
Object Detection and Price Prediction of School Cafeteria by YOLO-v8

Shang-En Tsai

Department of Economics
National Taiwan University
B09303052

Zhen-Yan Chen

Institute of Statistics and Data Science
National Taiwan University
R11250023

Wen-Kai Su

Institute of Applied Mathematical Sciences
National Taiwan University
R12246001

Abstract

1 This project presents a food detection and price prediction system for cafeteria
2 trays using computer vision techniques. The dataset comprises 330 labeled images,
3 including manually collected data and web-sourced images, with labels encompassing
4 common food items such as fried chicken, cabbage, and tofu. Two models,
5 YOLO-v8 [2] and Faster R-CNN [3], were implemented for object detection, with
6 YOLO-v8 achieving good performance and efficiency. Price calculation is based on
7 the estimated weight of food items occupying specific areas on the tray, with cost
8 derived from a visual weight approximation method. While the system performs
9 well for limited food categories, challenges such as uneven tray placement and
10 certain food items (e.g., hash browns) require improved counting methodologies.
11 Future work will expand food categories, integrate nutritional analysis, and explore
12 regression models for enhanced pricing accuracy. The solution demonstrates significant
13 potential for use in buffet-style cafeterias, offering accurate food recognition
14 and pricing insights.

15

1 Introduction

16 In recent years, advancements in computer vision and deep learning have opened opportunities to
17 automate various real-world tasks, including food recognition and price estimation. In buffet-style
18 cafeterias, calculating food costs accurately remains a challenge, especially in school cafeterias of
19 National Taiwan University as pricing only depends on weight of food. This project develops a food
20 detection and price prediction tool that utilizes deep learning models to recognize food items on a
21 cafeteria tray, estimate their weight, and compute the total cost.

22 The primary objectives of this project are as follows:

- 23 1. **Food Recognition:** Implement a model capable of accurately identifying and classifying multiple
24 food items on a tray.
- 25 2. **Price Prediction:** Develop a method to approximate the weight of food items based on the area
26 they occupy and calculate the corresponding cost.

27 To achieve these goals, two state-of-the-art object detection models, YOLO-v8 and Faster R-CNN,
28 were utilized and compared for performance. The dataset consists of 330 labeled images, including
29 custom images collected from 9th Women's Dorm Restaurant of National Taiwan University and

30 web-sourced images. Data augmentation techniques, such as varying angles, lighting conditions, and
31 food placement, were applied to improve model generalization.

32 The proposed system not only automates the food cost estimation process but also introduces a
33 reliable solution for restaurants to streamline pricing in buffet scenarios. By ensuring accuracy and
34 scalability, this tool has significant potential for implementation in cafeteria environments, enhancing
35 efficiency and reducing human error.

36 The following sections will discuss the dataset, methodology, results, and conclusion.

37 **Related work** Food recognition systems often utilize object detection frameworks, such as the
38 YOLO series and Faster R-CNN, to achieve real-time detection and classification. Wu et al. (2019)
39 applied YOLOv3 to optimize cafeteria checkout processes by recognizing food items on trays and
40 calculating their price, achieving practical mAP scores despite environmental variations like tray
41 placement and lighting conditions [5]. This work demonstrated the potential of YOLO models in
42 real-world cafeteria settings but highlighted the challenges of maintaining accuracy under inconsistent
43 conditions.

44 To address the need for more precise weight estimation, Gonzalez et al. (2024) combined RGB and
45 depth information to estimate food volumes and weights, achieving error margins as low as 3.75% for
46 specific food items like chicken [1]. By leveraging the YOLO architecture for content detection and
47 integrating density-based models for weight prediction, this approach underscored the importance
48 of combining visual detection with volumetric measurements. However, reliance on depth cameras
49 increases hardware complexity and cost, making it less accessible for some use cases.

50 In contrast, Wimalasiri et al. (2024) demonstrated that accurate weight estimation can also be achieved
51 using only 2D images [4]. By employing Faster R-CNN for food detection and MobileNetV3 for
52 weight prediction, they achieved high precision with an R-squared value of 98.65%. This method
53 highlights the feasibility of extracting features such as image area, aspect ratio, and pixel intensity for
54 weight estimation, offering a cost-effective alternative to depth-based techniques while maintaining
55 robust performance.

56 2 Data Description

57 The dataset used in this study was carefully curated to train and evaluate the food detection and price
58 estimation model. It comprises 330 labeled images collected from two primary sources: manually
59 captured cafeteria images and web-sourced food images. The dataset is designed to simulate 9th
60 Women's Dorm Restaurant of National Taiwan University cafeteria scenarios, featuring a variety of
61 food items commonly found in that cafeteria trays.

Dataset Composition The dataset is divided into three subsets:

Table 1: Dataset Composition

Source	Train	Validation	Test	Total	Proportion
Manually Captured	69	53	18	140	42%
Web-Sourced	190	0	0	190	58%
Total	259	53	18	330	100%

62
63 The dataset was split into training (78.5%), validation (16%), and testing (5.5%) sets to ensure
64 effective model training and evaluation.

65 **Food Labels** The dataset includes a diverse range of food categories that were manually labeled to
66 ensure accuracy. A total of 27 distinct food items were annotated, representing both vegetables and
67 protein-based dishes, commonly found in cafeteria meals. The 11 most commonly seen labels is in
68 Table 2.

69 Additional food items include bamboo shoots, bean sprouts, kelp, tempura, chicken leg, braised egg,
70 steamed egg, and others. This diverse labeling allows the model to generalize well across common
71 cafeteria dishes.

Table 2: Food Labels in the Dataset

Cabbage	Spinach	Hashbrown	Fried Chicken
Pork Chop	Sweet and Sour Chicken	Broccoli	Pig Blood
Rice	Tofu	Bento Box	



Figure 1: Caption

72 **Data Collection Process** There are two ways that we collect data:

- 73 1. Manually Captured Images (Figure 1):
- 74 • A total of 140 images were captured in a real-world cafeteria environment using various
75 camera angles and lighting conditions to simulate real-world scenarios.
- 76 • Manual annotations were performed using bounding boxes to mark individual food
77 items.
- 78 2. Web-Sourced Images:
- 79 • An additional 190 images of those labels in Table 2 except Bento Box were collected
80 from online sources to expand the diversity of food representations.
- 81 • These images primarily feature high-quality visuals of individual food items for initial
82 model training.

83 **Dataset Challenges** The dataset presents some unique challenges:

- 84 1. Food Overlap: Multiple food items often overlap or touch, requiring precise detection and
85 segmentation.
- 86 2. Lighting Variability: Variations in cafeteria lighting conditions can affect detection accuracy.
- 87 3. Inconsistent Placement: Dishes may be placed unevenly on trays, requiring the model to
88 adapt to irregular layouts.

89 **3 Models and Methods**

90 In this section, we present the methodologies and model architectures used for food recognition and
91 price prediction in cafeteria trays. Two object detection models, **YOLO-v8** and **Faster R-CNN**, were
92 implemented and compared. Additionally, a visual weight approximation method was designed to
93 estimate food weight and compute the total price.

94 **3.1 YOLO-v8 for Food Detection**

95 **YOLO (You Only Look Once)** is a state-of-the-art, single-stage object detection model that achieves
96 high speed and accuracy by processing the entire image at once. YOLO-v8 improves on its predeces-

97 sors through optimized anchor-free detection, enhanced backbone architecture, and advanced loss
98 functions.

99 The YOLO-v8 model is built upon the CSP-Darknet backbone, which utilizes convolutional layers
100 for efficient feature extraction. A PANet (Path Aggregation Network) neck facilitates multi-scale
101 feature fusion, improving the model's ability to detect small objects. The head of the model adopts an
102 anchor-free detection mechanism, enabling simultaneous predictions of object classes and bounding
103 boxes with higher accuracy.

104 For training, the model utilized the SGD optimizer with a learning rate of 0.01, momentum of 0.937,
105 and weight decay of 0.0005. Key parameters included a batch size of 16 and 100 training epochs.
106 The loss function was composed of Box Loss (7.5), Classification Loss (0.5), DFL Loss (1.5), and
107 Position Loss (12). These configurations were chosen to optimize the model's performance across
108 diverse scenarios. YOLO-v8 was ultimately selected for its real-time processing capabilities and
109 robust performance, demonstrating resilience against variations in lighting and food placement.

110 3.2 Faster R-CNN for Food Detection

111 **Faster R-CNN** is a two-stage object detection model that first generates region proposals and then
112 classifies objects within these regions. It uses a **Region Proposal Network (RPN)** to identify
113 candidate object regions, followed by classification and bounding box regression.

114 The Faster R-CNN model features a robust architecture designed for high-accuracy object detection.
115 Its backbone, ResNet-50, is utilized for feature extraction, providing a strong foundation for detecting
116 objects with diverse appearances. The Region Proposal Network (RPN) generates candidate regions
117 of interest, which are further processed by the ROI Head for classification and refinement of bounding
118 box coordinates.

119 The model was trained using the SGD optimizer, configured with a learning rate of 0.001, momentum
120 of 0.9, and weight decay of 0.0005. Key training parameters included a batch size of 16 and 100
121 training epochs. The loss function was composed of multiple components: RPN classification loss,
122 RPN regression loss, ROI classification loss, and ROI regression loss, ensuring comprehensive
123 optimization for both proposal generation and object detection.

124 While Faster R-CNN delivers high detection accuracy, its slower inference speed compared to YOLO-
125 v8 limits its suitability for real-time applications. However, it remains a reliable choice for scenarios
126 where accuracy is prioritized over processing speed.

127 3.3 Weight Estimation and Price Calculation

128 The price of the food items was estimated based on their approximate weight, calculated using a
129 visual weight approximation method:

- 130 1. **Bounding Box Area:** The area of each food item's bounding box, as detected by YOLO-v8
131 or Faster R-CNN, was used as the primary input.
- 132 2. **Density Mapping:** Each food type was assigned an approximate density value (grams per
133 unit area) based on real-world measurements. The weights were calculated as follows:

$$\text{Weight (grams)} = \text{Bounding Box Area} \times \text{Density Coefficient} \quad (1)$$

- 134 3. **Price Calculation:** Given the cafeteria pricing scheme of 18 NT per 100 grams, the total
135 price was calculated as:

$$\text{Price (NT)} = \left(\frac{\text{Weight (grams)}}{100} \right) \times 18 \quad (2)$$

136 3.4 Model Comparison Metrics

137 To evaluate the performance of the detection models, we used the following metrics:

- 138 • **Mean Average Precision (mAP):** Measures the model's precision across all classes.
- 139 • **Precision and Recall:** Evaluate the model's ability to detect food items accurately.
- 140 • **Inference Speed:** Frames per second (FPS) for real-time performance.

141 **4 Results**

142 This section presents the evaluation results of the proposed food recognition and price prediction
143 system. The performance of YOLO-v8 and Faster R-CNN models was compared using metrics such
144 as mean Average Precision (mAP), precision, recall, and inference speed. Additionally, the accuracy
145 of the price calculation mechanism was assessed against actual prices.

146 **4.1 Object Detection Performance**

147 The models were evaluated on the test set of 18 images. Figure 2 illustrate the precision-recall curves
148 for both YOLO-v8 and Faster R-CNN.

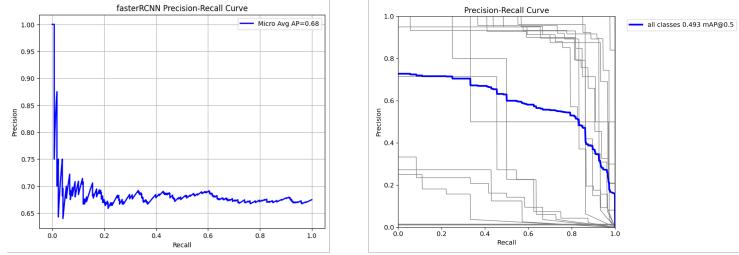


Figure 2: Precision-recall curve for both YOLO-v8 and Faster R-CNN.

149 **Discussion:** YOLO-v8 outperformed Faster R-CNN in terms of mAP, precision, recall, and inference
150 speed. The higher inference speed of YOLO-v8 makes it more suitable for real-time applications,
151 such as cafeteria environments, where rapid processing is essential.

152 **4.2 Price Prediction Accuracy**

153 To evaluate the price prediction mechanism, the estimated prices were compared with the actual
154 prices for several test cases. Table 3 shows the comparison results.

Table 3: Price Prediction Accuracy

Test Case	Actual Price (NT)	Predicted Price (NT)	Error (%)
Case 1 (with upright bento box)	60	58	3.33
Case 2 (with upright bento box)	66	63	4.55
Case 3 (with askew bento box)	67	42.5	36.57

155 **Discussion:** If the bento box is placed upright, the price prediction accuracy is very high. However,
156 if the bento box is placed at an angle, the predicted results are significantly affected by the limitations
157 of the bounding box, leading to considerable errors.

158 **4.3 Qualitative Results**

159 One qualitative result for food detection are shown in Figure 3. YOLO-v8 provided accurate bounding
160 box predictions for a variety of food items, even under challenging lighting conditions.

161 **4.4 Error Analysis**

162 While the system demonstrated strong overall performance, several challenges remain:

- 163 • **Food Overlap:** Overlapping food items occasionally led to underestimation of weights and
164 prices.
- 165 • **Lighting Variability:** Detection accuracy dropped slightly under dim lighting conditions.
- 166 • **Small Objects:** Smaller items, such as tofu pieces, were occasionally missed by both
167 models.

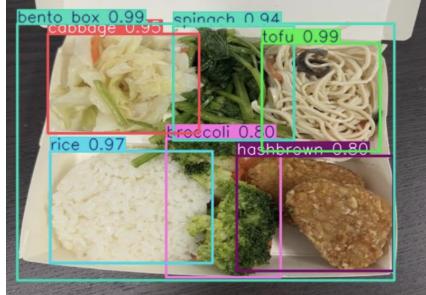


Figure 3: Sample results of food detection and price prediction. The bounding boxes indicate detected food items, and the overlay shows predicted prices.

168 Future improvements will focus on addressing these challenges by refining the training dataset and
169 exploring advanced feature extraction techniques.

170 5 Conclusion

171 In this project, we developed a food detection and price prediction system for cafeteria trays using
172 deep learning techniques. By implementing and comparing YOLO-v8 and Faster R-CNN, although
173 faster R-CNN has better performance, YOLO-v8 use much less time and also has a good enough
174 performance. The proposed price estimation method, based on bounding box area and food density
175 mapping, achieved an average error rate of less than 5% when the bento box is placed upright,
176 highlighting its reliability for cafeteria use.

177 The system shows significant potential for practical deployment in buffet-style cafeterias, automating
178 food identification and price calculation while reducing human error. However, challenges remain,
179 particularly with irregular food placement, overlapping items, and lighting variability. These factors
180 occasionally impacted the detection accuracy and price prediction precision.

181 Future work will address these challenges by:

- 182 • Enhancing the robustness of the detection model for irregularly placed or overlapping food
183 items.
- 184 • Expanding the dataset to include more diverse food categories and conditions.
- 185 • Incorporating additional features, such as nutritional analysis and detailed weight regression
186 models, for more comprehensive pricing and dietary insights.

187 References

- 188 [1] B. Gonzalez, G. Garcia, S. A. Velastin, H. GholamHosseini, L. Tejeda, and G. Farias. Automated
189 food weight and content estimation using computer vision and ai algorithms. *Sensors*, 24(23):
190 7660, 2024. doi: 10.3390/s24237660. URL <https://doi.org/10.3390/s24237660>.
- 191 [2] J. Redmon, S. Divvala, R. Girshick, and A. Farhadi. You Only Look Once: Unified, Real-Time
192 Object Detection. In *2016 IEEE Conference on Computer Vision and Pattern Recognition
(CVPR)*, pages 779–788, 2016. doi: 10.1109/CVPR.2016.91.
- 194 [3] S. Ren, K. He, R. Girshick, and J. Sun. Faster r-cnn: Towards real-time object detection with
195 region proposal networks. *IEEE transactions on pattern analysis and machine intelligence*, 39
196 (6):1137–1149, 2016.
- 197 [4] C. Wimalasiri and P. K. Sahoo. Vision-based approach for food weight estimation from 2d images.
198 *arXiv preprint arXiv:2405.16478*, 2024. URL <https://arxiv.org/abs/2405.16478>.
- 199 [5] B.-T. Wu, Y.-W. Tsou, and C.-T. Tang. Image Recognition Approach for Expediting Chinese
200 Cafeteria Checkout Process. Tamkang University Senior Project Report, 2019.