

End-to-End Digitization of Retail Receipts using a Hybrid AI Pipeline

Bhagwan Arsewad

Indian Institute of Technology Jodhpur

Department of Computer Science and Engineering

Jodhpur, India

B22AI010@iitj.ac.in

Dr. Bikash Santra

Indian Institute of Technology Jodhpur

School of Artificial Intelligence and Data Science

Jodhpur, India

bikash@iitj.ac.in

Abstract

This report details the successful implementation of a hybrid AI pipeline for retail receipt digitization. The system's objective is to extract all text from a receipt image and then parse that text to identify and categorize key semantic information (e.g., company, date, total).

The final system uses a 2-step "Scanner-Accountant" approach: (1) A YOLOv8 object detection model trained on the SROIE dataset to find all text blocks, and (2) A robust Python parser that reconstructs the bill's layout and uses advanced logic (regex) to extract key-value pairs.

The AI "Scanner" model proved to be extremely successful, achieving a **mean Average Precision (mAP50) of 96.4%** on the validation set. This confirms the model's ability to find text with high accuracy. The key finding of this project is that the primary challenge is not AI-based detection, but the rule-based parsing logic. An initial simple parser yielded low end-to-end accuracy (30-50%). By developing an upgraded, intelligent parser (v4.0), the system's performance was dramatically improved to over 89% accuracy across all key fields, proving the viability and success of the hybrid pipeline.

CCS Concepts

- Computing methodologies → Artificial intelligence; Computer vision;
- Applied computing → Document analysis.

Keywords

Receipt Digitization, YOLOv8, OCR, Information Extraction, Computer Vision

ACM Reference Format:

Bhagwan Arsewad and Dr. Bikash Santra. 2025. End-to-End Digitization of Retail Receipts using a Hybrid AI Pipeline. In . ACM, New York, NY, USA, 4 pages.

1 Introduction

The digitization of retail receipts represents a significant challenge in the field of document understanding and information extraction [2]. Manual data entry from physical receipts is not only time-consuming and expensive but also prone to human error. Automated

systems for receipt processing have substantial applications in expense management, accounting automation, and retail analytics.

1.1 Objective

The primary objective of this project is to design, build, and evaluate an automated system that can process raw images of retail bills and produce two distinct, high-value outputs:

- (1) **A structured JSON summary** containing key semantic information (Company, Date, Total Amount, Items).
- (2) **A "like-to-like" text reconstruction** that preserves the original bill's spatial layout and formatting.

1.2 Dataset Selection: CORD vs SROIE

A critical aspect of this project involved selecting an appropriate dataset that would enable effective model training and evaluation.

- **CORD Dataset (Rejected):** Despite its large size, the CORD dataset was deemed unsuitable due to extensive privacy blurring applied to critical fields such as `store_name` and `address`. This preprocessing made it impossible to train models for extracting these essential semantic fields.
- **SROIE Dataset (Selected):** The SROIE (Scanned Receipts OCR and Information Extraction) dataset emerged as the ideal choice [2]. It comprises 626 high-resolution, clean receipt scans with fully visible text. The dataset provides comprehensive annotations perfectly suited for our two-stage approach:

- (1) **box annotations:** Word-level bounding box coordinates for comprehensive text localization.
- (2) **entities annotations:** Key-value pair ground truth for semantic field extraction.

1.3 System Architecture: The "Scanner-Accountant" Model

We designed and implemented a hybrid 2-step pipeline that separates the tasks of visual detection and semantic understanding. This architectural decision enhances system accuracy, flexibility, and debuggability.

- **Step 1: The "Scanner" (AI Model):** A YOLOv8 model [4] trained with a single class ("text") dedicated exclusively to locating every text element within the receipt image. This focused

approach allows the model to achieve exceptional detection accuracy.

- **Step 2: The "Accountant" (Python Parser):** The extractor.py module processes the detected text regions through three sequential operations:
 - (1) **Recognition:** Utilizes EasyOCR [1] for optical character recognition within each detected bounding box.
 - (2) **Reconstruction:** Implements a custom reconstruct_lines function that groups text elements based on vertical alignment to reconstruct the original receipt layout.
 - (3) **Parsing:** Employs sophisticated Regular Expressions and rule-based logic in parse_extracted_text to identify and categorize key information.

2 Methodology

2.1 Data Preparation Pipeline

The SROIE dataset required substantial preprocessing to conform to YOLOv8's annotation format. We developed an automated Python script (`01_create_sroie_labels.py`) that:

- Processes word-level coordinate annotations from the train/box directory
- Generates corresponding .txt label files with class ID 0 for all text elements
- Implements a 90%-10% random split (563 training, 63 validation images)
- Creates the necessary data.yaml configuration file

2.2 Model Training Configuration

The training phase employed the following configuration:

- **Model Architecture:** YOLOv8n (nano variant) pre-trained on COCO dataset
- **Training Environment:** Google Colab with Tesla T4 GPU acceleration
- **Training Parameters:** 300 epochs with early stopping (patience=100)
- **Final Model:** best_text_detector.pt saved from epoch 129

2.3 Parsing Logic Evolution

The development of the parsing logic underwent significant iterative improvement:

- **v1.0 (Naive Parser):** Implemented basic heuristics:
 - Company identification from first text line
 - Total amount detection using "TOTAL" keyword matching
 - Date extraction using simple dd/mm/yy pattern matching
- **v4.0 (Intelligent Parser):** Enhanced with sophisticated processing:
 - Advanced date regex patterns (e.g., \d{1,2}\s+[A-Za-z]{3},\s+\d{4})
 - Total failed because it only looked for the "TOTAL" keyword.

- Comprehensive keyword lists for total detection (TOTAL, AMT, CASH)
- Look-ahead validation for amount fields
- Intelligent company identification with junk line filtering
- OCR error correction heuristics (character substitution mapping)

3 Results and Discussion

3.1 Result 1: AI "Scanner" Model Performance

The YOLOv8 text detection model demonstrated exceptional performance on the validation set, as summarized in Table 1.

Table 1: AI Model (Scanner) Validation Performance Metrics

Metric	Score	Percentage
mAP50	0.964	96.4%
mAP50-95	0.688	68.8%
Precision	0.960	96.0%
Recall	0.925	92.5%

Discussion: The achieved mAP50 score of 96.4% indicates outstanding detection capability. This performance level confirms that the computer vision component is not the system's limiting factor and validates the effectiveness of the single-class detection approach for text localization in receipt images.

3.2 Result 2: End-to-End System "Accountant" Performance

This is the true test of the entire system. An evaluation script (`run_evaluation.py`) was run on the **347 unseen images in the test set**. This script runs the full pipeline (Scanner + Accountant) and compares the final parsed JSON to the ground truth labels.

An initial test with a simple parser (v1.0) yielded the poor results seen in Table 2.

Table 2: End-to-End System Accuracy (v1.0 Simple Parser)

Key Information	Accuracy
Company Name	33.14% (115 / 347)
Date	59.94% (208 / 347)
Total Amount	47.55% (165 / 347)

Discussion of Findings: The low scores in Table 2 revealed the system's true bottleneck: the **"parsing logic"**. The 96.4% accurate AI model was finding the text, but the simple parser was not smart enough to understand it.

- **Company** failed because the "first line" rule is not always true.
- **Date** failed because it only looked for one date format.

To solve this, the parser was upgraded to version 4.0 with smarter regex and logic. The evaluation was run again. The final, much-improved scores (shown in Table 3) prove the success of the upgraded system.

Table 3: Final End-to-End System Accuracy (v4.0 Smart Parser)

Key Information	Final Accuracy
Company Name	89.34% (294 / 347)
Date	92.22% (320 / 347)
Total Amount	94.52% (328 / 347)

Discussion: The dramatic improvement in accuracy across all key fields demonstrates the effectiveness of the enhanced parsing logic. The v4.0 parser successfully addresses the limitations of the initial approach through:

- **Multi-format Date Recognition:** Handling various date formats beyond simple dd/mm/yy
- **Contextual Company Detection:** Intelligent filtering of non-company text lines
- **Robust Total Extraction:** Comprehensive keyword matching and look-ahead validation
- **OCR Error Correction:** Character substitution to handle common OCR mistakes

The final system achieves over 89% accuracy across all key semantic fields, proving the viability of the hybrid AI pipeline approach for retail receipt digitization.

4 Conclusion and Future Work

4.1 Project Summary

This project successfully demonstrates the effectiveness of a hybrid AI pipeline for retail receipt digitization. The "Scanner-Accountant" architecture proves particularly advantageous by:

- Achieving near-perfect text detection accuracy (96.4% mAP) through specialized YOLOv8 training
- Providing clear separation of concerns between detection and understanding tasks
- Enabling iterative improvement of parsing logic without re-training the vision model
- Facilitating system debugging and performance analysis

The key insight from this work is that the primary challenge in receipt information extraction lies not in computer vision-based detection, but in the domain-specific parsing and understanding logic. The iterative refinement of the parsing component from v1.0 to v4.0 resulted in dramatic improvements in end-to-end system accuracy, from approximately 47% to over 89% across key fields.

4.2 Future Work

Several promising directions emerge for extending this research:

- **Enhanced Item Parsing:** Develop sophisticated algorithms for reliably extracting complete line items (product names, quantities, individual prices) and handling multi-line product descriptions.
- **Domain Adaptation:** Address the domain shift problem by training on more diverse datasets (including CORD for real-world photos) and implementing data augmentation strategies for challenging conditions (blur, shadows, wrinkles).
- **End-to-End Architectures:** Explore modern transformer-based models (LayoutLM [5], Donut [3]) that unify detection and understanding in a single architecture, potentially simplifying the pipeline while maintaining accuracy.
- **Multi-modal Processing:** Incorporate additional signal sources such as layout analysis, font characteristics, and spatial relationships to improve parsing robustness.
- **Real-time Processing:** Optimize the pipeline for mobile deployment and real-time processing applications.

Acknowledgments

We would like to express our sincere gratitude to our project supervisor, Dr. Bikash Santra, for invaluable guidance, support, and encouragement throughout this B.Tech Project. The insights and feedback provided were instrumental in the development and successful completion of this work.

We also thank the Department of Computer Science and Engineering at IIT Jodhpur for providing the necessary resources and a conducive environment for research.

A Appendix

A.1 Project Repository

The complete source code, including data preprocessing scripts, training implementation, evaluation framework, and the Streamlit web application, is available at:

<https://github.com/bhagwan388/btp-receipt-extractor>

A.2 Technology Stack

- **Computer Vision:** YOLOv8 (Ultralytics) [4] for text detection
- **OCR Engine:** EasyOCR (JaideAI) [1] for text recognition
- **Dataset:** SROIE (ICDAR 2019 Competition Dataset) [2]
- **Web Framework:** Streamlit for application deployment
- **Development:** Python 3.10, Google Colab, OpenCV

References

- [1] JaideAI. 2020. EasyOCR: Ready-to-use OCR with 80+ supported languages. <https://github.com/JaideAI/EasyOCR>.
- [2] Zheng Huang, Kai Chen, Jianhua He, Xiang Bai, Dimosthenis Karatzas, Shijian Lu, and CV Jawahar. 2019. ICDAR2019 Competition on Scanned Receipts OCR and Information Extraction. In *2019 International Conference on Document Analysis and Recognition (ICDAR)*. IEEE, 1516–1520.

- [3] Geewook Kim, Teakgyu Hong, Moonbin Yim, Jeongyeon Nam, Jinyoung Park, Jinyoung Yim, Wonseok Hwang, Sangdoo Yun, Dongyoon Han, and Seunghyun Park. 2022. Donut: Document understanding transformer without OCR. *arXiv preprint arXiv:2111.15664* (2022).
- [4] Ultralytics. 2023. YOLOv8: The latest version of the award-winning YOLO model. <https://github.com/ultralytics/ultralytics>.
- [5] Yiheng Xu, Minghao Li, Lei Cui, Shaohan Huang, Furu Wei, and Ming Zhou. 2020. Layoutlm: Pre-training of text and layout for document image understanding. In *Proceedings of the 26th ACM SIGKDD International Conference on Knowledge Discovery & Data Mining*. 1192–1200.