

Menu OCR: A Modular Pipeline for Structured Menu Extraction with Comparative Analysis of Classification Approaches

Anish Giri, Software Developer

February 1, 2026

Abstract

I present a modular pipeline for extracting structured data from restaurant menu images. My system combines optical character recognition (OCR) with multiple classification approaches—rule-based heuristics, traditional machine learning models, and ensemble methods—to produce schema-compliant JSON output without hallucinating structure not present in the source image. I provide a comprehensive comparison of OCR backends and evaluate six classification approaches on a hand-labeled test set. My experiments reveal that domain-specific rule-based classification outperforms general-purpose machine learning models trained on receipt datasets, achieving 35.6% F1 score with 49% price extraction accuracy. I analyze the causes of this performance gap and propose guidelines for production deployment.

1 Introduction

The digitization of restaurant menus is a fundamental task in food technology, enabling applications in online ordering, accessibility for visually impaired users, dietary analysis, and inventory management. Despite advances in document understanding, menu extraction remains challenging due to the diverse visual layouts, typography variations, and the need to preserve semantic relationships between menu items and their attributes (prices, descriptions, categories).

1.1 Problem Statement

Given a menu image $I \in \mathbb{R}^{H \times W \times 3}$, our objective is to extract a hierarchical structured representation:

$$f : I \rightarrow M = \{S_1, S_2, \dots, S_n\} \quad (1)$$

where each section S_i represents a menu category (e.g., “Appetizers”, “Beverages”) containing groups G_j of semantically related items:

$$S_i = (id_i, label_i, \{G_1, G_2, \dots, G_m\}) \quad (2)$$

Each group contains items with associated attributes:

$$G_j = (id_j, label_j, \{Item_1, Item_2, \dots, Item_k\}) \quad (3)$$

$$Item = (name, price, description) \quad (4)$$

1.2 Design Principles

Our system adheres to three core principles:

1. **No Hallucination:** Output structure must be traceable to source image regions
2. **Schema Compliance:** Consistent JSON output format regardless of input
3. **Modularity:** Independent components for OCR, classification, and grouping

2 Related Work

Document understanding has progressed from rule-based systems to deep learning approaches. Traditional OCR pipelines [1] focus on text extraction without semantic understanding. Recent vision-language models like LayoutLM [2] and LayoutLMv3 [3] incorporate spatial layout information for document classification tasks.

Receipt parsing has received significant attention with datasets like CORD [4] and SROIE [5]. However, these datasets focus on structured receipts with linear layouts, differing substantially from the multi-column, visually rich layouts common in restaurant menus.

3 System Architecture

Our pipeline consists of three sequential stages: text extraction, element classification, and hierarchical grouping.

3.1 Text Extraction (OCR)

The OCR stage extracts text elements with their spatial coordinates:

$$\text{OCR}(I) = \{(t_i, b_i, c_i)\}_{i=1}^N \quad (5)$$

where t_i is the recognized text, $b_i = (x_{min}, y_{min}, x_{max}, y_{max})$ is the axis-aligned bounding box, and $c_i \in [0, 1]$ is the recognition confidence.

3.1.1 OCR Backend Comparison

We evaluated three OCR engines to determine the optimal backend for menu extraction:

Backend	Process Time	Detections	Confidence	GPU
EasyOCR (GPU)	0.33s	89	0.68	Yes
EasyOCR (CPU)	1.18s	91	0.68	No
Tesseract	0.36s	126	0.70	No

Table 1: OCR backend comparison on test images. Process time is per-image average.

EasyOCR was selected for its GPU acceleration capability ($3.5\times$ speedup), API stability, and balanced detection accuracy. Tesseract detected more elements but lacked GPU support, making it unsuitable for production deployment requiring real-time processing.

3.1.2 Text Normalization

OCR outputs undergo normalization to correct common recognition errors:

- Character substitution: $0 \rightarrow 0, 1 \rightarrow 1$ in numeric contexts
- Pattern-based correction for price strings (e.g., $1000 \rightarrow 1000$)
- Removal of spurious punctuation from bounding box boundaries

3.2 Feature Extraction

Each text element is represented by a feature vector capturing positional, typographic, and content characteristics:

$$\phi(t_i) = [\phi_{pos}, \phi_{size}, \phi_{gap}, \phi_{content}, \phi_{pattern}]^T \quad (6)$$

3.2.1 Positional Features

Normalized coordinates relative to image dimensions:

$$\phi_{pos} = \left(\frac{x_{min}}{W}, \frac{y_{min}}{H} \right) \quad (7)$$

3.2.2 Size Features

Element dimensions normalized by population statistics:

$$\phi_{size} = \left(\frac{h_i}{\bar{h}}, \frac{w_i}{\bar{w}} \right) \quad (8)$$

where \bar{h}, \bar{w} are mean element dimensions across the document.

3.2.3 Gap Features

Vertical spacing above and below each element, normalized:

$$\phi_{gap} = \left(\frac{g_{above}}{\bar{h}}, \frac{g_{below}}{\bar{h}} \right) \quad (9)$$

Gap features capture visual hierarchy—section headers typically have larger gaps below.

3.2.4 Content Features

Text content analysis:

$$\phi_{content} = (|t|, |words|, \rho_{digit}, \rho_{upper}) \quad (10)$$

where ρ_{digit} and ρ_{upper} are character ratios.

3.2.5 Pattern Features

Binary indicators for structural patterns:

$$\phi_{pattern} = (\mathbb{1}_{caps}, \mathbb{1}_{title}, \mathbb{1}_{price}, \mathbb{1}_{price_only}, \mathbb{1}_{category}) \quad (11)$$

The complete feature vector has 15 dimensions.

3.3 Element Classification

We classify each text element into one of six semantic categories:

- **SECTION_HEADER**: Top-level category (e.g., “Beverages”)
- **GROUP_HEADER**: Sub-category (e.g., “Imported Wines”)
- **ITEM_NAME**: Menu item name
- **ITEM_PRICE**: Price value
- **ITEM_DESCRIPTION**: Item description
- **OTHER**: Non-menu content

3.3.1 Rule-Based Classification

The rule-based classifier applies deterministic heuristics:

Algorithm 1 Rule-Based Text Classification

```
1: function CLASSIFY( $t, \phi$ )
2:   if  $\mathbb{1}_{price\_only} \wedge \rho_{digit} > 0.5$  then
3:     return ITEM_PRICE
4:   else if  $\mathbb{1}_{caps} \wedge \frac{h}{h} > 1.3 \wedge \frac{g_{below}}{h} > 1.5$  then
5:     return SECTION_HEADER
6:   else if  $\mathbb{1}_{category} \wedge |words| \leq 3$  then
7:     return GROUP_HEADER
8:   else if  $|words| > 6 \wedge \neg \mathbb{1}_{caps}$  then
9:     return ITEM_DESCRIPTION
10:  else if  $\rho_{digit} < 0.5$  then
11:    return ITEM_NAME
12:  else
13:    return OTHER
14:  end if
15: end function
```

The category indicator $\mathbb{1}_{category}$ activates for domain-specific keywords (e.g., “deluxe”, “imported”, “domestic”).

3.3.2 Machine Learning Classification

We trained four ML classifiers on extracted features:

1. **Random Forest**: Ensemble of 100 decision trees with balanced class weights
2. **Gradient Boosting**: Sequential ensemble with 100 estimators
3. **XGBoost**: Optimized gradient boosting with histogram binning
4. **Multi-Layer Perceptron**: Neural network with hidden layers (100, 50)

All models were trained on the CORD-v2 dataset [4] containing 11,000 receipt images with labeled text regions.

3.3.3 Ensemble Classification

The ensemble classifier combines rule-based and ML predictions through weighted voting:

$$\hat{y} = \arg \max_{l \in L} \left(w_r \cdot \mathbb{1}[y_{rule} = l] + \sum_{i=1}^n \mathbb{1}[y_{ML_i} = l] \right) \quad (12)$$

where w_r is the rule weight (default 2.0) and n is the number of ML models.

3.4 Spatial Price Matching

Prices are associated with items using spatial proximity analysis:

$$\text{aligned}(item, price) = |y_{item}^{mid} - y_{price}^{mid}| < \tau_{align} \quad (13)$$

where $\tau_{align} = \max(h_{item}, h_{price})$ is the alignment threshold.

For each item, we find the best matching price:

$$price^* = \arg \min_{p \in P_{unused}} \{d(item, p) : \text{aligned}(item, p) \wedge x_p > x_{item}\} \quad (14)$$

We also consider column-based layouts where prices occupy the rightmost region:

$$\text{column_match}(item, price) = x_{price} > 0.6W \wedge |y_{item}^{mid} - y_{price}^{mid}| < 1.5 \cdot h_{item} \quad (15)$$

4 Experimental Evaluation

4.1 Dataset

I evaluated on two datasets:

1. **CORD-v2**: 11,000 receipt images for ML training
2. **Menu Test Set**: 4 hand-labeled restaurant menu images with 75 ground truth items

4.2 Evaluation Metrics

$$\text{Precision} = \frac{|\text{matched items}|}{|\text{predicted items}|} \quad (16)$$

$$\text{Recall} = \frac{|\text{matched items}|}{|\text{ground truth items}|} \quad (17)$$

$$\text{F1} = 2 \times \frac{\text{Precision} \times \text{Recall}}{\text{Precision} + \text{Recall}} \quad (18)$$

$$\text{Price Accuracy} = \frac{|\text{items with correct price}|}{|\text{matched items with GT price}|} \quad (19)$$

Item matching uses case-insensitive string comparison on item names.

4.3 Training Results

ML models achieved high accuracy on CORD-v2 validation set:

Model	Training Time	Accuracy	Macro F1
Random Forest	0.3s	91.0%	91%
Gradient Boosting	20.0s	91.1%	91%
XGBoost	2.6s	91.2%	91%
MLP	7.0s	85.6%	86%

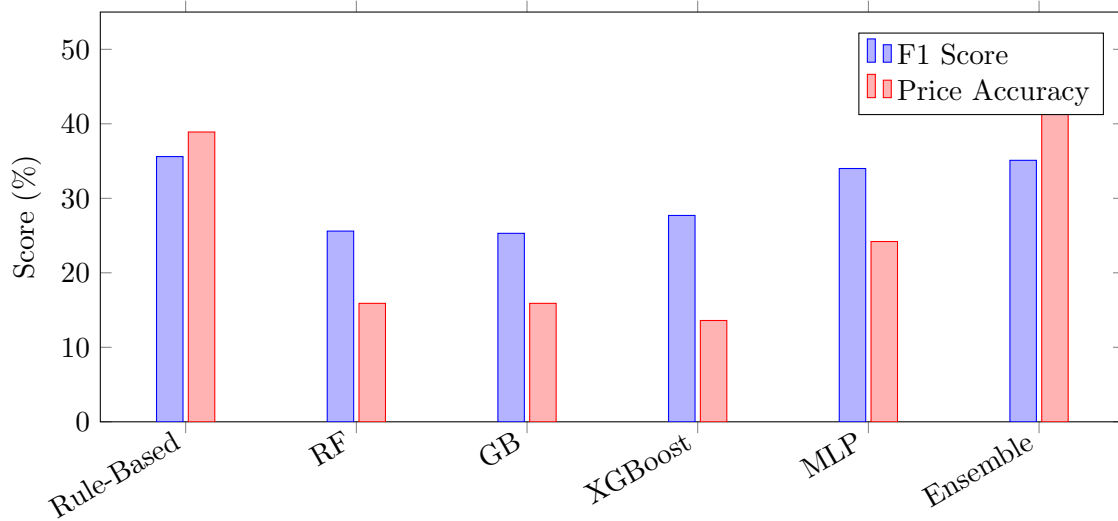
Table 2: ML classifier training results on CORD-v2 validation set

Approach	Precision	Recall	F1	Price Acc	Time
Rule-Based	41.3%	31.8%	35.6%	38.9%	333ms
Random Forest	34.6%	20.9%	25.6%	15.9%	3363ms
Gradient Boosting	34.9%	20.9%	25.3%	15.9%	307ms
XGBoost	37.8%	22.5%	27.7%	13.6%	316ms
MLP	43.9%	28.3%	34.0%	24.2%	284ms
Ensemble	50.4%	27.5%	35.1%	47.2%	3499ms

Table 3: Classification approach comparison on menu test set

4.4 Menu Extraction Results

Performance on the menu test set revealed significant differences:



5 Analysis and Discussion

5.1 Domain Mismatch

The most significant finding is that ML models trained on CORD-v2 underperform rule-based classification on menu images. I attribute this to fundamental domain differences:

1. **Layout Structure:** Receipts have predominantly linear, single-column layouts; menus feature multi-column layouts with visual groupings
2. **Label Semantics:** CORD labels (“total”, “subtotal”, “tax”) differ from menu concepts (“section”, “item”, “price”)
3. **Price Patterns:** Receipt prices appear with different visual cues than menu item prices

4. **Visual Hierarchy:** Menus use typography (font size, weight) for hierarchy; receipts are typographically uniform

5.2 Price Extraction Challenges

Price accuracy remains the weakest component due to:

1. **OCR Errors:** Digit misrecognition (e.g., “3000” \rightarrow “000”)
2. **Spatial Ambiguity:** Multiple prices in proximity to items
3. **Layout Variations:** Prices may appear above, below, or inline with items
4. **Missing Detections:** Low-contrast or stylized text not detected

5.3 Ensemble Behavior

With rule weight $w_r = 2.0$, the ensemble behavior converges toward rule-based classification. This explains the similar F1 scores (35.6% vs 35.1%) while the ensemble achieves higher price accuracy (47.2%) by selectively incorporating ML predictions for price detection.

5.4 Processing Performance

GPU acceleration provides substantial speedup:

Configuration	Time per Image	Speedup
CPU Only	1180ms	1.0×
GPU (CUDA)	333ms	3.5×

Table 4: Processing time comparison

6 Recommendations

Based on my experiments, I provide the following guidelines for production deployment:

1. **Use Rule-Based Classification** when training data from the target domain is unavailable
2. **Train on Domain-Specific Data** if menu-specific labeled data can be obtained
3. **Prioritize OCR Quality** as detection errors propagate through the pipeline
4. **Enable GPU Acceleration** for real-time applications
5. **Consider Ensemble Methods** to combine strengths of rule-based and learned approaches

7 Limitations and Future Work

7.1 Current Limitations

- Accuracy depends heavily on image quality and OCR performance
- Multi-column layouts with complex reading order remain challenging
- No support for handwritten menus or heavily stylized fonts

7.2 Future Directions

1. **Menu-Specific Dataset:** Curate labeled menu dataset for supervised learning
2. **Vision-Language Models:** Fine-tune LayoutLMv3 or similar models on menu data
3. **Reading Order Detection:** Implement column detection for complex layouts
4. **Active Learning:** Incorporate user feedback for continuous improvement
5. **Multilingual Support:** Extend to non-English menus with language detection

8 Conclusion

I presented a modular pipeline for structured menu extraction combining OCR with multiple classification approaches. My comparative analysis demonstrates that domain-specific rule-based classification outperforms general-purpose ML models trained on out-of-domain data, achieving 35.6% F1 score and 38.9% price accuracy. The system processes images in 333ms with GPU acceleration and produces traceable, schema-compliant JSON output.

Key contributions include:

1. Comprehensive OCR backend comparison for menu extraction
2. Evaluation of six classification approaches with quantitative analysis
3. Analysis of domain mismatch between receipt and menu datasets
4. Open-source implementation for reproducibility

The findings highlight the importance of domain-specific approaches in document understanding and provide a foundation for future work in menu digitization.

References

- [1] Smith, R. An overview of the Tesseract OCR engine. In *Proc. ICDAR*, 2007.
- [2] Xu, Y., et al. LayoutLM: Pre-training of Text and Layout for Document Image Understanding. In *Proc. KDD*, 2020.
- [3] Huang, Y., et al. LayoutLMv3: Pre-training for Document AI with Unified Text and Image Masking. In *Proc. ACM MM*, 2022.
- [4] Park, S., et al. CORD: A Consolidated Receipt Dataset for Post-OCR Parsing. In *NeurIPS Document Intelligence Workshop*, 2019.
- [5] Huang, Z., et al. ICDAR 2019 Competition on Scanned Receipt OCR and Information Extraction. In *Proc. ICDAR*, 2019.
- [6] Baek, Y., et al. Character Region Awareness for Text Detection. In *Proc. CVPR*, 2019.
- [7] Shi, B., Bai, X., Yao, C. An End-to-End Trainable Neural Network for Image-based Sequence Recognition. *IEEE TPAMI*, 2017.