# ScanWise

**24207010 – Aayush Balip**

**24207001 – Ashokkumar Bhati**

**24207014 – Sameer Saiyed**

**24207020 – Keval Shah**

**Project Guide**

**Ms. Sarala Mary**

# Outline

- Introduction

- Literature Survey of the existing systems

- Limitations of the existing systems

- Problem statement

- System Design

- Technologies and methodologies

- Implementation

- Conclusion

- References

# Introduction

- In most universities, exam results are published as PDF documents or Excel sheets containing student details, paper codes, marks, and result status.

- These PDFs are unstructured, making it difficult for students and institutions to quickly analyze results.

- Current practice: Manual checking, downloading, and searching — a time-consuming and error-prone process.

- While some digital result portals exist, they often lack:
  - Structured data export
  - Visualization of each student's batch

# Introduction

➢ **Objectives:**

- Automate the extraction of student results from unstructured PDFs using the PyPDF library and store the data in Excel using pandas.

- Reduce manual effort and minimize the chances of data entry errors.

- Extract and restructure data from Excel sheets using openpyxl, and apply formatting with openpyxl.fonts for better readability.

- Visualize the processed data from each Excel sheet using matplotlib for clearer insights.

# Literature Survey of the existing system

| Sr. No. | Title | Author(s) | Year | Outcomes | Methodology | Result |
|---|---|---|---|---|---|---|
| 1 | Extraction of User-Defined Information from PDF | R. Nadeem, T. Iqbal, N. Fatima, J. Altaf, A. Irshad and A. Farooq | 2024 | User-defined information extracted from PDFs using PyMuPDF library for automated data retrieval from unstructured documents. | PyMuPDF text extraction techniques applied to PDFs, utilizing page.get_text() methods for native text and OCR-based extraction. | Successful extraction of textual data from PDF documents with high accuracy using PyMuPDF's hybrid native and OCR approaches. |
| 2 | Python Data Analysis: Perform data collection, data processing, wrangling, visualization, and model building using Python | Avinash Navlani; Armando Fandango; Ivan Idris | 2021 | Comprehensive Python data analysis framework covering collection, manipulation, wrangling, visualization, and model building using Pandas and Matplotlib. | Pandas DataFrames for data structures, manipulation operations like head(), tail(), slicing, and integration with visualization libraries. | Pandas DataFrames for data structures, manipulation operations like head(), tail(), slicing, and integration with visualization libraries |

# Literature Survey of the existing system

| Sr. No. | Title | Author(s) | Year | Outcomes | Methodology | Result |
|---|---|---|---|---|---|---|
| 3 | Automating Data Analysis with Python: A Comparative Study of Popular Libraries and their Application | P. Bhardwaj, C. Choudhury and P. Batra | 2024 | Comparative evaluation of Python libraries including Pandas, Matplotlib, Seaborn, Scikit-learn for automated data analysis task performance. | Performance testing and comparison of popular Python libraries across various data analysis tasks including manipulation, visualization, and modeling. | Performance testing and comparison of popular Python libraries across various data analysis tasks including manipulation, visualization, and modeling. |
| 4 | Programming language Python for data processing | Z. Dobesova | 2021 | Python established as effective programming language for data processing automation in research and academic contexts with scripting capabilities. | Used matplotlib and seaborn in Python to generate bar charts and histograms representing pass/fail courts and average performance. | Successful implementation of automated data processing workflows using Python, reducing manual intervention and improving processing speed.. |

# Limitations of existing systems

- **Lack of Automation in PDF Extraction:** Raw university mark sheets in PDF format often need to be manually converted into Excel, as many systems lack reliable PDF parsing and data normalization features.

- **Limited Cross-Semester Integration:** Existing systems usually process one semester at a time without automatically merging and averaging results across multiple semesters.

- **Absence of Visual Analytics:** Many tools fail to provide visual insights such as pass/fail distribution charts or performance trends, which are valuable for faculty review and academic reports.

# Limitations of existing systems

- **Inconsistent Data Formats:** Variation in result templates and file headers causes difficulties in data cleaning, standardization, and automation.

- **No Intelligent Highlighting or Error Detection:** Systems often miss features like automatic detection and color highlighting of failed students or invalid entries.

# Problem statement

- **High Inefficiency and Error Rate:** The current process of manually checking, downloading, and entering data from result PDFs is extremely slow and prone to human error, leading to inaccurate records and wasted administrative hours.

- **Lack of Immediate Analytics:** Students and institutions have no way to quickly generate valuable insights, such as identifying subject toppers, calculating pass/fail percentages across different papers, or tracking individual academic progress over time.

- **No Visualization of Data:** Institutions do not have any visual data to get to know about student's record and help in understanding the performance of each student batch.

# System Design



**Fig 1: System Design for Hybrid Model**

# Technologies

- **PDF to Text:** PyPDF2
- **Reading and Generating Excel:** Pandas
- **Styling Excel:** openpyxl
- **Creating Charts:** Matplotlib
- **Frontend Development:**
    i.   ReactJS
    ii.  Tailwind CSS
- **Backend Development:**
    i.   Django 4.2.5
- **Hosting:**
    i.   Vercel
    ii.  Render

# Methodology

1.  **Extraction and Reading the Data:**

    **1. PDF:** Using PyPDF, the PDF files are extracted, and text data is retrieved.

    **2. Excel:** Using openpyxl, Excel sheets are opened and read.

2. **Processing the data:**

    **1. PDF:** After extracting the text from PDFs, regex operations are applied to process and extract specific data.

    **2. Excel:** After reading the Excel sheet, pandas is used to access specific columns and perform the desired operations.

3. **Result:**

    **1. PDF:** The processed data is stored in Excel sheets using pandas.

    **2, Excel:** The final results are saved to Excel sheets, and visualizations are created using matplotlib.

# Implementation

# Conclusion

- Developed an automated system to extract, process, and visualize student exam results from PDFs and Excel files.

- Utilized **PyPDF2**, **pandas**, **openpyxl**, and **matplotlib** to eliminate manual errors and inefficiencies.

- Implemented a **ReactJS + Tailwind CSS frontend** with a **Django backend**, hosted on **Vercel** and **Render** for scalability.

- Provides quick access to structured data, performance metrics, and visual analytics.

- Significantly reduces administrative workload, enhances data accuracy, and supports informed academic decision-making.

# References

[1] R. Nadeem, T. Iqbal, N. Fatima, J. Altaf, A. Irshad and A. Farooq, "Extraction of User-Defined Information from PDF," doi: 10.1109/DASA63652.2024.10836169.

[2] Avinash Navlani; Armando Fandango; Ivan Idris, Python Data Analysis: Perform data collection, data processing, wrangling, visualization, and model building using Python , Packt Publishing, 2021. http://ieeexplore.ieee.org/document/10163464

[3] P. Bhardwaj, C. Choudhury and P. Batra, "Automating Data Analysis with Python: A Comparative Study of Popular Libraries and their Application," doi: 10.1109/ICTACS59847.2023.10390032.

[4] Z. Dobesova, "Programming language Python for data processing," doi: 10.1109/ICECENG.2011.6057428.

# Thank You...!!