

# Technical Documentation for Invoice Data Extraction System

## Detailed Explanation of the Approach and Algorithms Used

---

### Approach Overview

The Invoice Data Extraction System is designed to process PDF invoices, extracting key information such as invoice numbers, dates, item details, and amounts. The approach integrates various libraries and techniques to ensure robustness and accuracy. The core components include:

- **Text Extraction:** Utilizing both structured text extraction (for digitally created PDFs) and Optical Character Recognition (OCR) for scanned documents.
- **Data Validation:** Implementing regex patterns and conditional checks to verify the accuracy of the extracted data.
- **Batch Processing:** Using multithreading to enhance performance and efficiency during the extraction process.

### Algorithms and Techniques Used

#### Text Extraction

- **pdfplumber:** This library is employed to extract text from PDF files directly when the documents are digitally generated. It provides reliable access to the text content and layout information.
- **pytesseract:** For scanned invoices where text extraction may not be straightforward, OCR is performed using pytesseract, which converts images of text into machine-encoded text.

**Regular Expressions (Regex)** Regex patterns are constructed to identify and extract specific fields from the text. The patterns are designed to accommodate variations in formatting across different invoices.

**Data Validation** Each extracted field undergoes a validation process to confirm its correctness. Validation includes:

- **Format checks:** Ensuring dates match the dd Month yyyy format.
- **Presence checks:** Ensuring mandatory fields are not empty.
- **Consistency checks:** Between related fields (e.g., due dates after invoice dates).

**Multithreading** The `concurrent.futures.ThreadPoolExecutor` is used to process multiple PDFs simultaneously, significantly reducing the time required for batch processing of invoices.

**Reporting** The system generates an extraction report, detailing the accuracy of the data captured and the number of invoices processed.

## Justification for Chosen Methods

---

### Cost-Effectiveness vs. Accuracy

**Cost-Effectiveness** The selected libraries (`pdfplumber`, `pytesseract`, and `pandas`) are open-source and widely used, minimizing licensing costs and providing a solid community support base. Additionally, the use of multithreading reduces processing time, allowing for more invoices to be handled without requiring additional computational resources.

**Accuracy** The combination of structured extraction and OCR ensures that the system can handle both types of documents effectively, maximizing the amount of usable data. By implementing strict validation criteria, the system can filter out inaccurate or incomplete data, enhancing overall reliability. The iterative improvement approach, supported by detailed logging of errors, allows for continuous refinement of the extraction and validation processes, which helps in maintaining a high accuracy level.

### Balance

The balance between cost-effectiveness and accuracy is achieved through the following strategies:

- **Robustness of Algorithms:** Using powerful extraction and validation techniques ensures high-quality data without incurring significant additional costs.
- **Efficiency in Processing:** Multithreading allows the system to maintain high throughput, enabling the processing of larger volumes of invoices within the same time frame, thus maximizing resource utilization.

## Specific Explanation of the Method Used to Achieve the 99% Trust Determination Requirement

---

To ensure a 99% trust determination requirement, the following specific methods were employed:

### Comprehensive Field Extraction

Each invoice is processed to extract critical fields using robust regex patterns. The fields extracted include Invoice Number, Invoice Date, Due Date, Taxable Amount, and Total Amount.

## Rigorous Validation Framework

Each extracted field is validated based on specific criteria:

- **Presence Checks:** Ensure that essential fields are not empty.
- **Format Checks:** Use regex to confirm that dates and amounts follow the expected formats.
- **Value Checks:** Verify that amounts are non-negative and correctly formatted.

## Trust Level Assignment

A trust level is assigned based on the validation results. If all critical fields meet the validation criteria, the invoice data is labeled as “Trusted.” Conversely, any failure in validation results in a “Not Trusted” label.

## Continuous Improvement through Feedback

The system logs any extraction errors, enabling developers to review and refine regex patterns and validation checks continuously. This iterative approach is critical in enhancing the accuracy of future extractions.

## Accuracy Metrics Tracking

The system keeps a count of total attempts and successful validations for each field. The trust level is determined by calculating the percentage of correct extractions against total attempts, ensuring that at least 99 out of every 100 invoices processed meet the trust criteria.

## Final Reporting

At the end of the extraction process, a detailed report is generated, summarizing the accuracy metrics, trust levels, and any anomalies encountered during processing. This report serves as a tool for stakeholders to assess data reliability and make informed decisions.

## Conclusion

---

The Invoice Data Extraction System combines effective algorithms, rigorous validation methods, and continuous feedback mechanisms to achieve a 99% trust determination requirement. This comprehensive approach ensures high-quality data extraction while maintaining cost-effectiveness, thereby supporting reliable business processes.