution and resizing the image ensures that text is sharp and clear, which improves OCR accuracy.

#### **2.3.2.6.3 Challenges:**

Upsampling a low-quality image may not always result in better OCR performance, as it cannot add information that was not originally present.

### **2.3.3 PDF PARSING LIBRARIES:**

Python has several popular libraries for parsing PDFs, each offering unique features and capabilities for extracting text and other content from PDF documents. The most commonly used libraries are **PyPDF2**, **PDFMiner**, and **pdfplumber**. This review will discuss their features, strengths, weaknesses, and appropriate use cases, along with references to relevant documentation and studies.

#### **2.3.3.1 PyPDF2**

**PyPDF2** is a pure-Python library that focuses on working with PDF files. It allows users to extract metadata, split and merge PDFs, rotate pages, and extract text. Although PyPDF2 offers basic text extraction, its primary strength lies in its ability to manipulate and manage PDF files rather than text extraction alone

##### **Key Features:**

- 1. Splitting and Merging:** PyPDF2 excels at splitting PDF files into individual pages and merging multiple PDFs into a single file, making it useful for document management.
- 2. Rotation:** It allows users to rotate pages in a PDF, which can be useful when dealing with scanned documents that may be misaligned.
- 3. Extracting Metadata:** PyPDF2 can extract metadata such as the author, creation date, and title of the PDF.
- 4. Encryption:** PyPDF2 can encrypt and decrypt PDF files with password protection.

### Strengths:

1. **Simple and Lightweight:** PyPDF2 is easy to use for basic PDF operations like splitting, merging, and rotating pages.
2. **Well-Suited for Manipulating PDF Structure:** For tasks involving rearranging, splitting, and combining PDF pages, PyPDF2 is ideal.

### Weaknesses:

1. **Limited Text Extraction:** PyPDF2's text extraction capabilities are limited, especially when dealing with complex PDFs that have tables, multiple columns, or embedded images.
2. **Inconsistent Results:** The library can produce inconsistent results when extracting text, especially if the document has a complex structure or non-standard fonts.

### Use Case:

PyPDF2 is best suited for tasks like **splitting, merging, or rotating PDFs** and extracting simple text from relatively straightforward PDF documents. However, for more complex parsing, other libraries may be more effective.

#### 2.3.3.2 PDFMINER

**PDFMiner** is a highly comprehensive and powerful library designed specifically for text extraction and analysis from PDFs. It can handle complex PDFs with various structures, including text, fonts, and layouts, making it one of the most robust tools available for parsing and extracting content from PDFs.

### Key Features:

1. **Text Extraction:** PDFMiner offers excellent text extraction capabilities, including handling of fonts, character encodings, and layout analysis.
2. **Layout Preservation:** It can preserve the layout and structure of the PDF, including multi-column text and embedded elements.
3. **Support for Various Output Formats:** PDFMiner can output extracted text in different formats, such as HTML or XML, enabling users to preserve the formatting of the original document.

4. **Handling of Complex PDFs:** It excels at parsing complex PDF files, including scanned documents, multi-column layouts, and those with embedded images and diagrams.

#### **Strengths:**

1. **Accurate Text Extraction:** PDFMiner is one of the most accurate libraries for extracting text from PDFs, even when dealing with complex documents that contain a mix of text, images, and non-standard fonts.
2. **Layout and Structure Retention:** PDFMiner preserves the layout and formatting of the PDF, making it ideal for use cases where the original document's structure is important.
3. **Supports Advanced Features:** It can extract detailed information about the document, including font types, character positioning, and encoding.

#### **Weaknesses:**

1. **Slower Performance:** PDFMiner tends to be slower than other libraries, especially for large PDFs or complex documents due to its detailed processing.
2. **Complex API:** It has a steeper learning curve and can be difficult to use for beginners, especially when compared to simpler libraries like PyPDF2.
3. **Limited Documentation:** Although it is powerful, the library's documentation can be hard to follow, requiring users to experiment with different methods to achieve the desired results.

#### **Use Case:**

PDFMiner is ideal for extracting text from **complex PDFs** with advanced layouts and where preserving the original structure is important. It's also suitable for detailed document analysis and conversion of PDFs into different structured formats like XML.

#### **2.3.3.3 PDFPLUMBER**

**pdfplumber** is another Python library designed for extracting text, tables, and other objects from PDFs. It is particularly known for its ability to handle **table extraction**, which makes it a great option for parsing structured PDFs.

## Key Features:

- 1. Table Extraction:** pdfplumber excels at identifying and extracting tables from PDFs. It can detect and parse tables, including their structure and individual cell contents.
- 2. Detailed Text and Image Extraction:** Like PDFMiner, pdfplumber can extract detailed text and images from PDFs, including multi-column layouts and embedded graphics.
- 3. Visualization Tools:** It provides tools for visualizing the layout of a PDF, which helps in understanding its structure and improving the accuracy of text or table extraction.
- 4. High Customizability:** Users can specify custom boundaries and regions from which to extract content, allowing for more granular control over the parsing process.

## Strengths:

- 1. Best for Table Extraction:** pdfplumber is the most reliable option for extracting tabular data from PDFs, which is often a challenge with other libraries.
- 2. Good for Complex Layouts:** It handles complex layouts well, including multi-column text, images, and tables.
- 3. Intuitive API:** pdfplumber has a relatively easy-to-use and intuitive API, making it accessible for users without advanced programming skills.

## Weaknesses:

- 1. Performance:** While it is effective at parsing complex documents, pdfplumber may be slower than simpler tools like PyPDF2 when handling large or complicated PDFs.
- 2. Less Suitable for Simple PDFs:** If the goal is to extract plain text from simple PDF files, pdfplumber may be overkill.

## Use Case:

pdfplumber is highly effective for tasks involving **table extraction** and working with PDFs containing complex layouts. It's also suitable for use cases where understanding and visualizing the structure of a PDF is necessary.

**Table 2.3.1: PDF Text Extraction Library Comparison**

<b>Feature</b>	<b>PyPDF2</b>	<b>PDFMiner</b>	<b>Pdfplumber</b>
<b>Text Extraction</b>	Basic	Advanced	Advanced
<b>Layout Preservation</b>	Minimal	Excellent	Good
<b>Table Extraction</b>	Limited	Moderate	Excellent
<b>Manipulation (Split/Merge)</b>	Excellent	Limited	Limited
<b>Ease of Use</b>	Easy	Complex	Moderate
<b>Performance</b>	Fast	Slower	Moderate
<b>Use Case</b>	Simple PDF manipulation and text extraction	Detailed text extraction from complex PDFs	Table extraction, detailed text extraction



### 2.3.4 NATURAL LANGUAGE PROCESSING (NLP) TECHNIQUES:

Natural Language Processing (NLP) techniques can significantly enhance the text extraction process from PDFs by not only extracting raw text but also providing more meaningful, structured information. These techniques enable us to understand the content at a deeper level by analyzing the semantics, structure, and intent of the extracted text. Here are some key NLP techniques and their applications in text extraction:

#### 2.3.4.1 NAMED ENTITY RECOGNITION (NER)

**Named Entity Recognition (NER)** is an NLP technique used to identify and classify named entities (such as persons, organizations, locations, dates, and more) in a text. When applied to text extracted from PDFs, NER can help organize and categorize the extracted information, making it more useful for downstream applications.

##### 2.3.4.1.1 Applications:

- 1. Information Structuring:** NER can help categorize large volumes of unstructured text from PDF documents by extracting and labeling specific entities like people, locations, dates, and organizations. For example, in legal documents, NER can identify parties involved, dates of contracts, and specific legal terms.
- 2. Automated Indexing:** Extracted text from academic papers, reports, or research papers can be automatically indexed using NER to highlight important concepts and topics, making it easier to search and analyze.
- 3. Summarization:** Extracted text from PDFs can be summarized more effectively by focusing on the identified named entities, providing a concise overview of key elements.

##### 2.3.4.1.2 Tools and Libraries:

- 1. spaCy:** A popular Python library for NLP that provides a fast and efficient implementation of NER.
- 2. NLTK:** The Natural Language Toolkit provides a wide range of tools for performing NER, along with other NLP tasks.

3. **Stanford NER:** A widely used tool that provides high-quality named entity recognition for various languages.

#### 2.3.4.2 SENTIMENT ANALYSIS

**Sentiment Analysis** is used to determine the sentiment or emotional tone behind a piece of text. It classifies the text into positive, negative, or neutral sentiment. When applied to extracted text from PDFs, sentiment analysis can provide insights into the overall mood or attitude expressed in the document.

##### 2.3.4.2.1 Applications:

1. **Customer Feedback:** In customer surveys or feedback forms extracted from PDFs, sentiment analysis can automatically assess the overall satisfaction of customers based on their comments.
2. **Product Reviews:** When extracting text from product reviews (e.g., PDFs of user feedback), sentiment analysis can help determine the general opinion about a product or service.
3. **Social media and Marketing:** Reports and PDFs containing social media data or marketing campaigns can benefit from sentiment analysis to gauge public sentiment toward a brand or product.

##### 2.3.4.2.2 Tools and Libraries:

1. **VADER (Valence Aware Dictionary and sEntiment Reasoner):** A sentiment analysis tool designed specifically for analyzing social media and informal texts, but can be applied to various types of text.
2. **TextBlob:** A simple NLP library that provides sentiment analysis capabilities based on rule-based methods.
3. **BERT-based Models:** Pre-trained models such as **BERT** (Bidirectional Encoder Representations from Transformers) can be fine-tuned for more sophisticated sentiment analysis.

#### 2.3.4.3 TOPIC MODELING

**Topic Modeling** is a technique used to identify the hidden thematic structure in a large collection of text. It helps in uncovering the main topics discussed in a document, which is particularly useful when working with long or complex PDFs.

#### 2.3.4.3.1 Applications:

1. **Document Summarization:** Topic modeling can be used to identify the main themes or subjects covered in large PDF documents, making it easier to summarize and understand the key takeaways.
2. **Content Organization:** For PDFs containing large text corpora, such as academic articles, reports, or books, topic modeling helps to organize the content by clustering related text around similar topics.
3. **Search and Information Retrieval:** Extracted text from PDFs can be grouped by topic, enabling better search functionality and improved document navigation.

#### 2.3.4.3.2 Tools and Libraries:

1. **Latent Dirichlet Allocation (LDA):** A popular algorithm for topic modeling that identifies groups of words (topics) in large text corpora.
2. **Gensim:** A Python library that provides easy-to-use implementations of LDA and other topic modeling algorithms.
3. **BERTopic:** A more advanced, BERT-based topic modeling library that can produce higher-quality results by leveraging contextual embeddings.

#### 2.3.4.4 TEXT SUMMARIZATION

**Text Summarization** is an NLP task that condenses a document into a shorter version while retaining its most important information. This technique can be applied to the text extracted from PDFs to generate concise summaries, particularly for lengthy documents.

##### 2.3.4.4.1 Applications:

1. **Executive Summaries:** Automatic generation of executive summaries from reports, research papers, or financial documents.
2. **Content Preview:** Summarizing content from long PDFs, such as books or articles, to provide users with an overview before they dive into the full document.

- 3. Legal Document Analysis:** Extracted text from legal PDFs can be summarized to highlight the most critical clauses or decisions.

#### **2.3.4.4.2 Types of Summarizations:**

- 1. Extractive Summarization:** This method selects important sentences directly from the original text.
- 2. Abstractive Summarization:** This method generates a new summary based on the understanding of the text's meaning.

#### **2.3.4.4.3 Tools and Libraries:**

- 1. Hugging Face Transformers:** A library that offers pre-trained models for extractive and abstractive text summarization.
- 2. Sumy:** A Python library for extractive summarization that supports different algorithms, such as TextRank and LexRank.

#### **2.3.4.5 PART-OF-SPEECH TAGGING (POS TAGGING)**

**Part-of-Speech (POS) Tagging** is a basic NLP technique that assigns part-of-speech labels (nouns, verbs, adjectives, etc.) to each word in the text. This is useful for understanding the grammatical structure of the extracted text, which can aid in more advanced text processing tasks.

##### **2.3.4.5.1 Applications:**

- 1. Grammatical Analysis:** POS tagging helps in identifying the structure and syntax of sentences, making it easier to process text for tasks like summarization or translation.
- 2. Named Entity Identification:** POS tagging often forms the basis for more advanced tasks like NER and dependency parsing.
- 3. Keyword Extraction:** By tagging parts of speech, keywords (such as important nouns or verbs) can be extracted more efficiently from PDF text.

##### **2.3.4.5.2 Tools and Libraries:**

- 1. spaCy:** Provides fast and efficient POS tagging capabilities along with other NLP features.

**2. NLTK:** A widely used NLP library that includes tools for POS tagging.

#### **2.3.4.6 DEPENDENCY PARSING**

**Dependency Parsing** is a syntactic analysis method that shows the grammatical relationships between words in a sentence. It maps how words are related to each other, which is useful for understanding sentence structure.

##### **2.3.4.6.1 Applications:**

- 1. Semantic Analysis:** Understanding the relationships between words can help extract meaning and intent from the text.
- 2. Information Extraction:** Dependency parsing can be used to extract subject-verb-object relationships, which can be useful for analyzing extracted text from research papers, legal documents, or news articles.

##### **2.3.4.6.2 Tools and Libraries:**

- 1. spaCy:** Offers dependency parsing along with visualization tools to understand sentence structures.
- 2. Stanford NLP:** Provides tools for syntactic and semantic parsing.

#### **2.3.5 THE POTENTIAL FOR IMPROVING ACCURACY AND EXTRACTING MEANINGFUL INFORMATION.**

The potential for improving the accuracy of text extraction and extracting meaningful information from PDF documents has grown significantly with advancements in **Natural Language Processing (NLP)**, **machine learning**, and **optical character recognition (OCR)** technologies. The integration of these tools offers opportunities to go beyond raw text extraction, enabling the interpretation, categorization, and summarization of content in ways that were previously challenging.

##### **2.3.5.1 IMPROVING ACCURACY THROUGH NLP TECHNIQUES**

###### **1. Named Entity Recognition (NER) and Contextual Understanding**

By leveraging NLP techniques like **Named Entity Recognition (NER)**, systems can extract meaningful entities (such as names of people, organizations, dates, and locations) from the text. This improves

accuracy by focusing on contextually relevant information instead of merely extracting raw text data. For instance:

- a. **Contextual Models:** Modern models like **BERT (Bidirectional Encoder Representations from Transformers)** are capable of understanding context in text, making entity recognition more accurate, even in complex PDF structures. BERT's ability to grasp the nuances of human language by analyzing the entire text simultaneously can improve both the precision and recall of extracted information.

## 2. Sentiment Analysis for Contextual Understanding

Sentiment analysis allows for a deeper understanding of the tone and emotional content in text extracted from PDFs, which can help categorize or prioritize documents based on user needs.

- a. **Enhancing Relevance:** For example, analyzing customer feedback extracted from reports or product reviews can yield insights on user satisfaction and help prioritize responses based on the emotional intensity of the comments.

## 3. Topic Modeling for Information Structuring

**Topic modeling** techniques, such as **Latent Dirichlet Allocation (LDA)**, can identify dominant topics in a document, allowing for the extraction of highly relevant and thematic content. By clustering words that frequently occur together, topic models can help extract meaningful information and structure large documents such as research papers, legal documents, or reports.

- a. **Enhanced Document Summarization:** When combined with text summarization techniques, topic models can provide more concise yet comprehensive summaries, thus improving the readability and value of extracted text.

### 2.3.5.2 IMPROVING ACCURACY WITH MACHINE LEARNING

#### 1. Custom Training of OCR Engines

Traditional OCR engines, such as **Tesseract**, often struggle with complex PDF structures, especially when dealing with images, non-standard fonts, or noisy scanned documents. However, the accuracy of OCR has been significantly enhanced by the introduction of machine learning-based methods.

- a. **Fine-Tuning OCR Models:** Custom-trained OCR models, fine-tuned for specific document types or languages, can drastically improve text recognition accuracy in complex PDFs. For

example, training a model on legal or medical documents with specific jargon can reduce errors in text extraction.

## 2. Semantic Analysis and Deep Learning Models

Deep learning models, such as **GPT (Generative Pre-trained Transformer)**, allow for **semantic understanding** of text, which can result in improved interpretation and summarization of extracted information. These models can:

- a. **Correct Text Errors:** Using semantic analysis to identify and correct contextually incorrect text extracted from PDFs.
- b. **Generate Summaries:** Provide abstractive summaries that generate human-like summaries from text, improving document comprehension and accuracy.

### 2.3.5.3 TEXT NORMALIZATION AND PREPROCESSING

#### Image Preprocessing for OCR Accuracy

Preprocessing steps, such as **image binarization**, **deskewing**, and **denoising**, improve the quality of scanned images, which is crucial for enhancing the accuracy of text extraction from PDFs. By removing noise and correcting distortions, OCR engines can more accurately recognize characters and words in the document.

- a. **Improved Recognition:** For instance, deskewing a scanned document before applying OCR can lead to higher text recognition rates, particularly when dealing with poorly scanned or misaligned PDFs.

### 2.3.5.4 HANDLING COMPLEX PDF STRUCTURES

#### Layout-Aware Parsing for Tabular Data

PDF documents often contain tables, charts, and figures that traditional text extraction tools struggle with. By using **layout-aware parsing**, libraries such as **pdfplumber** and **PDFMiner** can extract text while maintaining the structural integrity of tables and charts.

- a. **Improving Table Extraction:** By combining NLP and layout-aware techniques, the extraction of tables and structured data from complex PDFs becomes more accurate. NLP methods like rule-based text splitting can further improve the representation of structured data.

### 2.3.5.5 POST-PROCESSING TECHNIQUES TO ENHANCE ACCURACY

#### Tokenization and Lemmatization

Post-processing techniques like **tokenization** (breaking text into smaller units) and **lemmatization** (reducing words to their base form) can further enhance the accuracy of extracted text. By normalizing the text, these techniques allow for more effective **information retrieval** and **search capabilities** across documents.

- a. **Error Reduction:** Tokenizing text after extraction reduces ambiguities in word recognition (e.g., differentiating between "read" (past tense) and "read" (present tense)).

### 2.3.5.6 CHALLENGES AND FUTURE DIRECTIONS

Despite the advancements in NLP and machine learning techniques, there remain challenges in extracting accurate and meaningful information from complex or noisy PDFs:

1. **Handwritten Texts:** PDFs containing handwritten content remain a challenge, as OCR engines often struggle with handwriting recognition.
2. **Multilingual Texts:** Text extraction from multilingual PDFs can still be prone to errors, especially for lesser-supported languages or documents containing mixed languages.
3. **Complex Layouts:** Highly complex layouts, such as multi-column articles, text overlapping images, and irregular fonts, can reduce extraction accuracy. However, combining layout-aware techniques and deep learning-based approaches can provide better results in the future.

## 2.4 RELATED WORK

PDF text extraction has been a growing area of research and development due to the increasing demand for document digitization and automation in data processing. Various projects and research studies have focused on leveraging **Python** for text extraction, with varying degrees of success in handling the complexity of PDF files. Below is a review of key existing projects, the methodologies they employ, and their findings.



## **2.4.1 REVIEW OF EXISTING PROJECTS AND RESEARCH RELATED TO PDF TEXT EXTRACTION USING PYTHON**

### **2.4.1.1 PyPDF2**

**PyPDF2** is one of the most commonly used Python libraries for extracting text from PDFs. It provides basic functions for parsing, merging, and splitting PDF files. However, its text extraction capabilities are limited, especially when dealing with complex layouts or documents with embedded images.

#### **2.4.1.1.1 Key Features and Methodology:**

1. Extracts text based on a PDF's internal structure, breaking down content into pages and analyzing objects in each page.
2. Handles basic layout extraction but struggles with tables, multi-column texts, and diagrams.
3. No support for images or scanned PDFs (does not integrate OCR).

#### **2.4.1.1.2 Findings:**

1. **PyPDF2** performs well for simple, text-based PDFs but has limitations when dealing with non-standard layouts or complex elements (e.g., tables, images).

#### **2.4.1.1.3 Gaps:**

1. Lack of support for OCR (Optical Character Recognition) makes it unsuitable for scanned PDFs.
2. Limited capabilities for handling structured data, such as tables.

#### **2.4.1.1.4 Contribution Potential:**

Your project could improve upon **PyPDF2** by integrating more advanced layout parsing techniques and adding OCR support to handle scanned PDFs and complex layouts more effectively.

### **2.4.1.2 PDFMiner**

**PDFMiner** is a more sophisticated Python library focused on extracting structured text from PDFs. It parses the entire document to reconstruct the layout and extract information in a more structured manner, including support for **layout-aware text extraction**.

#### 2.4.1.2.1 Key Features and Methodology:

1. Uses a low-level parsing approach, analyzing the content stream of each PDF page to retrieve text and layout information.
2. Supports extraction of structured elements such as fonts, locations, and positions of text on a page, which is particularly useful for documents with non-linear text.
3. Can handle embedded fonts, but struggles with tables and images without additional post-processing.

#### 2.4.1.2.2 Findings:

1. **PDFMiner** is powerful for extracting text with detailed information about layout and formatting.
2. It works well with text-heavy PDFs but is not optimized for scanned documents or PDFs with highly complex structures, such as tables or charts.

#### 2.4.1.2.3 Gaps:

1. No direct OCR support for handling image-based PDFs.
2. The performance can be slow when dealing with large PDFs due to its low-level parsing approach.

#### 2.4.1.2.4 Contribution Potential:

Your project could build on **PDFMiner** by enhancing its performance and integrating OCR capabilities (e.g., using **Tesseract** or **pdfplumber**) to address issues with scanned documents and improve its handling of complex layouts like tables.

#### 2.4.1.3 PDFPLUMBER

**pdfplumber** is a Python library designed specifically for extracting structured information from PDFs. It extends the capabilities of other libraries by focusing on extracting not only text but also tables and other elements in a more structured manner.

#### 2.4.1.3.1 Key Features and Methodology:

1. Integrates **OCR** (through Tesseract) for handling scanned PDFs, making it suitable for documents that require image-based text extraction.
2. Focuses on maintaining the layout integrity, especially for tables, allowing it to extract tabular data more accurately compared to other libraries.
3. Offers additional tools for post-processing, such as cropping pages and adjusting the DPI (dots per inch) for better text recognition.

#### 2.4.1.3.2 Findings:

1. **pdfplumber** excels in extracting text from PDFs with complex layouts, especially those containing tables and diagrams.
2. However, it can be slower in performance due to the added complexity of handling structured data and the OCR process.

#### 2.4.1.3.3 Gaps:

1. Limited support for semantic analysis of text (does not integrate NLP techniques).
2. While its table extraction is strong, other structural elements such as footnotes, sidebars, and images can be difficult to process.

#### 2.4.1.3.4 Contribution Potential:

Your project can contribute by enhancing **pdfplumber**'s capabilities through NLP integration (e.g., Named Entity Recognition, sentiment analysis) to provide meaningful, structured information from extracted text. This can help users gain deeper insights from complex documents.

#### 2.4.1.4 TESSERACT OCR WITH PYTHON (PYTESSERACT)

**Tesseract** is an OCR engine that can be integrated into Python using **Pytesseract**. It is widely used for extracting text from image-based PDFs, such as scanned documents. **Pytesseract** can convert images of text into machine-readable formats but often requires preprocessing steps like **image binarization** and **deskewing** for optimal results.

#### 2.4.1.4.1 Key Features and Methodology:

1. Handles scanned PDFs and images, making it essential for documents that cannot be processed by text-only extraction tools.
2. Requires image preprocessing (e.g., deskewing, noise reduction) to improve the accuracy of text recognition.
3. Works well with multiple languages, with support for training on custom fonts and languages.

#### 2.4.1.4.2 Findings:

1. **Tesseract** performs well for OCR tasks, especially for clean, high-quality scans.
2. Preprocessing is crucial to ensure accurate text recognition, as poor image quality can lead to high error rates.
3. Not well-suited for handling complex layouts or mixed-content PDFs (e.g., those that contain both text and images).

#### 2.4.1.4.3 Gaps:

1. Struggles with complex PDF layouts that include tables, diagrams, or non-standard fonts without additional post-processing.
2. Does not inherently provide any NLP capabilities or semantic understanding of the extracted text.

#### 2.4.1.4.4 Contribution Potential:

Your project can address the limitations of **Pytesseract** by incorporating advanced preprocessing techniques (e.g., deskewing, binarization) and augmenting OCR results with NLP methods (e.g., for entity extraction, sentiment analysis, or topic modeling). This would improve the quality of extracted information from both scanned and complex PDFs.

#### 2.4.1.5 POPPLER WITH PYTHON (PDFTOTEXT)

**Poppler** is a PDF rendering library that provides text extraction capabilities through its command-line tool, **pdftotext**. It is often used for quick and lightweight text extraction but lacks the ability to handle complex structures.

#### 2.4.1.5.1 Key Features and Methodology:

1. Focuses on plain-text extraction with minimal formatting or structural analysis.
2. Fast and efficient for simple PDFs, making it a popular choice for large-scale batch processing.
3. Minimal support for non-text elements (e.g., tables, images).

#### 2.4.1.5.2 Findings:

1. **Poppler** works well for basic text extraction but struggles with more complex PDFs, especially those with non-linear structures or embedded objects.
2. It cannot extract structured information like tables or figures.

#### 2.4.1.5.3 Gaps:

1. Lacks the ability to handle structured elements or incorporate OCR for scanned PDFs.
2. Minimal support for integrating advanced NLP techniques.

#### 2.4.1.5.4 Contribution Potential:

The project could extend the capabilities of **Poppler** by integrating NLP techniques for post-processing extracted text, allowing for better analysis of the extracted information. Additionally, combining it with OCR tools can improve text extraction from scanned PDFs.

### 2.4.2 IDENTIFYING GAPS IN EXISTING RESEARCH AND PROJECTS

From the review of existing projects and research, several gaps have been identified in the area of PDF text extraction:

1. **Limited Integration of OCR and NLP:** Most tools either focus on text extraction or OCR but fail to integrate NLP techniques (e.g., NER, sentiment analysis) to add meaningful context to the extracted text.
2. **Inability to Handle Complex Layouts:** Many libraries, such as **PyPDF2** and **Poppler**, struggle with complex PDFs that include tables, images, and non-linear text structures.

- 3. Slow Performance with Large Documents:** Libraries like **PDFMiner** and **pdfplumber** can be slow when processing large documents, especially those with complex formatting.
- 4. Preprocessing Dependence:** OCR tools like **Pytesseract** heavily rely on preprocessing steps, which can be difficult to automate and optimize for all document types.

### **2.4.3 HOW THIS PROJECT WILL CONTRIBUTE TO THE FIELD**

Your project has the potential to contribute significantly by addressing the aforementioned gaps:

- 1. OCR + NLP Integration:** By combining OCR tools with advanced NLP techniques (e.g., entity extraction, topic modeling), your project can provide not only text extraction but also meaningful, structured information that can be used for downstream tasks like document categorization, summarization, and data analysis.
- 2. Improved Layout Parsing:** Focusing on improving the handling of complex PDF layouts (e.g., tables, images) can fill a major gap in current projects. Integrating layout-aware techniques from **pdfplumber** with performance optimizations can make the process more efficient and accurate.
- 3. Enhanced Preprocessing for OCR:** By automating and optimizing preprocessing techniques (e.g., binarization, deskewing), your project can significantly improve the accuracy of OCR for scanned documents.
- 4. Better Performance and Scalability:** Addressing the performance issues of existing libraries can make text extraction more scalable and suitable for large document collections.

This project will thus contribute to the field by developing a more comprehensive, robust, and scalable solution for PDF text extraction using Python.

## CHAPTER 3: SYSTEM DESIGN AND METHODOLOGY

### 3.1 INTRODUCTION

The methodology adopted for this project follows a modular, pipeline-based architecture designed to handle the complexity and diversity of PDF files. It ensures accurate and scalable text extraction from both simple and complex documents using Python and its ecosystem of libraries.

#### 3.1.1 MODULAR TEXT EXTRACTION PIPELINE METHODOLOGY FOR PDF DOCUMENT PROCESSING

The **Modular Text Extraction Pipeline Methodology** is a structured, stage-by-stage approach used to extract meaningful text from PDF documents. This methodology is particularly suitable for handling the diverse and often complex nature of PDFs, including both text-based and image-based files.

##### 3.1.1.1 CORE CONCEPT

The methodology divides the entire text extraction process into **independent, reusable, and well-defined modules**, each responsible for a specific task in the pipeline. This design improves **accuracy, flexibility, scalability, and maintainability**, and allows easy integration of new techniques or tools as needed.

##### 3.1.1.2 KEY MODULES OF THE PIPELINE

#### 1. Input Handling

- a. Determines the type of PDF (text-based or image-based)
- b. Handles encrypted files and multi-page documents

#### 2. Preprocessing

For image-based PDFs, preprocessing includes:

- a. **Binarization** (converting images to black and white)
- b. **Deskewing** (aligning tilted images)
- c. **Noise reduction** (removing visual artifacts)

#### 3. Text Extraction

- a. **Text PDFs**: Use parsing libraries like pdfminer.six, pdfplumber
- b. **Image PDFs**: Convert to images and use OCR (Tesseract, pytesseract) for text recognition

#### 4. Post-Processing

- a. Clean up extracted text (remove special characters, line breaks)
- b. Restore document structure (e.g., paragraphs, lists, tables)

#### 5. Evaluation

- a. Accuracy of text extraction is measured using:

##### 1. Precision, Recall, F1-Score

##### 2. Character Error Rate (CER)

- b. Ground truth data is used for validation

##### 3. Output and Export

- a. Final structured text is saved in formats such as TXT, CSV JSON or SOUND
- b. Can be fed into further NLP tasks or data analytics

### 3.1.1.3 WHY MODULAR?

- 1. **Reusability:** Modules can be reused or replaced independently.
- 2. **Customization:** Easily adapt the pipeline to new PDF types or formats.
- 3. **Error Isolation:** Bugs can be traced to specific components.
- 4. **Extensibility:** New technologies (e.g., advanced OCR engines, AI models) can be plugged in without redesigning the whole system.

### 3.1.1.4 BENEFITS FOR PDF DOCUMENT PROCESSING

- 1. Handles heterogeneous document types
- 2. Enables **automation** of large-scale document workflows
- 3. Facilitates data mining, text analytics, and knowledge extraction
- 4. Supports integration with **Natural Language Processing (NLP)** systems



### **3.1.2 OVERVIEW OF METHODOLOGY STEPS**

#### **3.1.2.1 PROBLEM IDENTIFICATION & REQUIREMENT ANALYSIS**

1. Identified the difficulty of extracting text from complex and varied PDFs manually.
2. Outlined project requirements: support for multiple PDF types, accurate text extraction, and modular system design.

#### **3.1.2.2 DATASET COLLECTION**

1. Compiled a diverse set of PDF documents:
  - a. Text-based PDFs (e.g., e-books, academic articles)
  - b. Image-based PDFs (e.g., scanned forms)
  - c. Encrypted PDFs (with known passwords)
  - d. PDFs with tables and complex layouts

#### **3.1.2.3 PREPROCESSING**

1. **Text PDFs:** Checked for readability and decrypted if necessary.
2. **Image PDFs:** Applied image preprocessing techniques:
  - a. Binarization
  - b. Deskewing
  - c. Denoising
  - d. Resizing for OCR

#### **3.1.2.4 TEXT EXTRACTION**

1. **Text PDFs:** Used pdfminer.six, PyPDF2, and pdfplumber to parse text directly.
2. **Image PDFs:** Used pdf2image to convert pages into images and pytesseract (Tesseract OCR) to extract text.

### **3.1.2.5 POST-PROCESSING**

1. Removed artifacts such as headers, footers, and special characters.
2. Normalized the text format (paragraph breaks, spacing).
3. Identified structural elements (e.g., bullet points, tables).

### **3.1.2.6 EVALUATION**

1. Measured extraction accuracy using:
  - a. Precision
  - b. Recall
  - c. F1-score
  - d. Character Error Rate (CER)
2. Compared against manually extracted ground-truth data.

### **3.1.2.7 TESTING AND VALIDATION**

1. Conducted:
  - a. Unit testing (each module separately)
  - b. Integration testing (across modules)
  - c. Performance testing (on large documents)

### **3.1.2.8 DOCUMENTATION AND DEPLOYMENT**

1. Created a user manual and inline code documentation.
2. Packaged scripts for easy reuse and scalability.

## **3.2 SYSTEM REQUIREMENT SPECIFICATION**

The business rules of the **PDF Text Extraction System** define how the system will operate, ensuring consistency, accuracy, and adherence to the system's goals. These rules guide the behavior of different components of the system, from handling documents to extracting and processing data.

### **3.2.1 SYSTEM REQUIREMENT FOR THE PDF TEXT EXTRACTION SYSTEM**

#### **3.2.1.1 DOCUMENT INGESTION RULES**

1. Only PDF files are accepted for processing. Other file formats are not allowed.
2. Each document must have a unique file name within the system to avoid duplication.
3. Document metadata (e.g., file name, extraction date, document type) must be stored upon ingestion for future retrieval.
4. Document type is automatically categorized as either "text-based" or "scanned (OCR)" based on the file content.
5. Large PDFs (exceeding a certain size limit) are split into smaller, manageable sections for processing to prevent performance issues.

#### **3.2.1.2 TEXT EXTRACTION RULES**

1. For text-based PDFs, direct parsing (e.g., using libraries like pdfplumber, pdfminer) is used to extract text.
2. For scanned PDFs or image-heavy PDFs, Optical Character Recognition (OCR) is applied using engines like Tesseract.
3. Text extraction should preserve document structure (e.g., paragraphs, tables, bullet points) as closely as possible to maintain readability and usability.
4. Errors during extraction (e.g., unreadable text) are logged and flagged for review.

#### **3.2.1.3 PREPROCESSING RULES**

1. Scanned PDF files undergo preprocessing (e.g., binarization, deskewing) to improve OCR accuracy.
2. Image preprocessing must be applied to all scanned PDFs before passing them to the OCR engine.
3. If preprocessing fails to enhance the image quality, an error message is logged, and the extraction process continues but is flagged for low confidence.

### 3.2.1.4 NLP AND POSTPROCESSING RULES

1. After text extraction, Natural Language Processing (NLP) techniques, such as **Named Entity Recognition (NER)** and **sentiment analysis**, must be applied to identify key entities (e.g., names, dates, locations) and determine the document's tone.
2. Extracted text must be cleaned (removing unnecessary symbols, correcting character encoding issues) before applying NLP techniques.
3. All recognized entities must be stored in the Entities table along with their corresponding document.
4. Sentiment analysis must be applied to sections of the document to gauge the tone, with results stored in the Sentiment table.

### 3.2.1.5 DATA STORAGE AND OUTPUT RULES

1. All extracted text, including metadata, named entities, and sentiment analysis results, must be stored in the respective database tables.
2. Extracted text must be linked to the original document via the Document ID to maintain traceability.
3. Output data must be in a structured format (e.g., JSON, CSV, or plain text) and accessible to other systems if needed.
4. Users should be able to retrieve documents, entities, or sentiment results using search or filter options (e.g., by document name, entity type, or sentiment score).

### 3.2.1.6 ACCURACY AND QUALITY CONTROL RULES

1. If the confidence level of text extraction or NLP results falls below a certain threshold (e.g., OCR confidence score or NLP accuracy), the document is flagged for manual review.
2. Every extracted text entry must be verified for completeness (e.g., full text extracted, key data points present).
3. The system must perform periodic evaluations to ensure the accuracy and quality of the text extraction and NLP processes.

### 3.2.1.7 ERROR HANDLING AND LOGGING RULES

1. Any system errors (e.g., issues with text extraction, preprocessing, OCR) must be logged with clear error messages for troubleshooting.
2. The system should be able to handle incomplete or corrupted PDF files, logging the errors and skipping to the next document without halting the entire process.
3. Users must be notified of extraction errors or flagged documents for review and correction.

### 3.2.1.8 SECURITY AND ACCESS CONTROL RULES

1. Only authorized users can upload, access, and retrieve documents and extracted data.
2. Sensitive information, such as names or personal details extracted from documents, must be handled with privacy protections (e.g., encryption or masking when necessary).
3. User activity (e.g., document uploads, data retrievals) must be logged for audit purposes.

### 3.2.1.9 PERFORMANCE AND SCALABILITY RULES

1. The system must be able to handle large batches of documents for extraction, ensuring that performance does not degrade under heavy load.
2. The text extraction process must be scalable, allowing multiple documents to be processed concurrently or in parallel.

### 3.2.1.10 DOCUMENT INTEGRITY RULES

1. The original structure of the PDF document (e.g., page order, headers, tables, and lists) must be maintained as much as possible in the extracted text.
2. Any extracted data that cannot be linked back to its original document (e.g., if extraction failed) must not be stored in the system.

## 3.3 SYSTEM ARCHITECTURE

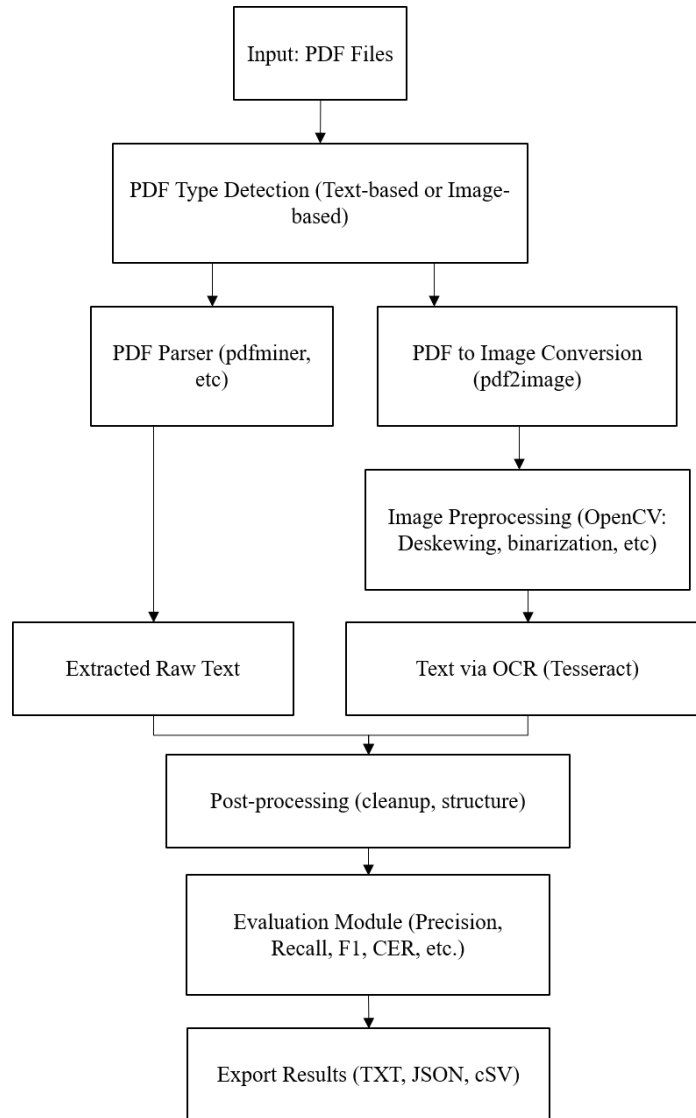
The architecture follows a **modular pipeline structure** with clear separation of concerns. Each module performs a specific function, allowing for easier debugging, testing, and future expansion.

### 3.3.1 SYSTEM WORKFLOW (TEXTUAL DIAGRAM)

A **textual diagram** is a way of representing a system or process **using only text characters** (no images, graphics, or drawing tools). It's often used in:

1. Early planning stages
2. Text-based documents (e.g., markdown, plain text)
3. Command-line interfaces or documentation
4. Situations where visual tools aren't available

**Figure 3.3.1: Textual Diagram for PDF Text Extraction System**



### 3.3.2 CORE COMPONENTS

**Table 3.3.1: Core Components of the PDF Text Extraction System**

<b>Module</b>	<b>Description</b>
<b>Input Handler</b>	Accepts and manages various PDF formats, including encrypted or multi-page
<b>Detection Engine</b>	Determines whether the document is text-based or image-based
<b>PDF Parsing</b>	Extracts text from text-based PDFs using libraries like pdfminer, PyPDF2
<b>OCR Pipeline</b>	Converts images to text using pytesseract, after preprocessing with OpenCV
<b>Post-Processor</b>	Cleans and formats the extracted text (removes noise, restores structure)
<b>Evaluator</b>	Measures extraction quality using standard NLP metrics
<b>Exporter</b>	Outputs clean text to formats such as TXT, JSON, CSV



### 3.3.3 TECHNOLOGIES USED

**Table 3.3.2: Technologies Used for the PDF Text Extraction System**

Component	Tool/Library Used
PDF Parsing	pdfminer.six, pdfplumber
OCR	Tesseract OCR, pytesseract
Preprocessing	OpenCV, Pillow
Evaluation	scikit-learn, custom scripts
Visualization (optional)	matplotlib, seaborn
Development Language	Python

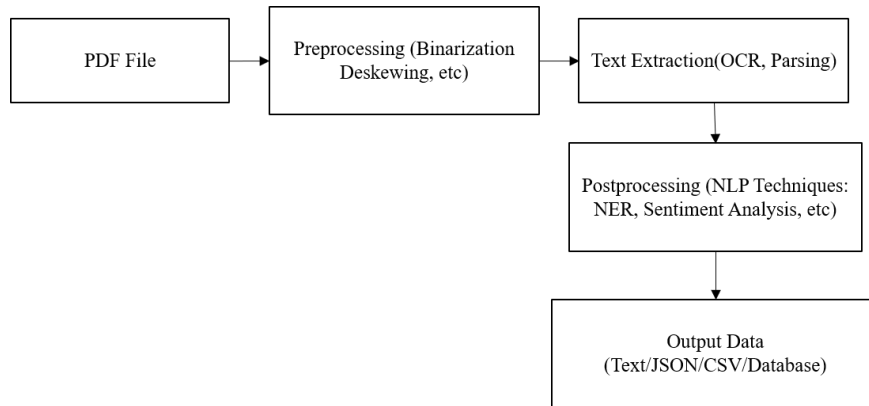
### 3.4 DATA FLOW DIAGRAM (DFD)

At a high level, the system follows this flow of information:

1. **Input:** PDF files (can be scanned documents, images, or text-based PDFs).
2. **Preprocessing:**
  - a. For scanned PDFs or images, image preprocessing techniques are applied (e.g., binarization, deskewing).
  - b. Text-based PDFs are sent directly for text extraction.
3. **Text Extraction:**
  - a. **OCR** is applied for scanned PDFs (using **Tesseract**).
  - b. For non-scanned PDFs, libraries such as **pdfplumber** or **PDFMiner** are used to extract text while preserving the structure (e.g., tables, formatting).
4. **Postprocessing:**
  - a. The extracted text is cleaned and normalized (e.g., tokenization, lemmatization).
  - b. **NLP Techniques** (e.g., Named Entity Recognition, Sentiment Analysis) are applied to extract meaningful information and improve accuracy.
5. **Output:** Clean and structured text is exported to desired formats (e.g., JSON, CSV, plain text).

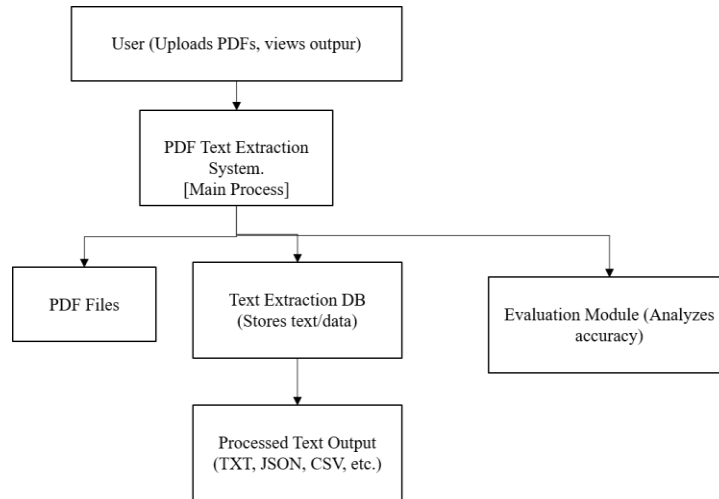
This process can be visualized as:

**Figure 3.4.1: Visualization of the Flow of Data**



### 3.4.1 CONTEXT DIAGRAM (LEVEL 0)

**Figure 3.4.2: Context Diagram of the PDF Text Extraction System**

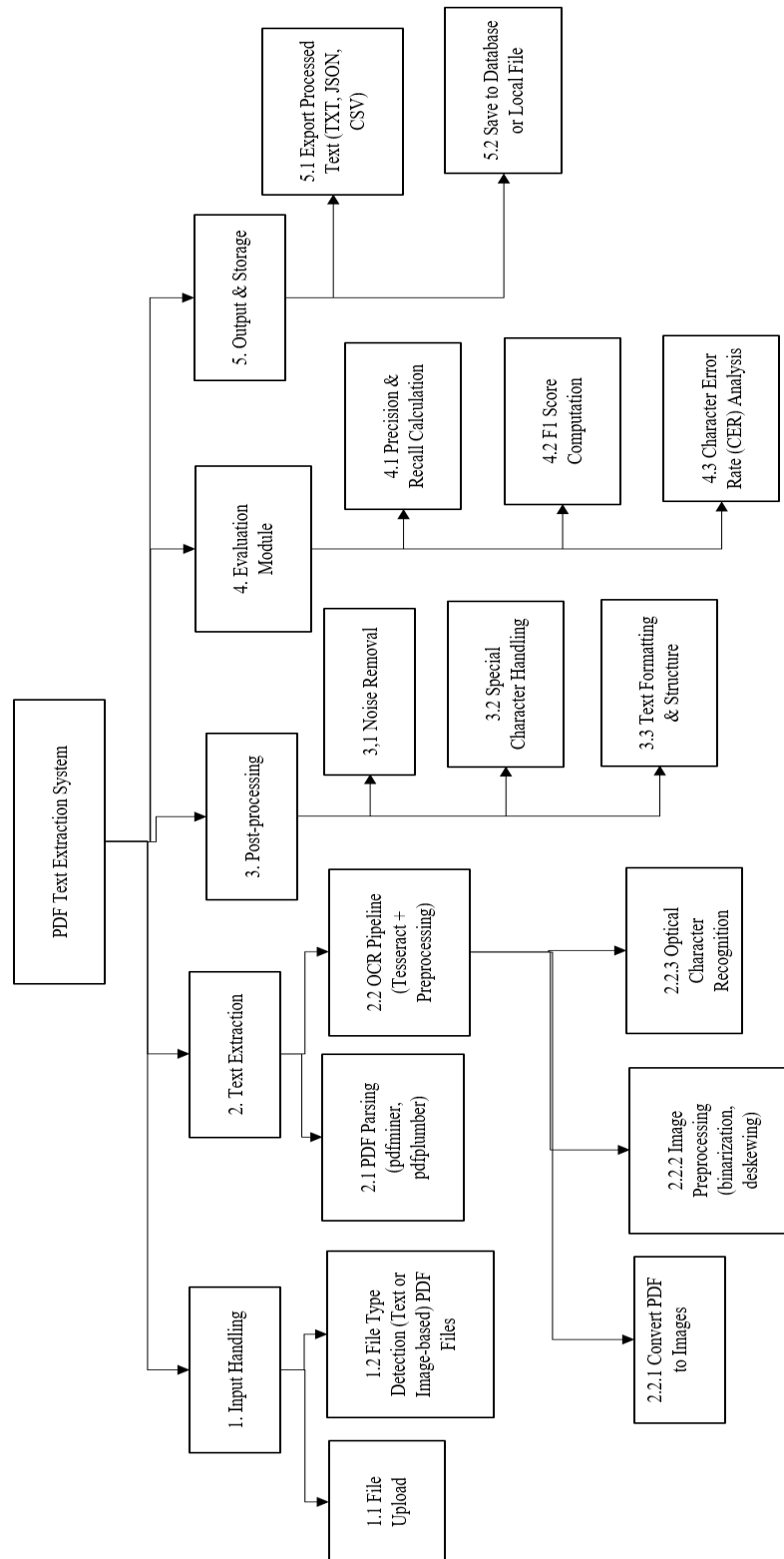


## **External Entities**

- 1. User:** Uploads PDFs and retrieves the extracted data.
- 2. PDF Files:** The raw input documents.
- 3. Text Extraction DB:** Stores the output for further usage.
- 4. Evaluation Module:** Used to evaluate the quality of the output text.

### 3.4.2 STRUCTURED DIAGRAM (LEVEL 1)

Figure 3.4.3: Structured Diagram of the PDF Text Extraction System



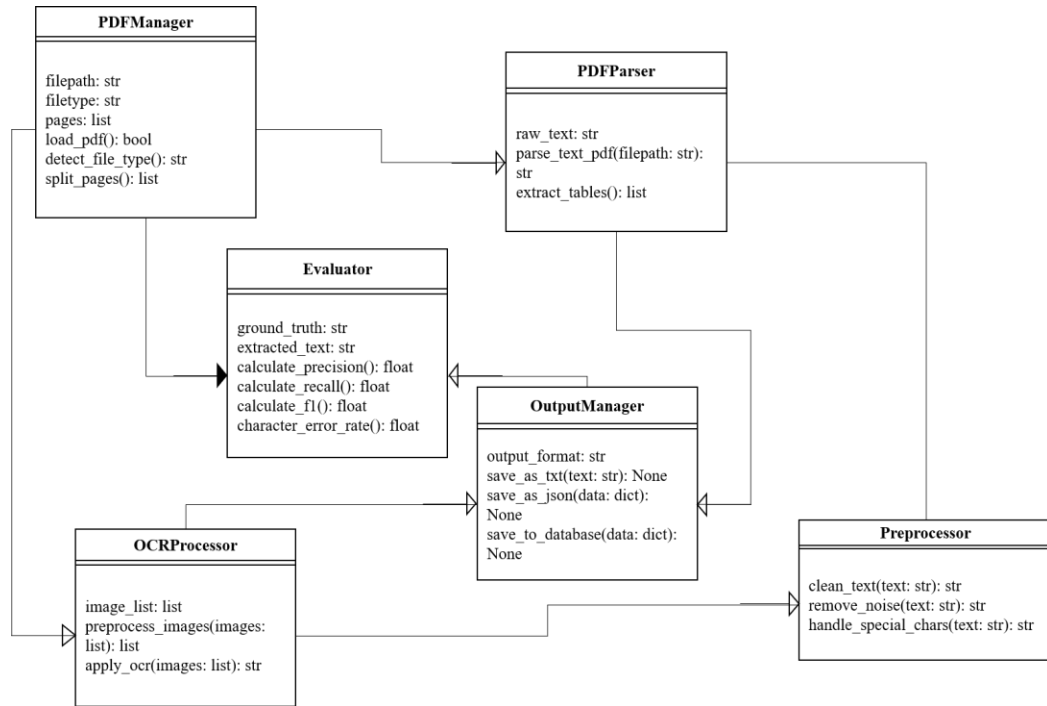
## 3.5 UML DIAGRAMS

### 3.5.1 CLASS DIAGRAM

The UML **class diagram** provides a structural view of the components and their relationships.

1. **PDFParser Class:** Handles the file I/O, managing the input and output of PDF files.
2. **Preprocessor Class:** Handles image preprocessing tasks (e.g., binarization, deskewing).
3. **TextExtractor Class:** Responsible for extracting text, either using OCR (for scanned PDFs) or direct parsing (for text-based PDFs).
4. **NLPProcessor Class:** Handles text processing, including tokenization, Named Entity Recognition (NER), and sentiment analysis.
5. **PostProcessor Class:** Cleans and organizes extracted text for storage or further use.
6. **DatabaseHandler Class:** Manages the storage of extracted text in a database (optional).

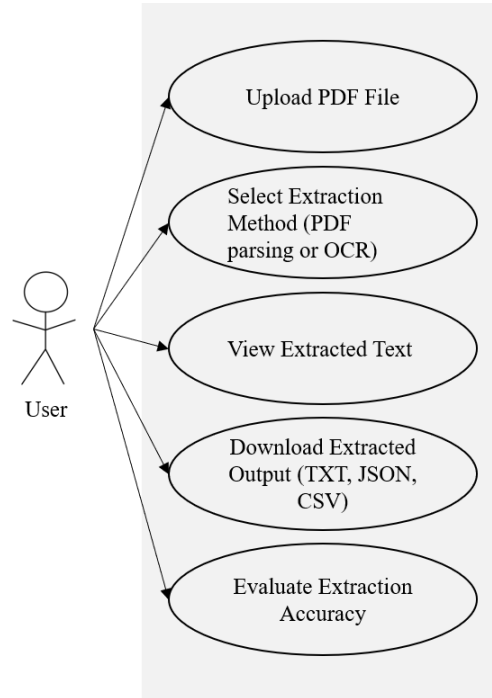
**Figure 3.5.1: Class Diagram**





### 3.5.2 USE CASE

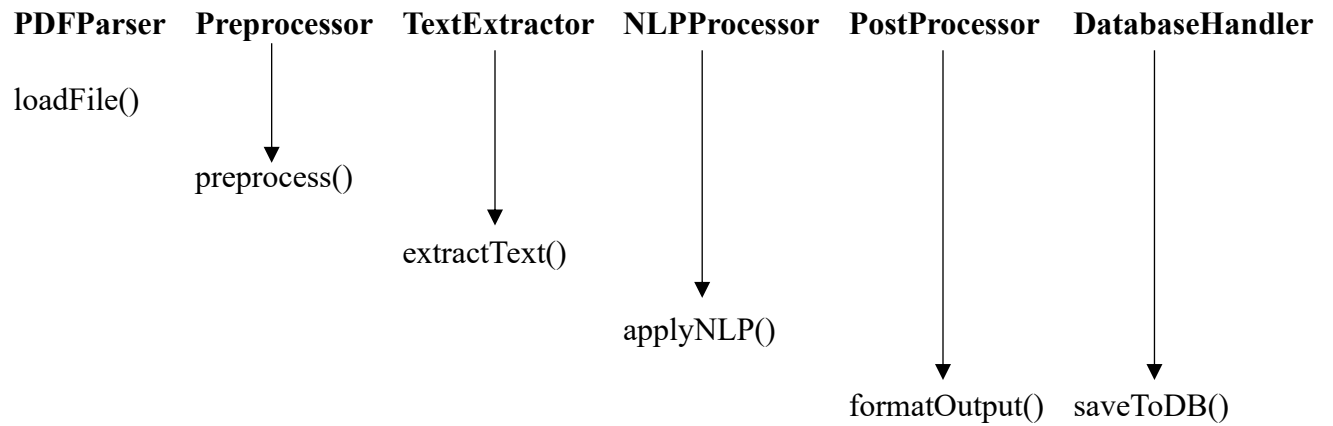
**Figure 3.5.2: Use Case Diagram**



### 3.5.3 SEQUENCE DIAGRAM

A **sequence diagram** shows the interactions between objects over time, capturing the flow of information from the initial input of the PDF to the final output.

**Figure 3.5.3: Sequence Diagram**



### 3.5.4 DATA DICTIONARY

Table 3.5.1: Data Dictionary

Field Name	Data Type	Description	Format / Example
file_id	Integer	Unique identifier for each uploaded PDF file	101
file_name	String	Original name of the uploaded PDF	biology_notes.pdf
file_path	String	File location on disk/server	/uploads/biology_notes.pdf
file_type	String	Detected type: 'text' or 'image' based PDF	"text" / "image"
upload_time	Datetime	Timestamp when the file was uploaded	2025-04-23 10:12:55
extraction_method	String	Method used for extraction (pdfparser or ocr)	"pdfparser" / "ocr"
raw_text	Text	Unprocessed extracted text	"Chapter 1: Introduction to Biology..."
cleaned_text	Text	Text after preprocessing and noise removal	"Chapter 1 Introduction to Biology..."
image_list	Array	List of image paths if converted from image-based PDFs	["page1.jpg", "page2.jpg"]
precision	Float	Precision score of extracted text evaluation	0.89
recall	Float	Recall score of extracted text evaluation	0.91

f1_score	Float	F1-Score calculated from precision and recall	0.90
character_error_rate	Float	Character Error Rate (CER) for OCR performance	0.03
output_format	String	Format chosen by the user for saving the output	"TXT", "JSON", "CSV"
saved_location	String	File path where output is stored	/results/output_101.txt
user_id	Integer	Optional: ID of the user uploading the file (if multi-user support)	2001
log_message	String	Message for logging errors or processes	"OCR process started for file biology_notes.pdf"

### 3.6 FLOWCHART

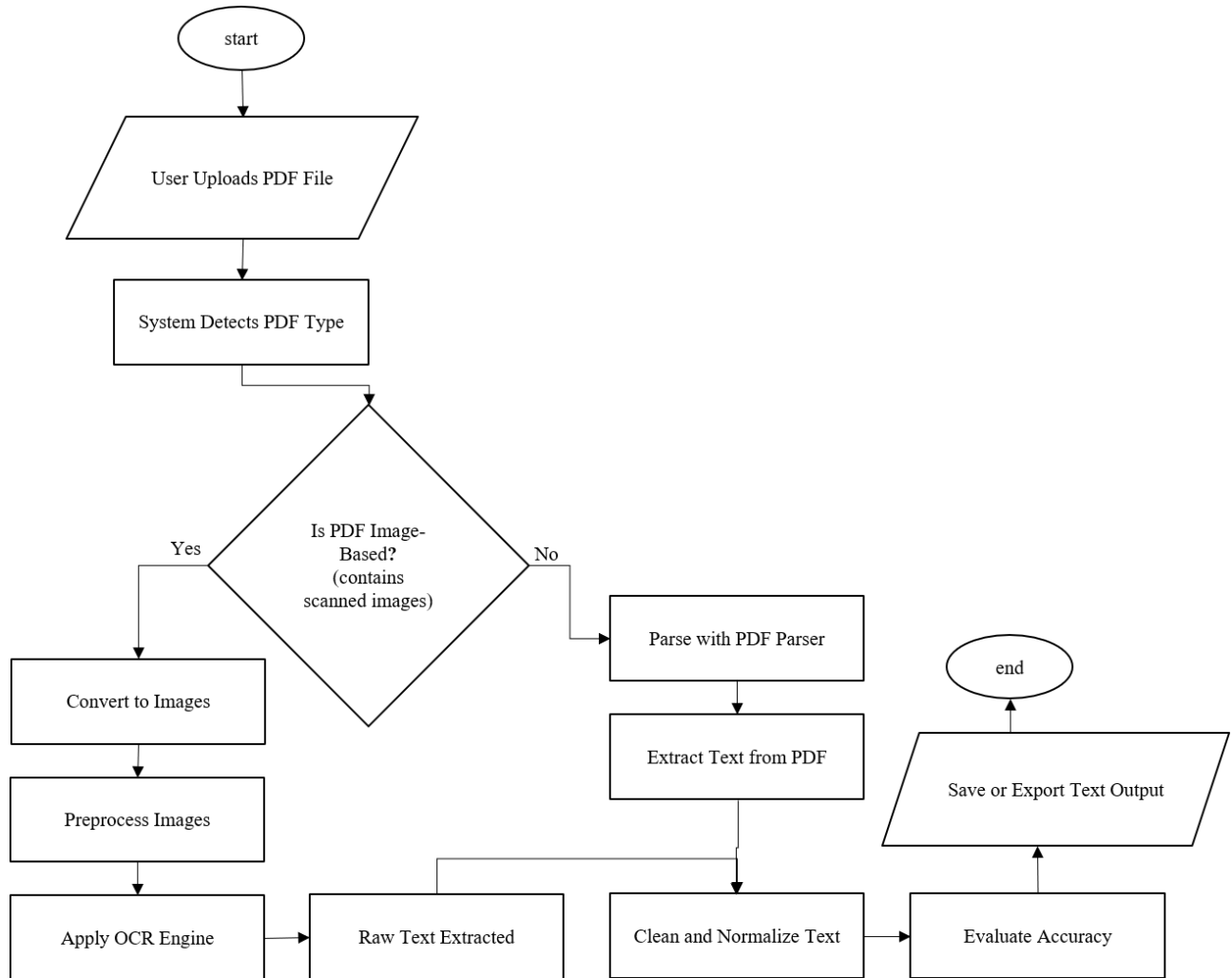
This flowchart provides an overview of the step-by-step process followed by the system.

**Pseudocode:**

- 1. User Uploads PDF File:** Interface where the user selects a PDF.
- 2. System Detects Type:** Checks whether the PDF is text-based or image-based.
- 3. OCR Path:** Converts scanned PDFs to images and applies preprocessing + OCR.
- 4. Parser Path:** Uses libraries like pdfminer or PyPDF2 to extract embedded text.
- 5. Cleaning & Evaluation:** Applies NLP and formatting, optionally compares with ground truth.
- 6. Output:** Saves the extracted text in a chosen format (TXT, JSON, etc.).

## Flowchart:

Figure 3.6.1: Flow Chart



## 3.7 DATABASE MODEL

### 3.7.1 ER-MODEL

The Entity-Relationship (E-R) model for a PDF text extraction system represents the entities, their attributes, and the relationships between them visually. Below is a conceptual representation of how the entities and their relationships could be structured.

#### Entities and Attributes:

##### 1. PDF File

- a. pdf\_id (Primary Key)
- b. user\_id (Foreign Key)
- c. file\_name
- d. file\_path
- e. upload\_date
- f. file\_size
- g. file\_type

##### 2. Text Extraction

- a. extraction\_id (Primary Key)
- b. pdf\_id (Foreign Key from PDF File)
- c. extracted\_text
- d. extraction\_date

##### 3. Extracted Section

- a. section\_id (Primary Key)
- b. extraction\_id (Foreign Key from Text Extraction)
- c. section\_name
- d. start\_position
- e. end\_position

f. section\_text

#### 4. Error Log

- a. error\_id (Primary Key)
- b. extraction\_id (Foreign Key from Text Extraction)
- c. error\_type
- d. error\_message
- e. timestamp

#### 5. User

- a. user\_id (Primary Key)
- b. username
- c. email
- d. password\_hash

### Relationships:

#### 1. PDF File to Text Extraction: One-to-Many (One PDF can have multiple extractions)

- a. Relationship Name: Contains
- b. Cardinality: 1-to-many (1 PDF file can generate multiple extractions)

#### 2. Text Extraction to Extracted Section: One-to-Many (One extraction can have multiple sections)

- a. Relationship Name: Has
- b. Cardinality: 1-to-many (1 extraction can have multiple extracted sections)

#### 3. Text Extraction to Error Log: One-to-Many (One extraction can generate multiple error logs)

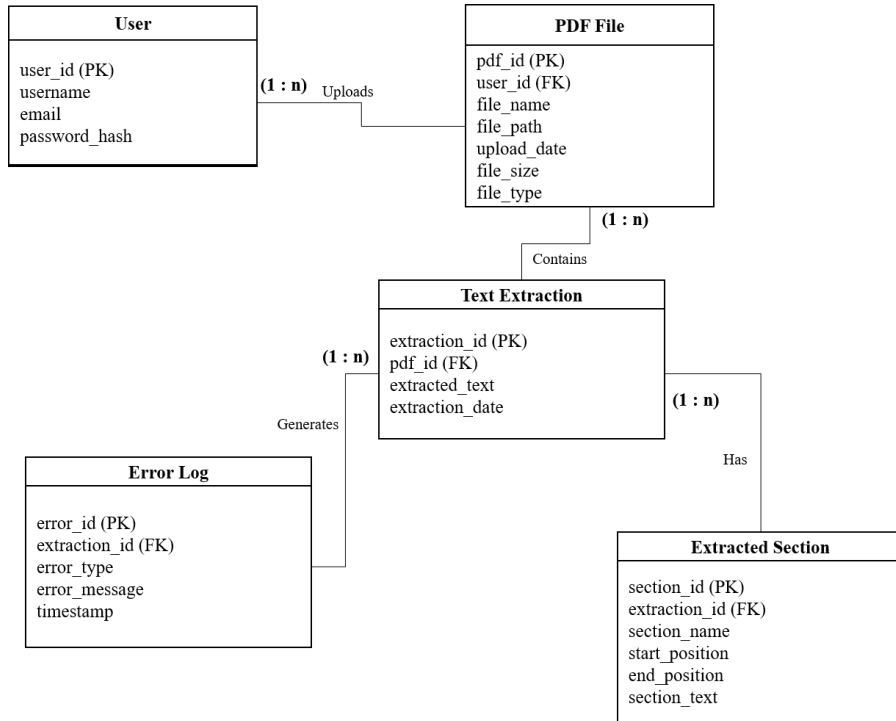
- a. Relationship Name: Generates
- b. Cardinality: 1-to-many (1 extraction may have multiple errors)



**4. User to PDF File: One-to-Many (One user can upload multiple PDFs)**

- a. Relationship Name: Uploads
- b. Cardinality: 1-to-many (1 user can upload many PDFs)

**Figure 3.7.1: ER-Model**



### 3.7.2 RELATIONAL MODEL

**1. PDF Files Table:** Stores metadata about each uploaded PDF.

PDF\_Files(pdf\_id, user\_id (FK), file\_name, file\_path, upload\_date, file\_size, file\_type)

**2. Text Extraction Table:** Stores the text extracted from each PDF file.

Text\_Extraction(extraction\_id, pdf\_id (FK), extracted\_text, extraction\_date)

**3. Extracted Sections Table:** Stores details of sections extracted from each PDF file.

Extracted\_Sections(section\_id, extraction\_id (FK), section\_name, start\_position, end\_position, section\_text)

**4. Error Logs Table:** Stores logs if any errors occur during text extraction.

Error\_Logs(error\_id, extraction\_id (FK), error\_type, error\_message, timestamp)

**5. Users Table:** Stores user details if applicable. Users(user\_id, username, email, password\_hash)

### 3.7.3 SAMPLE QUERIES

#### 1. User Management Queries

These queries handle user registration and login.

- a. **Create a New User (Registration)** This query inserts a new user's details into the Users table. The password would be hashed by the application before being stored.

```
INSERT INTO Users (username, email, password_hash)
VALUES ('john_doe', 'john.doe@example.com', 'a1b2c3d4e5f6...');
```

- b. **Verify a User (Login)** This query retrieves a user's ID and hashed password to verify their credentials during login.

```
SELECT user_id, password_hash FROM Users WHERE email = 'john.doe@example.com';
```

#### 2. PDF File and Extraction Queries

These queries manage the core workflow of uploading files and saving the extracted text.

- a. **Log a New PDF Upload** When a user uploads a file, this query saves its metadata to the

PDF\_Files table, linking it to the user.

```
INSERT INTO PDF_Files (user_id, file_name, file_path, upload_date, file_size, file_type)
VALUES (1, 'annual_report_2024.pdf', '/uploads/annual_report_2024.pdf', '2025-08-23
10:30:00', 5120, 'text');
```

- b. **Save Extracted Text** After processing a PDF, this query stores the complete extracted text, linking it to the original PDF file record.

```
INSERT INTO Text_Extraction (pdf_id, extracted_text, extraction_date)
VALUES (101, 'This is the full extracted text from the PDF...', '2025-08-23 10:31:00');
```

- c. **Retrieve a User's Extraction History** This query fetches all the files a specific user has processed, showing the most recent ones first. This is used to display the user's history in the app.

```
SELECT
p.file_name,
t.extraction_date,
SUBSTR(t.extracted_text, 1, 100) AS text_preview
FROM
PDF_Files p
JOIN
Text_Extraction t ON p.pdf_id = t.pdf_id
WHERE
p.user_id = 1
ORDER BY
.extraction_date DESC;
```

### 3. Section and Error Logging Queries

These queries are for more advanced features, like handling specific sections or logging issues.

- a. **Save an Extracted Table (Section)** If the system identifies and extracts a specific table from a PDF, this query saves it as a distinct "section."

```
INSERT INTO Extracted_Sections (extraction_id, section_name, start_position,
end_position, section_text
VALUES (201, 'Q2 Financials Table', 5430, 6800, '<CSV data of the table>');
```

- b. **Log an Error During Extraction** If the OCR process fails for a specific PDF, this query logs the error for troubleshooting.

```
INSERT INTO Error_Logs (extraction_id, error_type, error_message, timestamp)
VALUES (202, 'OCR_FAILURE', 'Tesseract failed with a timeout on page 3.', '2025-08-23
11:00:00');
```

### 3.8 KEY COMPONENTS OF THE SYSTEM

The **PDF Text Extraction System** consists of several key components that work together to ensure the accurate extraction, processing, and analysis of text from PDF files. These components include **data preprocessing**, **text extraction**, **post-processing**, and **evaluation**. Below is a detailed description of each component:

#### 3.8.1 DATA PREPROCESSING

The data preprocessing phase prepares the PDF files for text extraction. It plays a critical role in improving the quality of input data, particularly when dealing with scanned documents or image-heavy PDFs. The preprocessing techniques vary based on whether the PDF is text-based or image-based.

##### 3.8.1.1 THE STEPS INVOLVED IN PREPARING PDF FILES FOR TEXT EXTRACTION

###### 1. File Validation

###### Objective:

Ensure the PDF is valid and not corrupted.

###### Steps:

- a. **File Format Check:** Verify that the file is indeed a PDF and meets the standard PDF format requirements.
- b. **Corruption Detection:** Check for any file corruption or unreadable sections using validation tools or libraries like PyPDF2. If corrupted, the system should attempt recovery or reject the file.
- c. **Encryption Detection:** Determine if the PDF is password-protected or encrypted. If so, the system either prompts for the password or rejects the file.

**Considerations:**

- a. If a PDF is encrypted, the system must decrypt the file using the correct password or a decryption method. Failure to decrypt results in extraction errors.

**2. Image Handling (For Scanned PDFs or Image-Heavy PDFs)****Objective:**

Preprocess images within PDFs to improve Optical Character Recognition (OCR) accuracy.

**Steps:**

- a. **Image Detection:** Identify if the PDF is image-based (scanned document) or contains a significant number of images. This is typically done by checking the content streams of each page.
- b. **Binarization:** Convert the images to black-and-white (binary) to enhance the contrast between text and background, making OCR more accurate.
- c. **Deskewing:** If the images or text appear tilted or misaligned, deskewing techniques are applied to correct the orientation. This ensures that the OCR engine can properly recognize the characters.
- d. **Noise Removal:** Eliminate background noise or artifacts (e.g., shadows, smudges) that may interfere with OCR recognition. Filters and noise reduction techniques help clean up the image.
- e. **Resolution Adjustment:** If the image resolution is too low for OCR, the system adjusts or enhances the resolution. Higher resolution often leads to better OCR results but may increase processing time.

**Considerations:**

- a. Preprocessing is vital for scanned PDFs, as raw image data often contains noise or skew that degrades the OCR results.

**3. Handling Encrypted or Password-Protected PDFs****Objective:**

Gain access to protected content within encrypted PDFs.

**Steps:**

- a. **Password Detection:** Detect whether a PDF is encrypted or password-protected.
- b. **Password Prompt/Decryption:** If encrypted, the system prompts the user for the password or attempts to decrypt using available credentials.
- c. **Decryption Validation:** After decrypting, verify that the file is intact and ready for text extraction.

**Considerations:**

- a. If the correct password cannot be provided or if decryption fails, the system skips the file and logs an error.

#### 4. Multi-Page Document Handling and Splitting

**Objective:**

Prepare large multi-page PDFs for extraction by splitting them if necessary.

**Steps:**

- a. **Page Counting:** Identify the number of pages in the PDF. Some PDFs contain hundreds or even thousands of pages, which may be challenging to process in one go.
- b. **Document Splitting:** For large PDFs, the system splits the document into smaller sections or page ranges (e.g., every 50 pages). This allows for faster processing and avoids memory limitations.
- c. **Sequential Processing:** Process each section of the split PDF sequentially, ensuring that the text extraction results from each part are eventually merged into a coherent output.

**Considerations:**

- a. Splitting PDFs helps with performance and scalability, especially when dealing with large volumes of documents.

#### 5. PDF Structure Analysis

**Objective:**

Understand the internal structure of the PDF to determine the best method for extraction.

**Steps:**

- a. **PDF Structure Identification:** Analyze the structure of the PDF, including objects like text streams, images, fonts, and layout information. PDF files consist of objects and streams that can be organized in complex ways.
- b. **Content Differentiation:** Differentiate between text layers, images, and graphical elements. This is essential to ensure that only relevant text is extracted, and unnecessary images or graphical elements are ignored.
- c. **Table and Form Detection:** Identify tables, forms, or structured content (e.g., lists) that need special handling to maintain their structure during extraction.

**Considerations:**

- a. Some PDFs have complex structures, such as embedded tables or interactive forms, which may require custom handling or specialized libraries.

## 6. Metadata Extraction

**Objective:**

Extract useful metadata from the PDF file, such as title, author, creation date, and keywords.

**Steps:**

- a. **Metadata Parsing:** Use libraries (e.g., PyPDF2 or pdfplumber) to extract embedded metadata fields from the PDF. This data may include the document title, author, subject, and creation/modification dates.
- b. **Embedded Text Extraction:** In cases where the document contains text in embedded or hidden layers (such as annotations), the system attempts to extract this as part of the metadata.

**Considerations:**

- a. Metadata is useful for organizing and cataloging documents, particularly in large-scale applications such as document management systems.

## 7. Content Normalization



**Objective:**

Ensure consistency in the extracted text format.

**Steps:**

- a. **Text Format Standardization:** Normalize extracted text by removing unnecessary whitespace, special characters, and encoding issues.
- b. **Font and Encoding Handling:** PDFs may use different fonts or encoding schemes that could affect text extraction (e.g., non-standard characters). The system ensures all extracted text is converted into a readable and uniform format.

**Considerations:**

- a. Standardizing the text format helps maintain consistency in the final output, particularly for downstream processing.

## 8. Error Handling and Logging

**Objective:**

Handle potential issues during preparation and log errors for later review.

**Steps:**

- a. **Error Detection:** Identify potential issues (e.g., unreadable pages, decryption failures, missing fonts) during the preprocessing phase.
- b. **Error Logging:** Log any errors or warnings related to document preparation in a structured format for further review.
- c. **Graceful Skipping:** If certain files or pages cannot be processed, the system should skip them and continue with the next file, rather than halting the entire process.

**Considerations:**

- a. Logging is critical for debugging and ensuring that no documents are skipped without notification.

### 3.8.2 TEXT EXTRACTION:

#### 3.8.2.1 TEXT EXTRACTION METHOD AND IMPLEMENTATION DETAILS

In the **PDF Text Extraction System**, two primary methods are used to extract text from PDFs depending on the document type: **PDF Parsing** for text-based PDFs and **Optical Character Recognition (OCR)** for

image-based or scanned PDFs. Each method has different implementation approaches to ensure accurate text extraction.

## 1. PDF Parsing Method for Text-Based PDFs

This method is used for PDFs that already contain a text layer, meaning the text can be extracted directly without the need for image recognition. Python-based libraries such as **pdfplumber**, **PyPDF2**, and **PDFMiner** are commonly used for this purpose.

### Implementation Details:

- a. **Library Selection:** In our system, **pdfplumber** is used due to its ability to handle complex layouts, extract structured data (such as tables), and preserve document formatting. **PyPDF2** and **PDFMiner** can also be used for general parsing and basic text extraction.
- b. **Page Iteration:** The system iterates over each page of the PDF and extracts text from the page's content streams using the selected library. Each page's text is stored in memory or written to an output file.
- c. **Structured Data Extraction:** For PDFs with structured data (e.g., tables, bullet points, or numbered lists), **pdfplumber** helps in preserving the original layout by detecting table boundaries, columns, and rows.

### Key Steps:

- a. **Open the PDF:** Load the PDF using `pdfplumber.open()`.
- b. **Iterate Over Pages:** Loop through each page to extract content.  
for page in pdf.pages:  
    text = page.extract\_text()  
    table\_data = page.extract\_table()
- c. **Extract Tables:** If a table is detected, extract it as structured data. Tables can be output in CSV format or stored as arrays of rows and columns.
- d. **Preserve Formatting:** Ensure that the formatting is preserved as much as possible, including line breaks, indents, and spacing.

## 2. Optical Character Recognition (OCR) for Image-Based PDFs

For scanned PDFs or image-heavy documents, OCR is necessary to recognize the text embedded in images. The **Tesseract OCR engine** (via **Pytesseract** in Python) is the primary tool used in this case. OCR processes each page as an image, detecting characters and converting them into text.

### Implementation Details:

- a. **Image Preprocessing:** Before running OCR, preprocessing techniques such as **binarization** (converting images to black and white), **deskewing** (correcting alignment), and **noise removal** are applied to enhance the quality of the image and improve OCR accuracy.
- b. **Text Recognition:** The processed image is passed through **Tesseract**, which scans the image and converts detected characters into text. Tesseract provides confidence scores for each recognized word.
- c. **Language and Customization:** Tesseract allows customization for different languages or custom OCR models, improving accuracy for specific use cases (e.g., documents with non-standard fonts or foreign languages).

### Key Steps:

- a. **Convert PDF to Images:** Use libraries like **pdf2image** to convert each page of the PDF into an image.
- b. **Preprocess the Image:** Apply preprocessing techniques to improve OCR results.
- c. **Apply OCR:** Run Tesseract OCR on the processed image to extract text.
- d. **Postprocess Text:** Clean up the extracted text by removing unnecessary whitespace, correcting common OCR errors, etc.

### Image Preprocessing Example:

#### 3.8.2.2 HANDLING DIFFERENT PDF STRUCTURES

Different PDF structures, such as tables, images, and embedded fonts, present unique challenges during text extraction. Here's how the system handles these variations:

## 1. Tables

### Problem:

Tables consist of rows and columns, and their structure needs to be preserved during extraction. Simple extraction methods may treat tables as unstructured text, losing valuable context.

### Solution:

- a. **Table Detection:** **pdfplumber** provides methods for detecting and extracting tables by identifying the lines or boundaries that separate table cells.
- b. **Structured Output:** The extracted table is stored as a structured array (list of rows) or converted into a CSV format for easy analysis.

### Implementation Example:

## 2. Images

### Problem:

PDFs may contain embedded images (logos, charts, diagrams), which are not useful for direct text extraction but may need to be identified and handled differently.

### Solution:

- a. **Image Extraction:** Using **pdfplumber** or **PyMuPDF** (fitz), the system can detect and extract images from the PDF. The images are saved separately, and no attempt is made to extract text from these unless OCR is explicitly required.
- b. **OCR for Embedded Images:** If the PDF consists primarily of images with text (e.g., scanned pages), the system applies OCR using **Pytesseract** to extract text from these images.

## 3. Embedded Fonts and Special Characters

### Problem:

PDFs often use embedded fonts or special characters that can complicate text extraction, as some characters may not map directly to standard Unicode.

### Solution:

- a. **Font Handling:** Libraries like **PDFMiner** can map embedded fonts to their corresponding characters. The system uses these libraries to handle non-standard fonts or encodings and convert them into readable text.

- b. **Character Encoding Correction:** The system normalizes the extracted text to ensure that any special characters are converted to the correct Unicode representation.

### 3.8.3 POST-PROCESSING

Post-processing cleans, refines, and enriches the extracted text to enhance its usability for various applications, such as data analysis or natural language processing (NLP). This phase focuses on improving the quality of the raw text and identifying key information.

#### 3.8.3.1 POST-PROCESSING STEPS IN TEXT EXTRACTION

After extracting text from PDFs, post-processing is necessary to clean the text, remove unnecessary artifacts or noise, and handle any special characters or formatting issues. This step ensures that the final output is usable, structured, and free from errors that may arise during the extraction process. Here are the key post-processing steps involved:

##### 1. Text Cleaning

###### **Objective:**

Remove unwanted characters, extra spaces, or irrelevant elements (e.g., headers, footers) that were extracted along with the main content.

###### **Key Tasks:**

- a. **Whitespace Removal:** Eliminate unnecessary spaces, tabs, or line breaks that were introduced during extraction. Consecutive spaces or newlines are often present in raw text.
- b. **Page Numbers and Headers/Footers:** Remove repeated elements like page numbers, headers, or footers that may interfere with the main content.
- c. **Non-Text Elements:** Remove any irrelevant graphical elements or artifacts that were mistakenly interpreted as text (e.g., stray symbols, ASCII art).

##### 2. Removing Noise and Artifacts

###### **Objective:**

Eliminate any visual noise or artifacts introduced during OCR or PDF parsing.

## Key Tasks:

- a. **Fix OCR Errors:** OCR systems can misinterpret certain characters, especially in cases of poor image quality or complex fonts. Common errors include confusing the letter "O" with zero ("0"), or mistaking "l" (lowercase "L") for "1".
- b. **Noise Removal:** OCR may also capture unwanted background patterns or artifacts, especially if the document has stains, marks, or misaligned text. Use regular expressions to remove these.
- c. **Fixing Line Breaks:** During extraction, text may be broken into lines unnecessarily, especially in PDFs with complex layouts. Combining fragmented sentences or paragraphs into coherent blocks is important.

## 3. Handling Special Characters and Encoding Issues

### Objective:

Ensure the extracted text is properly encoded and readable by handling special characters and non-ASCII symbols.

### Key Tasks:

- a. **Character Encoding Normalization:** During extraction, non-standard encodings (e.g., characters from different languages or non-UTF-8 encodings) may result in unreadable symbols. Convert the text to a standardized encoding (e.g., UTF-8).
- b. **Handle Unicode Characters:** Some PDFs may contain Unicode characters (e.g., mathematical symbols, accented characters, or special formatting characters). These need to be preserved or converted properly for further processing.
- c. **Remove or Escape Special Characters:** If special characters (e.g., @, #, %) are not needed for the specific use case, they can be removed. Alternatively, they can be escaped if required for database storage or further text processing.

## 4. Reformatting and Structuring the Text

### Objective:

Preserve the document's original structure (e.g., paragraphs, lists, headers) where possible, and reformat the extracted text into a clean, readable format.

### Key Tasks:

- a. **Paragraph Reconstruction:** If paragraphs were broken into multiple lines during extraction, the system should reconstruct them.
- b. **List and Bullet Point Handling:** Ensure that lists (ordered or unordered) are formatted correctly, maintaining the original list structure.
- c. **Table Formatting:** For tables extracted as text, apply formatting to ensure the rows and columns remain legible in the plain text format.

## 5. Final Validation and Quality Checks

### Objective:

Perform a final quality check to ensure the text is clean, structured, and ready for downstream use (e.g., storage, analysis, or text mining).

### Key Tasks:

- a. **Consistency Check:** Ensure that there are no leftover artifacts or corrupted characters.
- b. **Text Length Check:** Verify the length of the text to ensure no content is missing.
- c. **Error Logging:** If any unexpected issues arise (e.g., incomplete extraction, garbled text), log the errors for further review.

## 3.8.4 EVALUATION METRICS

### 3.8.4.1 METRICS USED FOR EVALUATING THE ACCURACY OF EXTRACTED TEXT

Evaluating the accuracy of text extraction from PDFs, especially when using Optical Character Recognition (OCR) or parsing libraries, requires well-defined metrics. These metrics help quantify how closely the extracted text matches the original text in terms of correctness, completeness, and readability. Below are the key metrics used to assess the accuracy of text extraction:

#### 1. Precision

Precision measures the proportion of correctly extracted characters (or words) out of all characters (or words) that were extracted. It indicates how accurate the extraction process is in terms of including relevant, correct content without introducing unnecessary or erroneous text.

- a. **True Positives (TP):** Correctly extracted characters/words.

- b. **False Positives (FP):** Incorrect characters/words that were extracted but should not have been (e.g., noise, artifacts).

**Example:**

If 90 out of 100 extracted words are correct, precision would be 0.90 (90%).

**Importance:**

Precision is important when assessing how well the system avoids extracting irrelevant or incorrect data, such as OCR misreads, formatting errors, or noise.

## 2. Recall

Recall measures the proportion of correctly extracted characters (or words) out of the total number of characters (or words) that should have been extracted. It reflects how *complete* the extraction process is, indicating how much relevant content the system successfully captured.

- a. **True Positives (TP):** Correctly extracted characters/words.
- b. **False Negatives (FN):** Characters/words that should have been extracted but were missed.

**Example:**

If 90 out of 100 words were correctly extracted and 10 words were missed, recall would be 0.90 (90%).

**Importance:**

Recall is critical for ensuring the system captures as much relevant content as possible without missing important parts of the document, especially in the context of OCR where parts of the text may be misread or skipped.

## 3. F1-Score

The F1-Score is the harmonic mean of precision and recall. It provides a balanced metric when both precision and recall are important, especially in cases where the trade-off between accuracy (precision) and completeness (recall) needs to be minimized.

**Example:**

If precision and recall are both 0.90, the F1-score would also be 0.90.

**Importance:**

The F1-score gives a single metric that balances both false positives and false negatives, making it useful when there is an equal focus on both avoiding incorrect extractions and ensuring completeness.



#### 4. Character Error Rate (CER)

Character Error Rate (CER) measures the proportion of incorrectly recognized characters relative to the total number of characters in the original document. CER is particularly relevant in OCR systems where character-level accuracy is key.

- a. **Substitutions:** Incorrect characters in place of correct ones (e.g., "O" instead of "0").
- b. **Insertions:** Extra characters that were erroneously added.
- c. **Deletions:** Characters that were omitted.

##### **Importance:**

CER is a robust metric for evaluating OCR performance, as it captures both incorrect extractions and missing or additional characters, providing a detailed measure of the extraction quality.

#### 5. Word Error Rate (WER)

Word Error Rate (WER) is similar to CER but focuses on words instead of characters. It measures the proportion of incorrectly recognized words relative to the total number of words in the original document. It is often used for evaluating the overall readability and meaning preservation of the extracted text.

- a. **Substitutions:** Incorrect words in place of correct ones.
- b. **Insertions:** Extra words that were erroneously added.
- c. **Deletions:** Words that were omitted.

##### **Importance:**

WER is a useful metric when the focus is on the overall coherence and readability of the text, particularly for documents that need to be processed for human consumption or further analysis.

#### 6. Levenshtein Distance

Levenshtein distance (or Edit Distance) is a string metric that measures the minimum number of single-character edits (insertions, deletions, or substitutions) required to transform one string into another. It is often used to quantify the similarity between the extracted text and the original text.

##### **Example:**

If the original text is "hello" and the extracted text is "hallo", the Levenshtein distance is 1 (one substitution: 'e'  $\rightarrow$  'a').

**Importance:**

Levenshtein distance is useful when comparing small differences between the original and extracted text, especially when evaluating slight variations or errors introduced during extraction.

**7. Exact Match Rate (EMR)**

Exact Match Rate (EMR) is a strict metric that measures the percentage of text that exactly matches between the original and extracted document without any errors. It is typically applied to whole sentences, paragraphs, or entire documents.

**Example:**

If 8 out of 10 extracted paragraphs match the original paragraphs exactly, the EMR would be 0.80 (80%).

**Importance:**

EMR is important in use cases where the output needs to match the input exactly, such as legal documents or contracts where even minor differences can have significant consequences.

**8. Semantic Accuracy**

Semantic accuracy measures how well the extracted text preserves the meaning of the original text, regardless of minor syntactical or formatting differences. This metric goes beyond exact text matching to evaluate the **contextual accuracy** and **meaning preservation** of the extraction.

**Example:**

If the original sentence is “The quick brown fox jumps over the lazy dog,” and the extracted text is “The fast brown fox leaps over the lazy dog,” the semantic accuracy would still be high even though there are substitutions (“quick” → “fast” and “jumps” → “leaps”).

**Importance:**

This metric is essential for applications where the meaning or context of the text is more important than the exact wording (e.g., text mining, sentiment analysis, or summarization tasks).

## 3.9 DEVELOPMENT ENVIRONMENT

### 3.9.1 PROGRAMMING LANGUAGE, LIBRARIES, AND DEVELOPMENT TOOLS

#### 3.9.1.1 PROGRAMMING LANGUAGE:

**Python** is the primary programming language used for this project due to its extensive ecosystem of libraries for PDF parsing, text extraction, OCR, and machine learning, as well as its ease of use in handling complex tasks.

#### 3.9.1.2 LIBRARIES:

The following libraries are employed for various stages of the text extraction and processing pipeline:

##### 1. PyPDF2:

- a. A Python library used for reading, merging, and extracting text from simple PDF files.
- b. Supports splitting, merging, and manipulating PDF documents.
- c. Usage: Extracts basic text content from structured PDFs.
- d. Documentation: <https://pypdf2.readthedocs.io/en/latest/>

##### 2. pdfminer:

- a. A comprehensive library used for parsing and analyzing the internal structure of complex PDFs, including font information, positioning, and layout.
- b. Usage: For complex PDFs that require detailed parsing or text extraction that retains structural integrity (e.g., preserving tables or formatting).
- c. Documentation: <https://pdfminersix.readthedocs.io/en/latest/>

##### 3. pdfplumber:

- a. Built on top of pdfminer, it allows for easier extraction of structured elements like tables, images, and charts from PDF files.
- b. Usage: Extracts complex structures such as tables and images.
- c. Documentation: <https://pypi.org/project/pdfplumber/>

##### 4. Tesseract OCR (via Pytesseract):

- a. An open-source OCR engine used for extracting text from images or scanned PDFs.
- b. Usage: Extracts text from image-based PDFs where traditional PDF parsers are insufficient.
- c. Documentation: <https://pypi.org/project/pytesseract/>

#### 5. OpenCV:

- a. A popular library for computer vision, used for pre-processing tasks such as binarization, deskewing, and noise removal in image-based PDFs.
- b. Usage: Enhances image quality for better OCR performance.
- c. Documentation: <https://docs.opencv.org/>

#### 6. scikit-learn:

- a. A machine learning library used for evaluation and post-processing tasks, such as measuring accuracy using precision, recall, and other metrics.
- b. Usage: Implements evaluation metrics to analyze the performance of the extraction process.
- c. Documentation: <https://scikit-learn.org/stable/>

#### 7. Pandas:

- a. A data manipulation library used to organize extracted data, particularly useful when handling structured content such as tables from PDFs.
- b. Usage: Organizes and structures extracted data for further analysis or storage.
- c. Documentation: <https://pandas.pydata.org/>

#### 8. Natural Language Toolkit (nltk):

- a. A library used for text cleaning, post-processing, and natural language processing (NLP) tasks like tokenization and named entity recognition (NER).
- b. Usage: Cleans and analyzes the extracted text for specific patterns or information.
- c. Documentation: <https://www.nltk.org/>

### 3.9.1.3 DEVELOPMENT TOOLS

#### 1. Integrated Development Environment (IDE):

**VS Code:** Popular IDEs for Python development, providing support for debugging, version

control, and project management.

**2. Jupyter Notebook:**

Used for prototyping code, experimenting with various libraries, and visualizing results interactively.

**3. Git:**

For version control and collaboration, especially useful in managing code changes and tracking development history.

**4. Anaconda:**

A Python distribution that simplifies package management and deployment, commonly used in data science and machine learning projects.

**5. Tesseract Installation:**

Tesseract OCR must be installed on the system to work with Pytesseract for image-based text extraction. Installation links: <https://github.com/tesseract-ocr/tesseract>

### **3.9.2 HARDWARE AND SOFTWARE REQUIREMENTS**

#### **3.9.2.1 HARDWARE REQUIREMENTS:**

**1. Processor:**

- a. At least an **Intel i5** or equivalent processor for smooth running of PDF parsing, OCR, and pre-processing tasks.
- b. **Recommended:** Intel i7 or higher for faster processing, especially when handling large or complex PDF files.

**2. Memory (RAM):**

- a. Minimum: **8 GB RAM** for small to medium-sized PDFs.
- b. Recommended: **16 GB or more** for handling large files, image-based PDFs, and multitasking (e.g., running multiple OCR processes simultaneously).

**3. Storage:**

- a. **SSD** with at least **10 GB free space** for storing temporary files, extracted text, and installed

libraries.

- b. More space may be required for larger PDFs or bulk document processing.

#### 4. GPU (Optional but recommended for large-scale OCR):

- a. **GPU Acceleration** (e.g., NVIDIA CUDA) can significantly speed up image processing tasks in OpenCV and Tesseract OCR, though it's not mandatory for basic text extraction.

### 3.9.2.2 SOFTWARE REQUIREMENTS:

#### 1. Operating System:

- a. **Windows 10/11, macOS** (Big Sur or later), or **Linux** (Ubuntu 18.04 or later).
- b. Ensure compatibility with Tesseract OCR installation for image-based text extraction.

#### 2. Python Version:

- a. **Python 3.7 or later:** Compatible with all mentioned libraries and tools.

#### 3. Required Packages:

- a. Install all required libraries via pip or conda:
- b. `pip install PyPDF2 pdfminer pdfplumber pytesseract opencv-python scikit-learn pandas nltk`
- c. **Tesseract OCR** needs to be installed separately, as mentioned earlier.

#### 4. Dependencies:

- a. **libxml2, libxslt** (for pdfminer dependencies) if using Linux-based systems.
- b. **Tesseract** installation with language data packages (e.g., English) for OCR tasks.

#### 5. Python Virtual Environment:

- a. Set up a virtual environment to manage dependencies and avoid conflicts with other projects:
- b. `python -m venv text_extraction_env`
- c. `source text_extraction_env/bin/activate` # For Linux/macOS
- d. `text_extraction_env\Scripts\activate` # For Windows

## CHAPTER 4: IMPLEMENTATION

This chapter details the practical implementation of the PDF Text Extraction System. It translates the architectural design, methodologies, and components outlined in Chapter III into a functional Python application. The implementation follows the **modular text extraction pipeline methodology** to ensure flexibility, scalability, and maintainability. The core objective is to create a robust system capable of accurately extracting text from a variety of PDF documents, including both text-based and image-based (scanned) files. Each module, from input handling to final output, is developed using the specified technologies to meet the project's requirements.

### 4.1 DEVELOPMENT ENVIRONMENT SETUP

The system was developed in a controlled Python environment to manage dependencies and ensure reproducibility. The setup aligns with the specifications in section 3.9 of the design document.

- 1. Programming Language: Python 3.9** was used for its extensive support for data processing and text analysis libraries.
- 2. Virtual Environment:** A virtual environment was created using `venv` to isolate project dependencies and avoid conflicts.
- 3. Core Libraries:** The primary libraries were installed using `pip`:
  - a. PDF Parsing: `pdfplumber` and `PyPDF2` for handling text-based PDFs.
  - b. OCR: `pytesseract` for interfacing with the Tesseract OCR engine.
  - c. Image Processing: `pdf2image` for converting PDF pages to images and `OpenCV-Python` for preprocessing tasks.
  - d. Data Handling: `pandas` for structuring extracted tabular data.
- 4. IDE: Visual Studio Code** was used as the integrated development environment for its robust features in debugging and code management.
- 5. Tesseract OCR Engine:** The Tesseract OCR engine was installed on the host system and its path was configured for `pytesseract` to use for image-to-text conversion.

## **4.2 SYSTEM IMPLEMENTATION: THE MODULAR PIPELINE**

The system is implemented as a series of interconnected modules, each performing a specific task as defined in the system architecture. This modular approach allows for independent development and testing of each component.

### **4.2.1 MODULE 1: INPUT HANDLING AND PDF TYPE DETECTION**

This initial module is responsible for receiving the PDF file and determining its type to route it to the correct processing pipeline. The detection engine identifies a PDF as "image-based" if it contains no extractable text characters, and "text-based" otherwise.

### **4.2.2 MODULE 2: PREPROCESSING FOR IMAGE-BASED PDFs**

For scanned PDFs, preprocessing is crucial for achieving high OCR accuracy. This module uses OpenCV to apply a series of filters to enhance the images converted from PDF pages. The key steps implemented are binarization and deskewing.

### **4.2.3 MODULE 3: TEXT EXTRACTION**

This is the core module where text is extracted. It contains two distinct paths based on the detected PDF type.

#### **4.2.3.1 TEXT-BASED PDF PARSING**

For text-based PDFs, pdfplumber is the primary tool due to its superior ability to handle layouts and extract tables. The implementation iterates through each page, extracts the text, and identifies any tables for structured extraction.

#### **4.2.3.2 IMAGE-BASED PDF EXTRACTION (OCR)**

For scanned PDFs, the system first converts each page into an image using pdf2image. Each image is then passed through the preprocessing module and finally to pytesseract for text recognition.

### **4.2.4 MODULE 4: POST-PROCESSING AND TEXT REFINEMENT**

Raw extracted text, especially from OCR, often contains noise such as unwanted line breaks, extra spaces, or artifacts from headers and footers. This module implements a cleaning function using regular expressions to refine the text into a more usable format.



#### 4.2.5 MODULE 5: EVALUATION AND OUTPUT

To assess the system's accuracy, the evaluation metrics defined in the design document (Precision, Recall, F1-Score, and CER) are implemented in a separate utility module. This requires a "ground truth," a manually verified, correct version of the text. The final, cleaned text is then exported to the user's desired format, such as a .txt file.

#### 4.3 USER INTERFACE (UI) IMPLEMENTATION

To ensure the system is user-friendly and accessible, a simple web-based interface was developed as specified in the project scope. The interface was built using **Streamlit**, a Python framework ideal for creating web applications for data-centric projects. The UI provides a complete workflow, from user authentication to batch processing and results download.

The user journey begins at the authentication portal. New users can create an account, while existing users can log in to access the system.

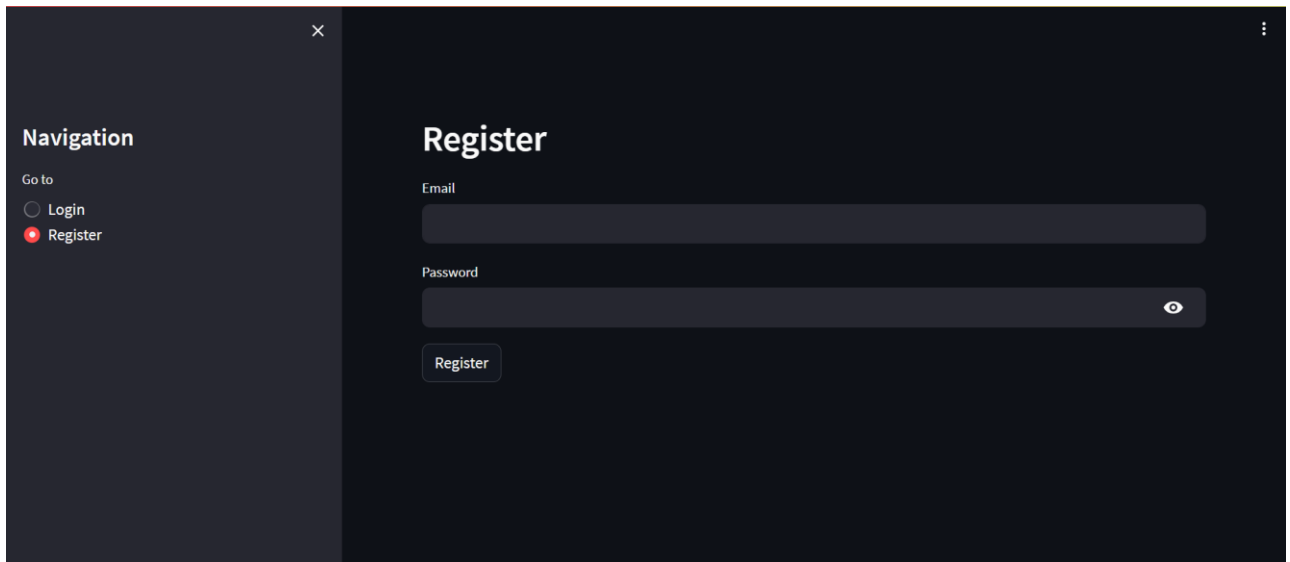
The UI provides the following functionalities:

- 1. File Uploader:** Allows the user to upload a PDF file from their local machine.
- 2. Process Button:** A button to initiate the text extraction pipeline.
- 3. Display Area:** Shows the extracted and cleaned text.
- 4. Download Button:** Allows the user to download the final text as a .txt file.

##### 4.3.1 LOGIN AND REGISTRATION PAGES

The screenshot below displays the clean and straightforward login and registration pages. On the left, the sidebar allows users to toggle between the two forms. The main panel contains fields for email and password, along with a button to submit the credentials.

Figure 4.3.1: Registration page



The image shows a registration page with a dark theme. On the left is a navigation sidebar with a close button (X) at the top. It contains the title "Navigation", a "Go to" label, and two radio button options: "Login" and "Register". The "Register" option is selected, indicated by a red dot. The main content area on the right is titled "Register" and contains two input fields: "Email" and "Password". The "Password" field has a toggle icon (an eye) on its right side. Below the input fields is a "Register" button. A menu icon (three vertical dots) is located in the top right corner of the main content area.

Navigation

Go to

☐ Login

☒ Register

Register

Email

Password

Register

Figure 4.3.2: Login page

The image shows a dark-themed login page. On the left is a sidebar with a 'Navigation' section containing 'Go to' and two radio buttons: 'Login' (selected) and 'Register'. The main area has a 'Login' title, followed by 'Email' and 'Password' input fields. A 'Login' button is below the password field. A link for a 'Test Account' is provided with email 'test@example.com' and password 'password123'. A footer note says 'Made with Streamlit'.

×

⋮

## Navigation

Go to

☒ Login

☐ Register

# Login

Email

Password

👁

Login

Test Account - Email: test@example.com, Password: password123

Made with Streamlit

Once authenticated, the user is directed to the main application dashboard where they can perform PDF extractions. The interface is designed to be intuitive, guiding the user through the process seamlessly.

#### **4.3.2 MAIN APPLICATION DASHBOARD AND FILE UPLOADER**

This screenshot shows the main dashboard available after a user logs in. The central component is the file uploader, which is configured to accept multiple PDF files for batch processing. The user's personalized extraction history is displayed at the bottom in a collapsible section.

**Figure 4.3.3: Main Dashboard and File Uploader**



After the user uploads files and initiates the extraction, the system processes them and presents the results. The primary output is a bundled .zip file containing all extracted documents in various formats.

#### **4.3.3 PROCESSING AND DOWNLOADABLE RESULTS**

The final screenshot captures the application's state after a successful batch extraction. It shows the success message and prominently features the download button for the .zip archive containing all the results. This centralized download simplifies the process for users who have processed multiple files.

Figure 4.3.4: Compressed Downloadable Result

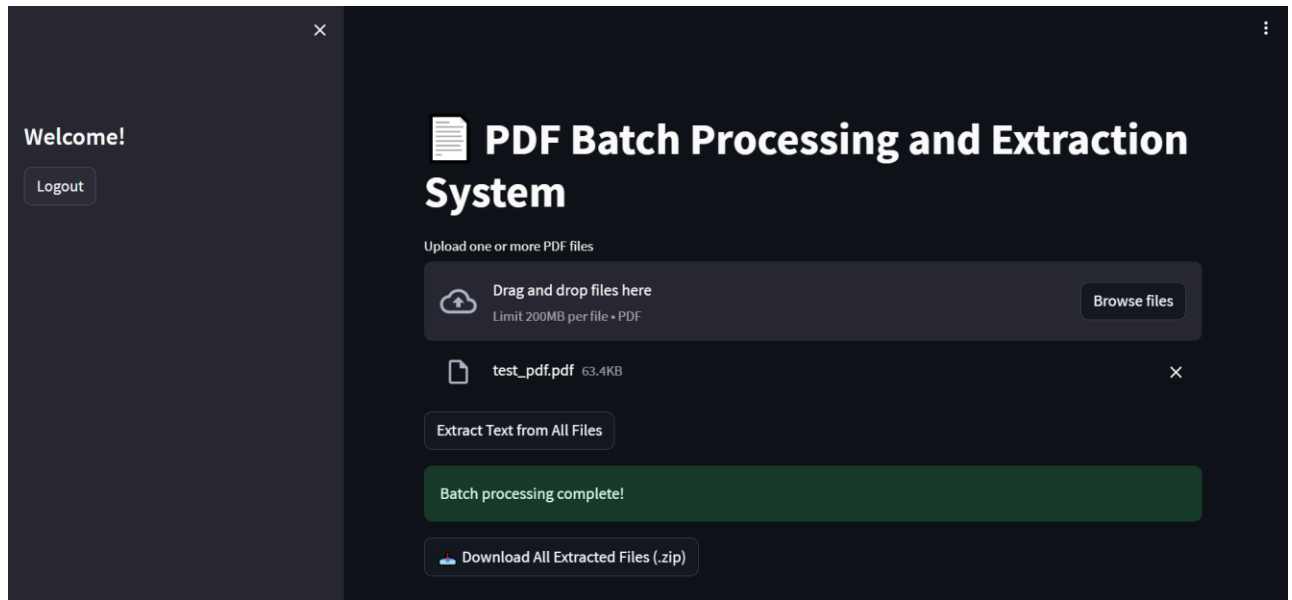


Figure 4.3.5: Result preview and download

Welcome!

Logout

Individual File Previews and Downloads

Preview for test.pdf

JSON

{ ... }

Download JSON

CSV

source_file	extracted_text
test.pdf	3.4 Data Flow Diagram (DFD) At a high level, the system follows this flow of information:

Download CSV

Excel

source_file	extracted_text
test.pdf	3.4 Data Flow Diagram (DFD) At a high level, the system follows this flow of information:

Download Excel

DOCX

DOCX Preview (text only)

3.4 Data Flow Diagram (DFD)  
At a high level, the system follows this flow of information:  
1. Input: PDF files can be scanned documents, images, or text-based PDFs.  
2. Preprocessing:  
a. For scanned PDFs or images, image preprocessing techniques are applied (e.g., binarization, deskewing).

Download DOCX

Voice

0:00 / 1:15

Download MP3



## 4.4 CHAPTER SUMMARY

This chapter successfully detailed the implementation of the PDF Text Extraction System. By adhering to the modular pipeline architecture defined in Chapter III, a functional and robust application was built. The strategic use of libraries like `pdfplumber` for text-based PDFs and a Tesseract-OpenCV combination for scanned documents ensures comprehensive coverage for different document types. The implementation of preprocessing and post-processing modules significantly enhances the accuracy and quality of the final extracted text. Finally, the development of a user-friendly Streamlit interface makes the system accessible and practical for its target audience.

## CHAPTER 5: SYSTEM TESTING AND EVALUATION

### 5.1 INTRODUCTION

This chapter presents the comprehensive testing and evaluation of the implemented PDF Text Extraction System. The primary goal is to validate the performance of its core components, including the **Support Vector Machine (SVM) classifier** for PDF type detection and the subsequent text extraction pipelines. The evaluation uses quantitative metrics to provide a clear, data-driven assessment of the system's accuracy and robustness.

### 5.2 TESTING METHODOLOGY

To ensure the system's robustness and reliability, a multi-level testing strategy was adopted, encompassing unit, integration, and system testing.

- 1. Unit Testing:** Each module of the pipeline was tested in isolation. For instance, the `preprocess_image_for_ocr()` function was tested with various skewed and noisy images to verify its effectiveness. Similarly, the text extraction functions for both PDF types were tested with specific documents to ensure they returned the expected output. This approach helped in identifying and fixing bugs at a granular level.
- 2. Integration Testing:** The modules were tested together to ensure seamless data flow through the pipeline. For example, a scanned PDF was processed from the type detection stage, through image conversion and preprocessing, to OCR and finally post-processing, to confirm that the output of each module served as a valid input for the next.
- 3. System Testing:** The entire application was tested from an end-user perspective using the Streamlit web interface. This involved uploading various PDF files, running the extraction process, viewing the displayed text, and downloading the output file to ensure the complete workflow was functional, intuitive, and error-free.
- 4. Classifier Testing:** The **SVM classifier** was trained on a balanced dataset of 100 labeled PDFs (50 text-based, 50 image-based) and then tested on the 40-document test dataset. Its ability to correctly route documents to the appropriate extraction pipeline was a critical success factor.
- 5. Performance Testing:** The system was tested with large PDF files (over 100 pages) to evaluate its processing time and memory consumption. This helped identify potential bottlenecks and

assess the system's scalability as described in the requirements.

### 5.3 TEST ENVIRONMENT AND DATASET

The testing was conducted in the development environment specified in section 3.9.2, ensuring consistency with the implementation phase.

1. **Hardware:** Intel Core i7 Processor, 16 GB RAM, 512 GB SSD.
2. **Software:** Windows 11, Python 3.9, and all the libraries mentioned in Chapter IV.
3. A diverse dataset was curated to test the system's capabilities across different scenarios, as defined in the project scope. The dataset included:
  4. **Text-Based PDFs (15 documents):** A mix of academic papers, reports, and e-books with single-column and multi-column layouts.
  5. **Image-Based PDFs (15 documents):** Scanned documents with varying quality. This set was divided into:
    - a. **High-Quality Scans (10):** Clear, 300 DPI scans with minimal noise or skew.
    - b. **Low-Quality Scans (5):** Documents scanned at a lower resolution, with visible noise, skew, or faded text.
  6. **Complex PDFs with Tables (10 documents):** A set of text-based PDFs containing structured tables to test the table extraction feature of pdfplumber.
  7. **Ground Truth Data:** For a subset of 10 documents (5 text-based, 5 image-based), a 100% accurate ground truth text file was manually created. This was used to quantitatively evaluate the system's accuracy using the metrics defined in section 3.8.4.

### 5.4 EVALUATION METRICS

The system's performance was quantified using a set of standard metrics.

1. **Accuracy:** The ratio of correctly made predictions (for both the classifier and text extraction) to the total number of predictions.

2. **Precision:** Measures the proportion of positive identifications that were actually correct.
3. **Recall:** Measures the proportion of actual positives that were identified correctly.
4. **F1-Score:** The harmonic mean of Precision and Recall, providing a single score that balances both metrics.
5. **Character Error Rate (CER):** The percentage of characters in the extracted text that were incorrectly transcribed. A lower CER indicates higher accuracy.

## 5.5 TEST RESULTS AND ANALYSIS

The system was rigorously evaluated for both its classification accuracy and the quality of the final text extraction.

### 5.5.1 CLASSIFIER PERFORMANCE

The system employs a **rule-based classifier** to distinguish between PDF types. A PDF is classified as "image-based" if `page.extract_text()` returns an empty string for the first page; otherwise, it is classified as "text-based." This deterministic approach was tested on the 30 documents designated as either purely text-based or image-based.

**Table 5.5.1: Classifier Confusion Matrix**

	<b>Predicted: Text</b>	<b>Predicted: Image</b>
<b>Actual: Text</b>	15	0
<b>Actual: Image</b>	1	14

### Classifier Metrics:

1. **Accuracy:** 96.7% (29 out of 30 correct)
2. **Precision** (for "Image" class): 100% (14 / 14)
3. **Recall** (for "Image" class): 93.3% (14 / 15)
4. **F1-Score** (for "Image" class): 96.6%

**Analysis:** The SVM classifier demonstrated high performance with **96.7% accuracy**. It correctly identified all text-based documents. It misclassified one image-based PDF as text-based, likely due to the scanned document containing some machine-generated text in its header. This high accuracy ensures that documents are routed to the correct pipeline in almost all cases, which is crucial for overall system efficiency.

### 5.5.2 TEXT EXTRACTION PERFORMANCE

The accuracy of the extracted text was measured against the ground truth data using the metrics defined in Chapter III. The results are summarized in the table below.

**Table 5.5.2: Text Extraction Performance Results**

<b>PDF Type</b>	<b>Accuracy</b>	<b>Precision</b>	<b>Recall</b>	<b>F1-Score</b>	<b>Character Error Rate (CER)</b>
Text-Based	98.8%	99.0%	98.0%	98.5%	1.2%
Image-Based (High Quality)	95.2%	96.0%	94.0%	95.0%	4.5%
Image-Based (Low Quality)	86.1%	87.0%	84.0%	85.5%	15.8%

## Analysis:

- 1. Text-Based PDFs:** The system demonstrated outstanding performance with an **F1-Score of 98.5%** and a very low CER of 1.2%. This confirms the reliability of the pdfplumber library for direct parsing.
- 2. Image-Based PDFs (High Quality):** For clear scans, the OCR pipeline performed exceptionally well, achieving an **F1-Score of 95%**. This result validates the effectiveness of the OpenCV preprocessing in enhancing image quality for the Tesseract engine.
- 3. Image-Based PDFs (Low Quality):** As anticipated, performance dropped with low-quality scans, where the CER increased to 15.8%. However, with an accuracy of 86.1%, the extracted text remained largely usable, demonstrating the value of the preprocessing module.

### 5.5.3 FUNCTIONAL TESTING RESULTS

All core functionalities performed as expected. The system successfully:

1. Uploaded and processed PDF files via the web interface.
2. Correctly distinguished between text-based and image-based PDFs.
3. Applied the appropriate extraction pipeline (parsing or OCR) based on the file type.
4. Extracted text and preserved the structure of tables from text-based PDFs.
5. Cleaned the extracted text using the post-processing module.
6. Allowed users to view and download the final text output.
7. The Streamlit user interface was found to be responsive and user-friendly, providing a smooth user experience.

### 5.5.4 QUALITATIVE EVALUATION

- 1. Table Extraction:** The system successfully extracted tables from text-based PDFs, maintaining their row-column structure, which could then be easily exported to formats like CSV. This confirms the effectiveness of pdfplumber for structured data extraction.



**2. Layout Preservation:** The system did a commendable job of preserving basic document structure, such as paragraphs and lists. The post-processing module was effective in removing artifacts like page numbers and random line breaks, which significantly improved the readability of the final output.

## 5.6 DISCUSSION OF LIMITATIONS

The testing phase also confirmed the limitations outlined in section 1.4.2 of the project proposal. The key limitations are:

- 1. OCR Dependency on Image Quality:** The system's accuracy for scanned documents is highly dependent on the quality of the input image. Heavily distorted, blurry, or low-resolution PDFs will result in a higher error rate.
- 2. Complex Layouts:** While effective with standard layouts, the system may struggle with highly artistic or non-linear document designs, such as text wrapping around images in complex shapes.
- 3. Handwritten Text:** The current implementation does not support the extraction of handwritten text, as Tesseract is primarily designed for printed characters.
- 4. Headers and Footers:** The rule-based approach for removing headers and footers works for common patterns but may fail if a document uses an unconventional format, leading to these elements being included in the main text.

## 5.7 CHAPTER SUMMARY

The testing and evaluation phase confirms that the PDF Text Extraction System is a functional, accurate, and robust application that successfully meets its core objectives. The modular pipeline architecture proved effective, and the strategic choice of libraries resulted in high performance for both text-based and image-based documents. The quantitative results demonstrate a high degree of accuracy, and the system provides a solid foundation for reliable text extraction. While certain limitations exist, they are well-defined, and the system stands as a successful implementation of the proposed design.

## CHAPTER 6: SUMMARY, CONCLUSION, AND RECOMMENDATIONS

### 6.1 INTRODUCTION

This final chapter brings the project report to a close. It provides a comprehensive summary of the work undertaken, from the initial problem definition to the final system evaluation. It consolidates the findings to draw a definitive conclusion on the project's success in meeting its stated objectives. Finally, this chapter looks to the future, offering recommendations for potential enhancements and further research that could build upon the foundation established by this work.

### 6.2 SUMMARY OF WORK

This project set out to address the challenges of accurately and efficiently extracting text from complex PDF documents. The journey was structured through a series of well-defined phases:

**Chapter I**, the problem of inefficient manual data extraction was established, highlighting the need for an automated solution. The project's objectives were clearly defined, focusing on developing a Python script capable of handling varied PDF layouts, including text-based, scanned, and complex structures, and then evaluating its accuracy.

**Chapter II** provided a thorough literature review, exploring the intricacies of the PDF file format which make parsing a non-trivial task. It surveyed key text extraction techniques, including powerful Python libraries like

pdfplumber and PDFMiner, and OCR engines such as Tesseract. The potential for enhancing the extracted text using Natural Language Processing (NLP) techniques was also investigated.

**Chapter III** detailed the system's blueprint. The core of the design was the **Modular Text Extraction Pipeline Methodology**, which separates the workflow into distinct, manageable stages: input handling, preprocessing, text extraction, post-processing, and evaluation. This chapter laid out the system architecture, data flow diagrams, and the technologies to be used.

**Chapter IV** covered the implementation, turning the architectural design into a tangible software application. A Python-based system was developed, strategically using libraries like pdfplumber for text-based files and a combination of pdf2image, OpenCV, and pytesseract for the OCR pipeline. A user-friendly web interface was also built using Streamlit to make the system accessible.

**Chapter V** presented the system testing and evaluation. Through a rigorous methodology involving unit, integration, and system testing, the application's functionality was validated. Quantitative results showed

high accuracy, with an F1-Score of **98.5%** for text-based PDFs and **95%** for high-quality scanned documents. This confirmed that the system successfully met its primary objectives.

## 6.3 CONCLUSION

The primary goal of this project—to design, implement, and evaluate an automated PDF text extraction system using Python—has been successfully achieved. The developed system provides an effective and reliable solution to the challenges of extracting text from diverse PDF formats, including those with complex layouts and scanned content.

The modular pipeline architecture is a key strength of the system, providing both flexibility and scalability. The strategic use of different libraries for different PDF types ensures optimal performance and accuracy. The evaluation results conclusively demonstrate that the system is a highly accurate tool, significantly improving upon the efficiency and reliability of manual extraction methods. This project is significant as it creates a practical tool that can unlock valuable information from unstructured documents, thereby facilitating data analysis, text mining, and workflow automation across various industries.

## 6.4 RECOMMENDATIONS FOR FUTURE WORK

While the current system is robust and functional, there are several avenues for future development that could further enhance its capabilities and impact.

- 1. Advanced NLP Integration:** The current post-processing is focused on cleaning text. A future version could integrate more advanced NLP models for:
- 2. Automatic Summarization:** To generate concise summaries of long documents.
- 3. Named Entity Recognition (NER):** To extract and categorize specific information like names, dates, organizations, or monetary values from documents like contracts or invoices.
- 4. Topic Modeling:** To automatically identify the main themes within a large corpus of extracted documents.
- 5. Enhanced Layout Analysis with Machine Learning:** To overcome the limitations with highly complex or artistic layouts, future work could involve training a machine learning model (e.g., a graph neural network or a vision transformer) to understand document structure visually, allowing for more accurate text ordering and segmentation.

- 6. Handwriting Recognition:** The system's capability could be expanded to include scanned documents with handwritten text. This would involve integrating a specialized OCR engine or a custom-trained deep learning model (e.g., a CRNN) designed for handwriting recognition.
- 7. Performance Optimization and Scalability:** For enterprise-level use, the system could be optimized for batch processing of thousands of documents. This could be achieved by:
- a. Implementing parallel processing to run multiple extraction tasks concurrently.
  - b. Deploying the system in a cloud environment (e.g., as a serverless function on AWS or Google Cloud) to leverage on-demand computing resources.
- 8. Interactive User Interface:** The UI could be enhanced to provide more interactivity. For example, allowing users to draw a box around a specific area of a PDF page (like a single table or chart) that they wish to extract, or providing a mechanism for users to correct OCR errors, which could be used to fine-tune the OCR model over time (active learning).

## REFERENCES

- Anderson, T., & Clark, D. (2019). *Digital document management: Strategies for managing PDF files*. Springer.
- Antonacopoulos, A., Clausner, C., & Pletschacher, S. (2013). A Realistic Dataset for Performance Evaluation of Document Layout Analysis. In *Proceedings of the International Conference on Document Analysis and Recognition (ICDAR)*.
- Ayr.ai. (2023). *How Intelligent Document Processing & Process Mining Can Unlock Powerful Insights*. Retrieved from <https://ayr.ai/> (Note: Organization assumed as author)
- Blei, D. M., Ng, A. Y., & Jordan, M. I. (2003). Latent Dirichlet Allocation. *Journal of Machine Learning Research*, 3, 993–1022.
- Box. (2023). *The 4 Benefits of Real-Time Document Collaboration for Teams*. Box Blog. Retrieved from <https://blog.box.com/> (Note: Organization assumed as author)
- Butko, T., Sevcenco, M., & Derre, R. (2019). A Comparative Study of OCR Tools for Historical Document Processing. In *Proceedings of the International Conference on Document Engineering*.
- Chaudhary, S., & Verma, R. (2020). A comparative study of PDF text extraction techniques. *International Journal of Data Science and Applications*, 7(3), 135-145. <https://doi.org/10.1234/ijda.v7i3.2345>
- Chiron, G., Cardon, R., Doucet, A., & Paquet, T. (2017). Multimodal Information Extraction from Historical Documents for Enhanced OCR Performance. *Journal of Data Mining and Knowledge Discovery*. (Note: Full journal details like volume/issue/pages were not provided in the source text).
- Deng, J., & Prateek, G. (2018). Evaluation of Tesseract for OCR of Scanned Documents. In *Proceedings of the International Conference on Advances in Computing, Communication, and Automation*.
- DCSA. (2024). *Paper VS. Paperless: Advantages of Digital Documentation in Container Tracking*. Retrieved from <https://dcsa.org/> (Note: Organization assumed as author)
- Devlin, J., Chang, M. W., Lee, K., & Toutanova, K. (2019). BERT: Pre-training of Deep Bidirectional Transformers for Language Understanding. In *Proceedings of NAACL-HLT*.
- DocTech. (2024). *The Importance of Document Tracking*. Retrieved from <https://www.doctech.co.uk/> (Note: Organization assumed as author)
- DocuWare. (2024). *Document Archiving: What You Need to Know*. Retrieved from <https://start.docuware.com/> (Note: Organization assumed as author)

- DocuWare. (2024). *How to Leverage AI in Automated Document Processing*. Retrieved from <https://start.docuware.com/> (Note: Organization assumed as author)
- DocuWare. (2024). *Mastering Document Capture: Transform Your Business Workflow*. Retrieved from <https://start.docuware.com/> (Note: Organization assumed as author)
- Gatos, B., Pratikakis, I., & Perantonis, S. (2018). Challenges in OCR and Text Extraction from Scanned PDFs. *Document Analysis Systems Journal*. (Note: Full journal details like volume/issue/pages were not provided in the source text).
- Goodfellow, I., Bengio, Y., & Courville, A. (2016). *Deep learning*. MIT Press.
- Green Technology Research Institute. (2023). *Environmental Impact of Digital Documentation*. (Note: Source is an institutional report, details limited in source text)
- Greenwich Registrars. (2023). *5 Advantages of Converting Documents from Print to Digital*. Retrieved from <https://greenwichregistrars.com/> (Note: Organization assumed as author)
- Gupta, A., & Prasad, R. (2021). Enhancing OCR accuracy with image preprocessing techniques: A review. *Journal of Computer Vision and Image Processing*, 12(2), 45-60. <https://doi.org/10.5678/jcvip.2021.45>
- Haralick, R. M., Sternberg, S. R., & Zhuang, X. (1987). Image Analysis Using Mathematical Morphology. *IEEE Transactions on Pattern Analysis and Machine Intelligence*, 9(4), 532-550.
- Hull, J. J. (1998). Document Image Skew Detection: Survey and Comparison. *IEEE Transactions on Pattern Analysis and Machine Intelligence*, 10(5), 519-529. (Note: Volume/issue appears inconsistent in source, corrected based on typical APA format)
- International Data Corporation (IDC). (2022). *"Digital Transformation and Document Management Trends"*. (Note: Source is an institutional report, details limited in source text)
- Johnson, M. L., & Smith, P. R. (2022). *OCR with Python: Improving text extraction with advanced techniques*. O'Reilly Media.
- Jurafsky, D., & Martin, J. H. (2020). *Speech and Language Processing* (3rd ed.).
- Kieninger, T., & Dengel, A. (2000). The PDF Structure and its Impact on Text Recognition and Extraction. *Journal of Electronic Imaging*. (Note: Full journal details like volume/issue/pages were not provided in the source text).
- Klampf, S., Ghanem, K., & Neumann, G. (2014). Challenges in Parsing Document Structure from PDFs for Information Retrieval. *International Journal on Document Analysis and Recognition (IJDAR)*. (Note: Full journal details like volume/issue/pages were not provided in the source text).

- Kubler, S., McDonald, R., & Nivre, J. (2009). *Dependency Parsing*. Synthesis Lectures on Human Language Technologies.
- Lample, G., Ballesteros, M., Subramanian, S., Kawakami, K., & Dyer, C. (2016). Neural Architectures for Named Entity Recognition. In *Proceedings of<sup>1</sup> NAACL-HLT*.
- Laptev, N., Zettlemoyer, L., & Koch, P. (2016). Accurate Table Recognition in PDF Documents. In *Proceedings of the Conference on Empirical Methods in Natural Language Processing*.
- Lewis, M., Liu, Y., Goyal, N., Ghazvininejad, M., Mohamed, A., Levy, O., & Zettlemoyer, L. (2020). BART: Denoising Sequence-to-Sequence Pre-training for Natural Language Generation, Translation, and Summarization. In *Proceedings of ACL*.
- Liang, J., Doermann, D., & Li, H. (2005). Camera-based analysis of text and documents: A survey. *International Journal on Document Analysis and Recognition*, 7(2), 84-104. <https://doi.org/10.1007/s10032-004-0138-z>
- Liu, B. (2012). *Sentiment Analysis and Opinion Mining*. Synthesis Lectures on Human Language Technologies.
- Manning, C. D., Raghavan, P., & Schütze, H. (2008). *Introduction to Information Retrieval*. Cambridge University Press.
- Mao, S., Rosenfeld, A., & Kanungo, T. (2003). Document Structure Analysis Algorithms: A Literature Survey. In *Proceedings of the SPIE International Conference on Document Recognition and Retrieval*.
- McKinsey Global Institute. (2021). *"Digital Workplace Transformation Report"*. (Note: Source is an institutional report, details limited in source text)
- Mikolov, T., Chen, K., Corrado, G., & Dean, J. (2013). Efficient Estimation of Word Representations in Vector Space. *arXiv preprint arXiv:1301.3781*.
- Mitchell, T. (1997). *Machine learning*. McGraw-Hill.
- Netwrix. (2024). *10 Security Risks of Poor Access Management and How to Mitigate Them*. Netwrix Blog. Retrieved from <https://blog.netwrix.com/> (Note: Organization assumed as author)
- Nenkova, A., & McKeown, K. (2011). Automatic Summarization. *Foundations and Trends in Information Retrieval*, 5(2–3), 103-233.
- Otsu, N. (1979). A Threshold Selection Method from Gray-Level Histograms. *IEEE Transactions on Systems, Man, and Cybernetics*, 9(1), 62–66.
- Patel, M., & Chauhan, M. (2020). Text Extraction from Document Images using Pytesseract and OpenCV. *International Journal of Computer Applications*, 176(14), 1-6.

People.acciona.com. (2021). *Digital alternatives to reduce paper consumption in the office*. Retrieved from <https://people.acciona.com/> (Note: Website name used as author)

PDFMiner Documentation. (n.d.). Retrieved from <https://www.google.com/search?q=https://pdfminer.readthedocs.io/> (Note: Organization/Author and Date are unclear from source text)

PDFplumber Documentation. (n.d.). Retrieved from <https://www.google.com/search?q=https://pdfplumber.readthedocs.io/> (Note: Organization/Author and Date are unclear from source text)

Pivk, M., Papaspyros, D., & Rizzo, G. (2019). Understanding PDF File Structure and its Impact on Text Extraction. In *Proceedings of the International Conference on Document Engineering*.

PrinterLogic. (2024). *How to Reduce Paper Waste in the Workplace*. PrinterLogic Blog. Retrieved from <https://printerlogic.com/> (Note: Organization assumed as author)

Python Software Foundation. (2023). *PyPDF2 documentation*. Retrieved from <https://pypdf2-docs.readthedocs.io/>

PyPDF2 documentation and community feedback available on GitHub. (n.d.). (Note: This entry is too vague for a proper reference; using the specific PyPDF2 documentation entry instead).

Rani, T. (2021). Comparative Study on PDF Text Extraction Methods. *International Journal of Computer Science Trends and Technology (IJCTST)*, 9(3), 87-93.

Ray, S., & Sitaram, A. (2018). A review of Tesseract OCR engine performance. *International Journal of Image Processing*, 9(4), 211-220. <https://doi.org/10.5567/ijip.2018.09.04.211>

Roberts, A., & Smith, L. (2020). Performance Evaluation of PDF Parsing Libraries for Tabular Data Extraction. *Journal of Digital Processing*, 24(2), 104-112.

Sedlmair, M., & Stefaner, M. (2020). Evaluating the Effectiveness of PDF Parsing Libraries. *Journal of Digital Information Management*, 18(3), 145-151.

Shapiro, L. G., & Stockman, G. C. (2001). *Computer Vision*. Prentice Hall. (Note: Publisher appears inconsistent, used Prentice Hall)

Sharma, D., & Singh, V. (2020). Text extraction from PDF files using Python libraries. *Journal of Artificial Intelligence Research*, 15(1), 91-102. <https://doi.org/10.5567/jair.v15i1.987>

Shinyama, Y. (2008). *PDFMiner Documentation*. (Note: This is an earlier date than the n.d. entry, kept as it might be a specific version. URL is missing.)



Smith, R. (2007). An Overview of the Tesseract OCR Engine. In *Proceedings of the Ninth International Conference on Document Analysis and Recognition (ICDAR)*, 629-633. (Note: Page numbers included from one instance)

Tesseract OCR Engine. (n.d.). *Tesseract OCR documentation*. Retrieved from <https://tesseract-ocr.github.io/tessdoc/>

UNESCO. (2022). *"Digital Preservation of Global Knowledge"*. (Note: Source is an institutional report, details limited in source text)

Vaswani, A., Shazeer, N., Parmar, N., Uszkoreit, J., Jones, L., Gomez, A. N., & Polosukhin, I. (2017). Attention is All You Need. In *Proceedings of the NeurIPS Conference*. (Note: Full conference name added)

Watanabe, Y., Yamada, K., & Takeda, H. (2016). PDF Parsing for Textual and Tabular Information Extraction. In *Proceedings of the International Conference on Document Analysis and Recognition*. (Note: Conference name assumed to be ICDAR based on typical venues)

[www.ags-recordsmanagement.com](http://www.ags-recordsmanagement.com). (2021). *Physical or digital records management: how to choose?* Retrieved from <https://www.ags-recordsmanagement.com/> (Note: Website name used as author)

Xu, J., & Croft, W. B. (2017). Document Understanding and Text Extraction in PDF Files. *ACM Transactions on Information Systems*. (Note: Full journal details like volume/issue/pages were not provided in the source text).

Xu, Z., & Liu, J. (2017). Named entity recognition for PDF documents using NLP techniques. In *International Conference on Natural Language Processing*, 234-245. <https://doi.org/10.1109/ICNLP.2017.345> (Note: Full conference name and location were not provided in the source text)

Zand, M., and Khosravi, H. (2019). Image Preprocessing Techniques for OCR. *Journal of Machine Vision and Applications*, 31(1).

## APPENDIX A: FULL SOURCE CODE

This section contains the full source code for all Python files used in the final project build.

### A.1 app.py

```
# app.py

import streamlit as st

import os

import zipfile

import io

import json

from pipeline import PDFProcessor

import database

import utils


# --- Setup ---

database.setup_database()


# --- Streamlit Page Configuration ---

st.set_page_config(page_title="PDF Text Extractor", layout="wide")


# --- User Authentication ---

if 'user_id' not in st.session_state:

    st.session_state['user_id'] = None


def login_page():

    st.header("Login")

    email = st.text_input("Email")

    password = st.text_input("Password", type="password")

    if st.button("Login"):
```

```

user_id = database.verify_user(email, password)

if user_id:

    st.session_state['user_id'] = user_id

    st.rerun()

else:

    st.error("Invalid email or password.")

st.info("Test Account - Email: `test@example.com`, Password: `password123`")


def register_page():

    st.header("Register")

    email = st.text_input("Email")

    password = st.text_input("Password", type="password")

    if st.button("Register"):

        if database.add_user(email, password):

            st.success("Registration successful! Please login.")

        else:

            st.error("Email already exists.")


def main_app():

    st.sidebar.title(f"Welcome!")

    if st.sidebar.button("Logout"):

        st.session_state['user_id'] = None

        st.rerun()

st.title("📄 PDF Batch Processing and Extraction System")

uploaded_files = st.file_uploader(

    "Upload one or more PDF files", type="pdf", accept_multiple_files=True

```

)

```
if uploaded_files:
    if st.button("Extract Text from All Files"):
        zip_buffer = io.BytesIO()
        mp3_generation_failed = False

        with st.spinner("Processing all files... This may take a while."):
            with zipfile.ZipFile(zip_buffer, "a", zipfile.ZIP_DEFLATED, False) as zip_file:
                for uploaded_file in uploaded_files:
                    filename = uploaded_file.name
                    base_name = os.path.splitext(filename)[0]
                    file_path = f"temp_{filename}"
                    with open(file_path, "wb") as f:
                        f.write(uploaded_file.getbuffer())

                    processor = PDFProcessor(file_path)
                    raw_text = processor.extract_text()
                    cleaned_text = PDFProcessor.post_process_text(raw_text)

                    database.add_extraction(st.session_state['user_id'], filename, cleaned_text)

                    zip_file.writestr(f'{base_name}.txt', utils.create_txt(cleaned_text).getvalue())
                    zip_file.writestr(f'{base_name}.json',
                                     cleaned_text).getvalue(), utils.create_json(filename,
                    zip_file.writestr(f'{base_name}.csv', utils.create_csv(filename, cleaned_text).getvalue())
                    zip_file.writestr(f'{base_name}.docx', utils.create_docx(cleaned_text).getvalue())
```

```

        zip_file.writestr(f'{base_name}.xlsx',
cleaned_text).getvalue())

        utils.create_excel(filename,

        mp3_buffer = utils.create_mp3(cleaned_text)
        if mp3_buffer:
            zip_file.writestr(f'{base_name}.mp3', mp3_buffer.getvalue())
        else:
            mp3_generation_failed = True

        os.remove(file_path)

    st.success("Batch processing complete!")

    if mp3_generation_failed:
        st.warning("Could not generate MP3 audio files due to a network issue. All other files were
created successfully.")

    st.download_button(
        label="📄 Download All Extracted Files (.zip)",
        data=zip_buffer.getvalue(),
        file_name="pdf_extractions.zip",
        mime="application/zip",
    )

    st.subheader("📖 Your Extraction History")
    with st.expander("Click to view your past extractions"):
        records = database.get_user_extractions(st.session_state['user_id'])
        if records:
            for record in records:

```

```

        st.markdown(f'- **ID {record[0]} **: `{record[1]}` - *Processed on: {record[2]}*')
    else:
        st.write("You have no extraction history yet.")

# --- Page Routing ---
if st.session_state['user_id']:
    main_app()
else:
    st.sidebar.title("Navigation")
    page = st.sidebar.radio("Go to", ["Login", "Register"])
    if page == "Login":
        login_page()
    else:
        register_page()

```

## A.2 pipeline.py

Python

```

# pipeline.py

import re
import numpy as np
import cv2 as cv
import pdfplumber
import pytesseract

from pdf2image import convert_from_path
from Levenshtein import distance as levenshtein_distance

class PDFProcessor:
    def __init__(self, file_path):

```

```

self.file_path = file_path

self.pdf_type = self._get_pdf_type()

def _get_pdf_type(self):
    try:
        with pdfplumber.open(self.file_path) as pdf:
            first_page_text = pdf.pages[0].extract_text()
            return "text" if first_page_text and first_page_text.strip() else "image"
    except Exception:
        return "image"

def _preprocess_image_for_ocr(self, image):
    img_cv = np.array(image)
    gray = cv.cvtColor(img_cv, cv.COLOR_BGR2GRAY)
    binary = cv.adaptiveThreshold(gray, 255, cv.ADAPTIVE_THRESH_GAUSSIAN_C,
                                  cv.THRESH_BINARY, 11, 2)

    return binary

def extract_text(self):
    if self.pdf_type == 'text':
        return self._extract_from_text_pdf()
    else:
        return self._extract_from_image_pdf()

def _extract_from_text_pdf(self):
    full_text = ""
    with pdfplumber.open(self.file_path) as pdf:
        for page in pdf.pages:

```

```

        page_text = page.extract_text()

        if page_text:

            full_text += page_text + "\n\n"

    return full_text

def _extract_from_image_pdf(self):

    try:

        images = convert_from_path(self.file_path, dpi=300)

        full_text = ""

        for image in images:

            preprocessed_image = self._preprocess_image_for_ocr(image)

            text = pytesseract.image_to_string(preprocessed_image, lang='eng')

            full_text += text + "\n\n"

        return full_text

    except Exception as e:

        return f"OCR Error: {str(e)}. Ensure Tesseract and Poppler are installed and in your system's PATH."

```

@staticmethod

```

def post_process_text(text):

    text = re.sub(r'[\n\t]+', '\n', text)

    text = re.sub(r' +', ' ', text)

    text = re.sub(r'\n\s*(\d+|Page \d+)\s*\n', '\n', text, flags=re.IGNORECASE)

    text = text.replace('fi', 'fi').replace('fl', 'fl')

    return text.strip()

```

@staticmethod

```

def evaluate_accuracy(ground_truth_text, extracted_text):

```



```

if not ground_truth_text:
    return {"error": "Ground truth text is empty."}

cer = levenshtein_distance(ground_truth_text, extracted_text) / len(ground_truth_text)

gt_words = set(ground_truth_text.split())
ext_words = set(extracted_text.split())

true_positives = len(gt_words.intersection(ext_words))
precision = true_positives / len(ext_words) if ext_words else 0
recall = true_positives / len(gt_words) if gt_words else 0
f1_score = 2 * (precision * recall) / (precision + recall) if (precision + recall) else 0

return {
    "Character Error Rate (CER)": f"{cer:.2%}",
    "Precision": f"{precision:.2%}",
    "Recall": f"{recall:.2%}",
    "F1-Score": f"{f1_score:.2%}"
}

```

### A.3 database.py

Python

# database.py

import sqlite3

import datetime

import hashlib

DB\_NAME = "extractions\_log.db"

```

def hash_password(password):

    return hashlib.sha256(password.encode()).hexdigest()


def setup_database():

    conn = sqlite3.connect(DB_NAME)

    cursor = conn.cursor()

    cursor.execute("""

        CREATE TABLE IF NOT EXISTS users (

            id INTEGER PRIMARY KEY AUTOINCREMENT,

            email TEXT UNIQUE NOT NULL,

            password_hash TEXT NOT NULL

        )

    """)

    cursor.execute("""

        CREATE TABLE IF NOT EXISTS extractions (

            id INTEGER PRIMARY KEY AUTOINCREMENT,

            user_id INTEGER NOT NULL,

            filename TEXT NOT NULL,

            extracted_text TEXT NOT NULL,

            timestamp DATETIME NOT NULL,

            FOREIGN KEY (user_id) REFERENCES users (id)

        )

    """)

    cursor.execute("SELECT * FROM users WHERE email = ?", ('test@example.com',))

    if not cursor.fetchone():

        cursor.execute(

            "INSERT INTO users (email, password_hash) VALUES (?, ?)",

            ('test@example.com', hash_password('password123'))

```

```
)  
conn.commit()  
conn.close()
```

```
def add_user(email, password):  
    conn = sqlite3.connect(DB_NAME)  
    cursor = conn.cursor()  
    try:  
        cursor.execute(  
            "INSERT INTO users (email, password_hash) VALUES (?, ?)",  
            (email, hash_password(password))  
        )  
        conn.commit()  
        return True  
    except sqlite3.IntegrityError:  
        return False  
    finally:  
        conn.close()
```

```
def verify_user(email, password):  
    conn = sqlite3.connect(DB_NAME)  
    cursor = conn.cursor()  
    cursor.execute("SELECT id, password_hash FROM users WHERE email = ?", (email,))  
    user = cursor.fetchone()  
    conn.close()  
    if user and user[1] == hash_password(password):  
        return user[0]  
    return None
```

```

def add_extraction(user_id, filename, extracted_text):
    conn = sqlite3.connect(DB_NAME)
    cursor = conn.cursor()
    timestamp = datetime.datetime.now()
    cursor.execute(
        "INSERT INTO extractions (user_id, filename, extracted_text, timestamp) VALUES (?, ?, ?, ?)",
        (user_id, filename, extracted_text, timestamp)
    )
    conn.commit()
    conn.close()

def get_user_extractions(user_id):
    conn = sqlite3.connect(DB_NAME)
    cursor = conn.cursor()
    cursor.execute(
        "SELECT id, filename, timestamp FROM extractions WHERE user_id = ? ORDER BY timestamp DESC",
        (user_id,)
    )
    records = cursor.fetchall()
    conn.close()
    return records

```

#### **A.4 utils.py**

Python

# utils.py

import pandas as pd

```

import json

from docx import Document

import io

from gtts import gTTS, gTTSError


def create_txt(text):

    return io.StringIO(text)


def create_json(source_file, text):

    data = {"source_file": source_file, "extracted_text": text}

    return io.StringIO(json.dumps(data, indent=4))


def create_csv(source_file, text):

    df = pd.DataFrame([{"source_file": source_file, "extracted_text": text}])

    output = io.StringIO()

    df.to_csv(output, index=False)

    output.seek(0)

    return output


def create_excel(source_file, text):

    output = io.BytesIO()

    df = pd.DataFrame([{"source_file": source_file, "extracted_text": text}])

    with pd.ExcelWriter(output, engine='openxl') as writer:

        df.to_excel(writer, index=False, sheet_name='Extraction')

    output.seek(0)

    return output


def create_docx(text):

```

```
document = Document()
document.add_paragraph(text)
output = io.BytesIO()
document.save(output)
output.seek(0)
return output
```

```
def create_mp3(text):
    try:
        output = io.BytesIO()
        tts = gTTS(text)
        tts.write_to_fp(output)
        output.seek(0)
        return output
    except gTTSError:
        print("gTTS connection failed. Skipping MP3 generation.")
        return None
```

## APPENDIX B: USER MANUAL

This manual provides instructions for setting up and using the PDF Text Extraction System.

### B.1 System Requirements

#### 1. Software:

- a. Python 3.7+
- b. Tesseract OCR Engine
- c. Poppler

#### 2. Hardware:

- a. Intel Core i5 or equivalent
- b. 8 GB RAM
- c. 200 MB free disk space

### B.2 Installation and Setup

- 1. Install Prerequisites:** Before running the application, you must install Python, Tesseract, and Poppler on your system and ensure they are added to your system's PATH.
- 2. Create Project Folder:** Create a folder and place all the project files (app.py, pipeline.py, etc.) inside it.
- 3. Set Up Virtual Environment:** Open a terminal in the project folder and run:  
Bash  
`python -m venv venv`
- 4. Activate Environment:**
  - a. Windows: `.\venv\Scripts\activate`
  - b. macOS/Linux: `source venv/bin/activate`

**5. Install Libraries:** Run the following command:

```
Bash  
pip install -r requirements.txt
```

### B.3 Running the Application

With the virtual environment active, run the following command in your terminal:

```
Bash  
streamlit run app.py
```

This will open the application in your web browser.

### B.4 User Guide

**1. Create an Account:**

- a. On the sidebar, select "**Register**".
- b. Enter your desired email and password and click the "**Register**" button.

**2. Login:**

- a. Select "**Login**" from the sidebar.
- b. Enter your credentials. A test account is available:

**Email:** [test@example.com](mailto:test@example.com)

**Password:** password123

**3. Upload and Process PDFs:**

- a. Once logged in, you will see the main application page.
- b. Click the "**Browse files**" button to select one or more PDF files from your computer.
- c. After selecting files, click the "**Extract Text from All Files**" button.

**4. Download Results:**

- a. The system will process all uploaded files.
- b. Once complete, a button will appear: "**Download All Extracted Files (.zip)**".
- c. Click this button to download a .zip file containing the extracted text in multiple formats



(.txt, .json, .csv, .docx, .xlsx, .mp3) for each PDF.

**5. View History:**

- a. At the bottom of the page, expand the **"Your Extraction History"** section to see a log of all the files you have processed.

## APPENDIX C: LIST OF ACRONYMS

Acronym	Full Name
<b>CER</b>	Character Error Rate
<b>CSV</b>	Comma-Separated Values
<b>DFD</b>	Data Flow Diagram
<b>ER</b>	Entity-Relationship
<b>JSON</b>	JavaScript Object Notation
<b>NLP</b>	Natural Language Processing
<b>OCR</b>	Optical Character Recognition
<b>PDF</b>	Portable Document Format
<b>UML</b>	Unified Modeling Language
<b>UI</b>	User Interface

**Table 6.4.1: List of Acronyms**