# Vision-Based Approach for Food Weight Estimation from 2D Images

Chathura Wimalasiri
*Dept. of Computer Engineering*
*University of Peradeniya*
Peradeniya, Sri Lanka
e18402@eng.pdn.ac.lk

Prasan Kumar Sahoo
*Dept. of Computer Science and Information Engineering*
*Chang Gung University*
Guishan, Taiwan
pksahoo@mail.cgu.edu.tw

*Abstract*—In response to the increasing demand for efficient and non-invasive methods to estimate food weight, this paper presents a vision-based approach utilizing 2D images. The study employs a dataset of 2380 images comprising fourteen different food types in various portions, orientations, and containers. The proposed methodology integrates deep learning and computer vision techniques, specifically employing Faster R-CNN for food detection and MobileNetV3 for weight estimation. The detection model achieved a mean average precision (mAP) of 83.41%, an average Intersection over Union (IoU) of 91.82%, and a classification accuracy of 100%. For weight estimation, the model demonstrated a root mean squared error (RMSE) of 6.3204, a mean absolute percentage error (MAPE) of 0.0640%, and an R-squared value of 98.65%. The study underscores the potential applications of this technology in healthcare for nutrition counseling, fitness and wellness for dietary intake assessment, and smart food storage solutions to reduce waste. The results indicate that the combination of Faster R-CNN and MobileNetV3 provides a robust framework for accurate food weight estimation from 2D images, showcasing the synergy of computer vision and deep learning in practical applications.

*Index Terms*—food weight estimation, food recognition, computer vision, deep learning

## I. Introduction

In today's digital world, food is more important than just a source of strength and nutrients. With the increasing use of technology, 2D food images are appearing everywhere and influencing our food choices. There are several uses of 2D food images that can be used to visualize and analyze data, food recognition and classification, monitor food quality and defect detection, assess food safety and hygiene practices, and create visually appealing and appetizing representations of food. With these facts, human life is getting easier and safer.

Artificial intelligence has emerged as a trans-formative force in the modern world, impacting virtually every aspect of our lives. Deep learning is a subset of artificial intelligence. Deep learning can perform complex tasks such as image recognition, natural language processing, and decision-making. A convolutional neural network (CNN) is a common deep learning architecture. CNNs recognize patterns and extract features from grid-like data, making them the go-to solution for various computer vision tasks. Computer vision is a subfield of AI that is used for classification, object detection, and segmentation. Combining both computer vision and deep learning concepts can be done with the most useful real-world applications, such as self-driving cars, medical image analysis, facial recognition, object detection and tracking, and image and video captioning.

Food weight estimation using artificial intelligence is an important step in real life. This concept can be applied to automatically estimate calorie intake based on pictures of meals by estimating food weight. Smart food storage containers can use this concept to reduce food waste, store them in proper conditions, maximize their shelf life, and reduce spoilage. Restaurants can use this to accurately measure the food weight using 2D images. This can be used in medical applications to monitor food and calorie amounts easily.

In this paper, we propose a method to estimate food weight from 2D images using deep learning and computer vision techniques. We use a novel dataset, which includes 2380 images with fourteen types of food in 2D with different portions, orientations, and containers. Food types and information are presented in Table I. After a food image is acquired, the first challenge is to recognize it. After that, using these results, weight estimation is performed.

The paper is structured as follows: In the next section, a brief overview of the related works in this area will be provided. Next, we describe our food recognition and weight estimation methods. Next, we discuss our results in food recognition and weight estimation. Finally, we summarize our work.

## II. Related Works

This paper is mainly focused on estimating food weight from 2D images using a vision-based approach. This aligns with various food-related tasks, including calorie estimation. Research by [1] proposed a method to find calories in Indonesian street food using 2D images. Before finding the calorie amount [1], find the weight of the food item using multiple linear regression. Then, [1] used a weight-calorie scale to find the amount of calories in food.

Research by [2] proposed a method to find the weight of lettuce using a 3D stereoscopic technique. [2] proposed a new image preprocessing technique to estimate weight using a 3D stereoscopic technique and estimated the weight of

fresh lettuce from a 3D spatial domain. Research by [3] proposed a method to localize the picking point of strawberries and estimate their weight. [3] proposed two novel datasets annotated with picking points, key points, weight, and size of strawberries. [3] used the concept that the weight of a strawberry is proportional to its shape, size, and density in their weight estimation. Research by [4] proposed a method to estimate the weight of meals in self-service lunch line restaurants using 2D images. [4] capture the top view of the meal, identify each food item, and estimate the weight of each item on the plate. Research by [5] proposed a method to estimate the weight of harumanis mango from an RGB image. All the mango image dataset background is an A4 sheet to reduce background noise and to function as a visual cue for CNNs. [5] defined a CNN architecture and the final layer is a dense layer with rectified linear unit (ReLU) activation function for regression task. Research by [6] proposed a method to estimate the weight of melons from unmanned aerial vehicle images. [6] detected all melons in an image and estimated their weight individually, finally estimating the weight of the melons yield. Research by [7] proposed a smart harvesting decision system to estimate fruits types, maturity level, and weight. [7] used a powerful supervised machine learning algorithm called support vector machine (SVM) to estimate weight. Weight isn't a discrete number; therefore, [7] used a regression technique called support vector regression (SVR). Research by [8] proposed a method to analyze food calories and nutrition. For that task, [8] estimated the weight of food using a linear regression model and, from that, estimated food calories and nutrition. Research by [9] proposed a method of dietary assessment for Chinese children. [9] estimated the weight of food using 2D images in the system in order to achieve their goal. They stored density information in their database and calculated the volume and estimated the weight using the relationship between weight, volume, and density. Research by [10] proposed a method to estimate food portions from a single view using geometric models. 3D reconstruction using a single view is a very challenging task; therefore, they have used geometrical models such as the shape of the container. [10] estimated the volume of each food with the help of the geometric models, then relationships between weight, volume, and density were used to estimate the weight of the food item. Research by [11] proposed a method to estimate food intake using a depth camera. First, the empty tray, before eating and after eating were captured, respectively. Using these three image information, the volume of food intake was estimated. After that, the weight was estimated using a specific gravity function. Finally, weight was converted to calories using food calorie database information. Research by [12] proposed a dietary assessment method to record daily food intake using a food image of a meal. First, the system identified and segmented food in the image. Then, they used a 3D reconstruction method for regular-shaped foods and an area-based weight estimation method for irregular-shaped foods.

In addition to application in food estimation, there are methods that have been proposed to estimate the weight of animals using 2D images. Research by [13] proposed a method to estimate fish weight without contact. [13] estimated fish weight using the perimeter of the fish. The fish contour information and 3D coordinates of the fish were taken to estimate the perimeter of the fish. After that, perimeter and fish weight information were used to develop a weight estimation model.

TABLE I
INFORMATION ON THE FOOD DATASET

| Food Type | Number of Images |
|---|---|
| Cherry Tomato | 80 |
| Oatmeal | 50 |
| Steamed Rice | 40 |
| Stir Fried Spinach | 122 |
| Sweet Corn | 252 |
| Grape | 80 |
| Guava | 240 |
| Orange | 86 |
| Papaya | 160 |
| Pineapple | 160 |
| Red Apple | 260 |
| Steamed Bun with Meat | 52 |
| Sweet Potato | 120 |
| Toast Bread | 478 |

## III. PROPOSED METHODOLOGY

It is challenging to divide the dataset for training, validation, and testing because the dataset is an imbalance; features like food weight, container, and orientation must be fairly divided. Therefore, the dataset divides 60%, 20%, and 20% for training, validation, and testing, respectively.
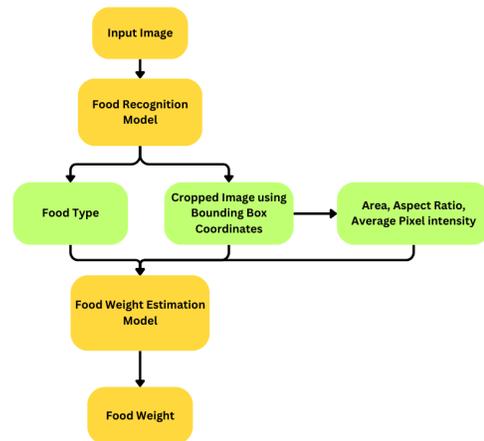


Fig. 1. The high-level architecture of the proposed solution

Fig. 1 shows a data flow diagram illustrating the high-level architecture of the proposed solution. The first task is to detect and recognize the food from the image. We use Faster R-CNN [14] which is a powerful two-stage CNN model for object detection. The two stages are the regional proposal network (RPN), which goes through the input image and proposes candidate regions where objects might be located, and the

other stage, which takes proposed regions from the RPN, refines the bounding box coordinates, and predicts the class probabilities for each region. For this task, ResNet [15] is used as backbone or feature extractor because Resnet is already pretrained on a large dataset, such as ImageNet and it allows them to capture general knowledge about visual patterns.

Before training, 'LabelImg' is used to annotate all the images. After annotating, we resize all the images to (224, 224), which is the standard size for ResNet [15] and was developed based on computer efficiency, model performance and compatibility with available pretrained models. We have fourteen food types, and the output node is set to fifteen because one node is for the background. The optimizer is used as Adam [16], the learning rate is 0.0001, the batch size is 1, and the number of epochs is 10.

After training Faster RCNN, which can be used to detect food items with bounding box coordinates and recognize them. All the images are detected, cropped using bounding box coordinates, and also recognized. These data are used to develop the weight estimation model. Fig. 2 provides details of the food weight distribution of training, validation, and testing data.
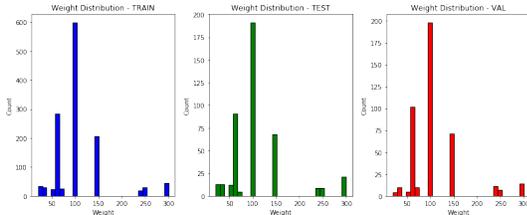


Fig. 2.  Food weight distribution of training, testing and validation data

Features of the weight estimation model are cropped image, food type, image area, aspect ratio, and average pixel intensity. Image area, aspect ratio, and average pixel intensity are calculated using cropped images. Image area means cropped image area. Equation (1) is used to calculate it.

$$Image\ Area = Image\ Height \times Image\ Width \quad (1)$$

Aspect ratio is the proportional relationship between width and height. Equation (2) is used to calculate it.

$$Aspect\ Ratio = \frac{Image\ Width}{Image\ Height} \quad (2)$$

Average pixel intensity means the average value of all individual pixel intensities. Equation (3) is used to calculate it.

$$Average\ pixel\ intensity = \left(\frac{1}{N}\right) \sum_{i=1}^{N} p_i \quad (3)$$

Where $N$ means total number of pixels and $p_i$ means intensity of pixel $i$.

Fig. 3 shows the proposed weight estimation model. It takes cropped images and goes through a backbone or feature

extractor. MobileNetV3 [19] is used as the backbone for this task. The output of the backbone is a single output value that concatenates food type, image area, aspect ratio, and average pixel intensity. These five values are used as the input vector of shape (5,1) to the dense layer 1, whose output is a vector of shape (64,1). The output vector from dense layer 1 goes through a Rectified Linear Unit (ReLU) [17] activation function, which introduces non-linearity into the network by replacing negative values with zero and positive values unchanged. Likewise, the input vector of dense layer 2 is (64,1) and the output vector is (32,1), and the output vector from dense layer 2 goes through a ReLU [17] activation function. Finally, the dense layer 3 input vector is (32,1) and the output vector is (1,1), which is the predicted food weight value.
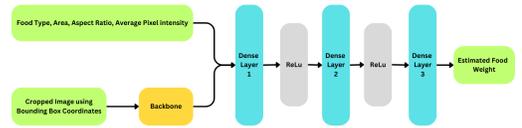


Fig. 3.  The proposed weight estimation model

We can derive an equation for our weight estimation model. The input X is a vector of size 5 to the dense layer 1.

$$X = [\alpha, FT, A, AR, API]$$

Where $\alpha$ is the output of the backbone, FT is food type, A is area, AR is aspect ratio, and API is average pixel intensity.

The output vector $\mathbf{Z}_i$ for dense layer $i$ is obtained by applying a linear transformation to the input vector $\mathbf{A}_{i-1}$.

$$Z_i = W_i A_{i-1} + b_i \quad (4)$$

Where $\mathbf{W}_i$ is a weight matrix and and $\mathbf{b}_i$ is a bias vector for layer $\mathbf{i}$.

The output vector $\mathbf{A}_i$ for ReLU [17] activation $\mathbf{i}$ is applied element-wise to the elements of $\mathbf{Z}_i$.

$$A_i = max(0, Z_i) \quad (5)$$

$\mathbf{A}_0$ is the input vector $\mathbf{i}$. Therefore, the output of the weight estimation model $\mathbf{Y}$ can be written using equations (4) and (5).

$$Y = W_3 \cdot \max(0, (W_2 \cdot \max(0, (W_1[\alpha, FT, A, AR, API] + b_1)) + b_2)) + b_3 \quad (6)$$

Where $\mathbf{W}_1$, $\mathbf{W}_2$, and $\mathbf{W}_3$ are weight matrices of sizes (64, 5), (32, 64) and (1, 32) respectively. $\mathbf{b}_1$, $\mathbf{b}_2$, and $\mathbf{b}_3$ are bias matrices of sizes (64,1), (32,1) and (1,1) respectively.

As preprocessing tasks, after finding the necessary values from the cropped image, all cropped images resize to (224, 224). Then the data augmentation technique, which is the random horizontal flip, is used for training data to increase the diversity of the training dataset. Next, training, validation, and testing data are normalized to the values that are presented

in Table II. The learning rate, optimizer, and loss function are 0.0001, Adam [16], and mean squared error, respectively. For training, batch size and number of epochs are used as 32 and 10, respectively.

| Data Type | Mean Value ($\mu$) | Standard Deviation Value ($\sigma$) |
|---|---|---|
| Training | 0.4890 | 0.2301 |
| Validation | 0.4844 | 0.2305 |
| Testing | 0.4917 | 0.2287 |

## IV. RESULTS AND DISCUSSIONS

The first task is to detect food items with a bounding box and recognize them. Faster R-CNN [14] and RetinaNet [18] object detection algorithms were used and selected one according to their performance. The performance of models was evaluated using mAP (mean average precision), including mAP@0.5 and mAP@0.75. mAP, a higher value indicates better performance in object detection. Additionally, we used two metrics: classification accuracy (8) and average Intersection over Union (9). The results are presented in Table III.

$$mAP = \frac{1}{N} \sum_{i=1}^{N} AP_i \qquad (7)$$

Where $AP_i$ is the Average Precision for class $i$ and $N$ is the number of classes.

$$Classification\ Accuracy = \frac{Num.\ of\ Correct\ Predictions}{Total\ Number\ of\ Predictions} \times 100 \quad (8)$$

$$Average\ IoU = \frac{1}{N} \sum_{i=1}^{N} IoU_i \qquad (9)$$

Where $IoU_i$ is the IoU for the $i$th prediction and $N$ is the total number of predictions.

| Method | Dataset | mAP | mAP$_{0.5}$ | mAP$_{0.75}$ | Classification Accuracy | Average IoU |
|---|---|---|---|---|---|---|
| Faster R-CNN | Train | 0.8522 | 0.9999 | 0.9999 | 1.0000 | 0.9154 |
| | Val | 0.8423 | 0.9997 | 0.9667 | 1.0000 | 0.9095 |
| | Test | 0.8341 | 1.0000 | 1.0000 | 1.0000 | 0.9182 |
| RetinaNet | Train | 0.7053 | 0.7909 | 0.7878 | 0.8716 | 0.9314 |
| | Val | 0.7042 | 0.7942 | 0.7912 | 0.8693 | 0.9266 |
| | Test | 0.7028 | 0.7986 | 0.7892 | 0.8716 | 0.9256 |

Faster R-CNN [14] has a higher mAP in all dataset categories than RetinaNet [18]. We also considered the mAP at 0.5 and 0.75 Intersection over Union (IoU) thresholds. Faster R-CNN [14] has higher values than RetinaNet [18] in every combination. Food type prediction with detection is also very important in this problem because food type is an input to the weight estimation model. Faster R-CNN [14] classification accuracy is 100% in all datasets, but RetinaNet

[18] has lower classification accuracy in all datasets. Both algorithms have good value for average IoU. That means the predicted bounding box with high probability score is closely aligned with the actual bounding box. RetinaNet [18] mAP values are low but have high average IoU values. The reasons might be: the high average IoU might be misleading because it only considers the predicted bounding box with the highest probability score. However, RetinaNet [18] might still generate other incorrect bounding boxes and misclassify the food name, generating bounding boxes with low-confidence results and also generating a high False Positive rate and a high False Negative rate. Fig. 4 shows the low performance of RetinaNet [18] with examples of low confidence levels, wrong food name predictions, and wrong bounding box predictions of RetinaNet [18].
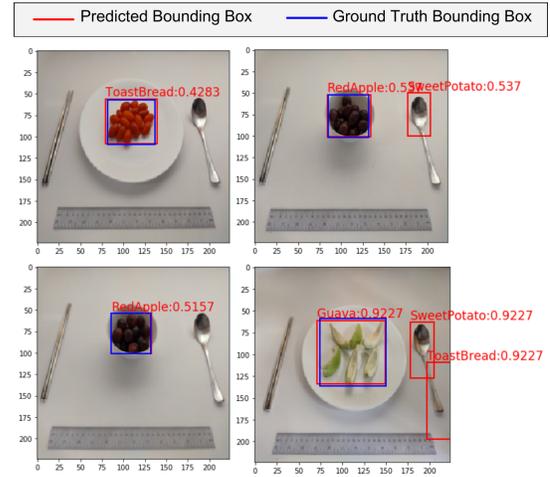


Fig. 4. Wrong prediction results of RetinaNet

Therefore, Faster R-CNN [14] is used for our food detection and recognition. Fig. 5 shows sample images showing the detection and recognition results.

Next task is to estimate the weight of the food. We have experimented with task performance with different backbones. We used Mean Squared Error (MSE) (10), Root Mean Squared Error (RMSE) (11), Mean Absolute Error (MAE) (12), Mean Absolute Percentage Error (MAPE) (13) and Coefficient of Determination (R-Squared) (14) to evaluate the performance of weight estimation model.

$$MSE = \frac{1}{n} \sum_{i=1}^{n} (y_i - \hat{y}_i)^2 \qquad (10)$$

$$RMSE = \sqrt{\frac{1}{n} \sum_{i=1}^{n} (y_i - \hat{y}_i)^2} \qquad (11)$$

$$MAE = \frac{1}{n} \sum_{i=1}^{n} |y_i - \hat{y}_i| \qquad (12)$$

$$MAPE = \frac{1}{n} \sum_{i=1}^{n} \left| \frac{y_i - \hat{y}_i}{y_i} \right| \times 100 \qquad (13)$$
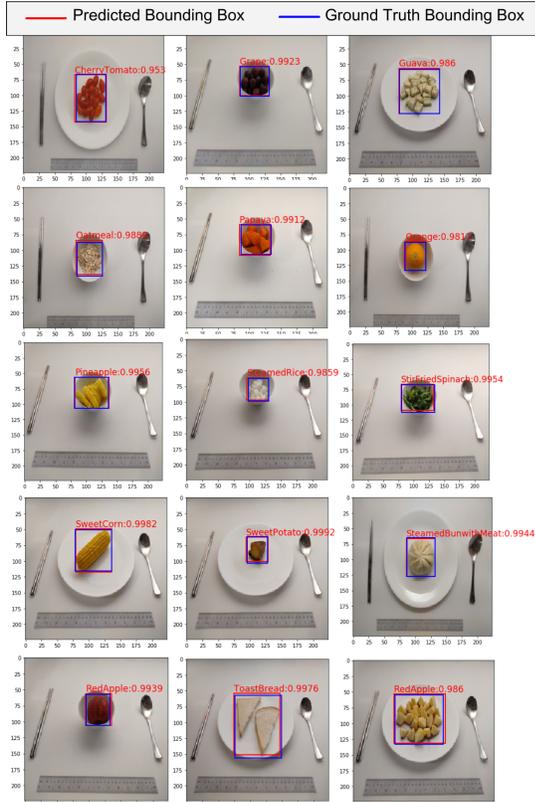
Fig. 5. Faster R-CNN Food Detection and Recognition Results

Where $\mathbf{N}$ is the number of samples, $\hat{\mathbf{y}}_i$ is the actual value, and $\mathbf{y}_i$ is the predicted value.

$$R^2 = 1 - \frac{\sum_{i=1}^{n}(y_i - \hat{y}_i)^2}{\sum_{i=1}^{n}(y_i - \bar{y})^2} \qquad (14)$$

Where $\bar{\mathbf{y}}$ is the mean of the actual values.

MSE (Mean Squared Error) and RMSE (Root Mean Squared Error) are suitable when you want to penalize larger errors more. MAE (Mean Absolute Error) is suitable when we prefer a metric that is less sensitive to outliers. MAPE (Mean Absolute Percentage Error) is valuable when interpretability is crucial, and you want errors represented as percentages. R-Squared (Coefficient of Determination) helps understand the proportion of variance explained by the model. Using multiple metrics to evaluate your model is a good practice as it provides a more comprehensive understanding of its performance. Each metric may capture different aspects of model performance, and using a variety of them can help you gain insights into how well your model is addressing various aspects of the problem.

We tested five different models, which are MobileNetV3 [19], AlexNet [20], ResNet50 [15], ResNet101 [15], and DenseNet121 [21]. Testing results are the most important result because they represent unseen data for the model. MobileNetV3 [19] has better values in MSE, RMSE, MAE, and R-Squared. AlexNet [20] has a better MAPE value. The MAPE value between MobileNetV3 [19] and AlexNet [20] has not a large difference. DenseNet121 [21] has larger errors

when compared with other models. Therefore, MobileNetV3 [19] was used as the backbone of the weight estimation model. Table V presents the average confidence scores, average actual weight, average predicted weight, average weight error, and average absolute weight error of each class, and the final row represents overall. Fig. 6 shows the overall system summary. First, it detects the food item with a bounding box, and using a bounding box weight estimation model, it predicts the weight and shows the actual and predicted weight in the figure.

TABLE IV
RESULT OF THE FOOD WEIGHT ESTIMATION MODEL

| Backbone | Dataset | MSE | RMSE | MAE | MAPE | R-Squared |
|---|---|---|---|---|---|---|
| MobileNetV3 | Train | 23.5711 | 4.8550 | 3.5700 | 0.0502% | 0.9933 |
| | Val | 51.7263 | 7.1921 | 4.8986 | 0.0605% | 0.9848 |
| | Test | 39.9473 | 6.3204 | 4.8219 | 0.0640% | 0.9865 |

TABLE V
AVERAGE CONFIDENCE, ACTUAL AND PREDICTED WEIGHTS, WEIGHT ERROR, AND ABSOLUTE WEIGHT ERROR FOR EACH FOOD CLASS AND TOTAL IMAGES

| class | Average Confidence Scores | Average Actual Weight | Average Predicted Weight | Average Weight Error | Average Absolute Weight Error |
|---|---|---|---|---|---|
| Cherry Tomato | 0.8738 | 100 | 99.5539 | 0.4461 | 0.4461 |
| Grape | 0.9927 | 100 | 100.6670 | -0.6670 | 0.6670 |
| Guava | 0.9920 | 166.6667 | 165.7642 | 0.9024 | 0.9024 |
| Oatmeal | 0.9814 | 20 | 22.25811 | -2.2581 | 2.2581 |
| Orange | 0.9842 | 150 | 151.1659 | -1.1659 | 1.1659 |
| Papaya | 0.9864 | 100 | 99.6441 | 0.3559 | 0.3559 |
| Pineapple | 0.9828 | 100 | 106.3429 | -6.3429 | 6.3429 |
| Red Apple | 0.9675 | 150 | 153.3526 | -3.3526 | 3.3526 |
| Steamed Bun with Meat | 0.9944 | 30 | 31.5970 | -1.5970 | 1.5970 |
| Steamed Rice | 0.9585 | 50 | 53.1813 | -3.1813 | 3.1813 |
| Stir Fried Spinach | 0.9779 | 100 | 102.1309 | -2.1309 | 2.1309 |
| Sweet Corn | 0.9915 | 126.19 | 127.7443 | -1.5543 | 1.5543 |
| Sweet Potato | 0.9968 | 146.6667 | 145.0690 | 1.5977 | 1.5977 |
| Toast Bread | 0.9963 | 60 | 66.5558 | -6.5558 | 6.5558 |
| TOTAL | 0.9769 | 99.9660 | 101.7876 | -1.8217 | 1.8217 |

## V. CONCLUSION

This study successfully demonstrates a novel approach for estimating food weight from 2D images using deep learning and computer vision techniques. By utilizing Faster R-CNN for food detection and MobileNetV3 as the backbone for weight estimation, the methodology achieved high accuracy in both detection and weight prediction tasks. The results highlighted significant metrics, including a mean average precision of 83.41% and an R-squared value of 98.65%, emphasizing the robustness and reliability of the proposed system. This approach can revolutionize applications in various domains such as healthcare, where accurate dietary assessments are crucial, and in the food industry, where efficient weight estimation can optimize processes and reduce waste. The integration of advanced AI techniques in everyday applications marks a significant step towards leveraging technology for better health and efficiency, paving the way for further innovations in the field.
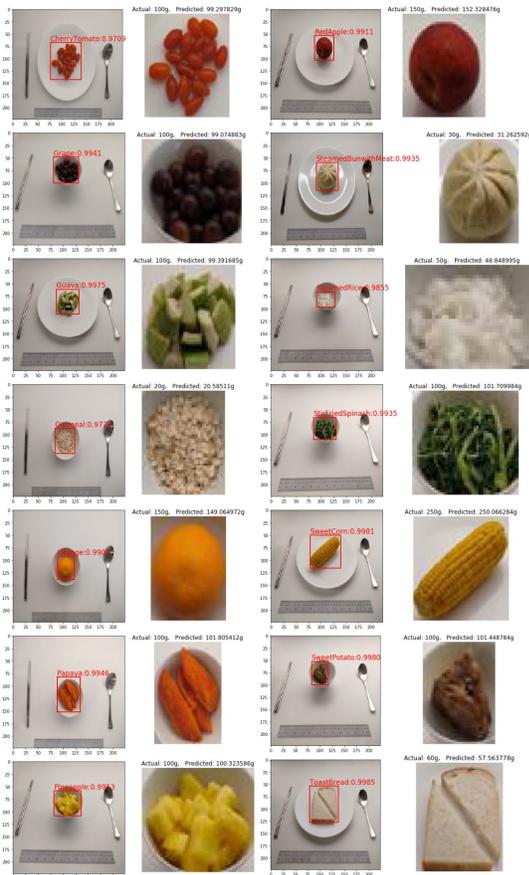
Fig. 6. Food Detection and Weight Estimation Results

# REFERENCES

[1] N. Aditama and R. Munir, "Indonesian Street Food Calorie Estimation Using Mask R-CNN and Multiple Linear Regression," 2022 Second International Conference on Power, Control and Computing Technologies (ICPC2T), Raipur, India, 2022, pp. 1-6, doi: 10.1109/ICPC2T53885.2022.9776804.

[2] N. Chimwai, S. Saiyod and R. Varakulsiripunth, "Fresh Weight Estimation of Lettuce Using 3D Stereoscopic Technique," 2023 8th International Conference on Business and Industrial Research (ICBIR), Bangkok, Thailand, 2023, pp. 108-113, doi: 10.1109/ICBIR57571.2023.10147436.

[3] A. Tafuro, A. Adewumi, S. Parsa, G. E. Amir and B. Debnath, "Strawberry picking point localization ripeness and weight estimation," 2022 International Conference on Robotics and Automation (ICRA), Philadelphia, PA, USA, 2022, pp. 2295-2302, doi: 10.1109/ICRA46639.2022.9812303.

[4] T. Sarapisto, L. Koivunen, T. Mäkilä, A. Klami and P. Ojansivu, "Camera-Based Meal Type and Weight Estimation in Self-Service Lunch Line Restaurants," 2022 12th International Conference on Pattern Recognition Systems (ICPRS), Saint-Etienne, France, 2022, pp. 1-7, doi: 10.1109/ICPRS54038.2022.9854056.

[5] M. H. Bin Ismail, M. N. Wagimin and T. R. Razak, "Estimating Mango Mass from RGB Image with Convolutional Neural Network," 2022 3rd International Conference on Artificial Intelligence and Data Sciences (AiDAS), IPOH, Malaysia, 2022, pp. 105-110, doi: 10.1109/AiDAS56890.2022.9918807.

[6] Aharon Kalantar, Yael Edan, Amit Gur, Iftach Klapp, A deep learning system for single and overall weight estimation of melons using unmanned aerial vehicle images, Computers and Electronics in Agriculture, https://doi.org/10.1016/j.compag.2020.105748.

[7] M. Faisal, F. Albogamy, H. Elgibreen, M. Algabri and F. A. Alqershi, "Deep Learning and Computer Vision for Estimating Date Fruits Type, Maturity Level, and Weight," in IEEE Access, vol. 8, pp. 206770-206782, 2020, doi: 10.1109/ACCESS.2020.3037948.

[8] M. -L. Chiang, C. -A. Wu, J. -K. Feng, C. -Y. Fang and S. -W. Chen, "Food Calorie and Nutrition Analysis System based on Mask R-CNN," 2019 IEEE 5th International Conference on Computer and Communications (ICCC), Chengdu, China, 2019, pp. 1721-1728, doi: 10.1109/ICCC47050.2019.9064257.

[9] P. Xu, D. Chen, X. Liu and J. Loo, "Image-based Dietary Assessment System for Chinese Children," 2018 IEEE International Conference on Big Data (Big Data), Seattle, WA, USA, 2018, pp. 5471-5473, doi: 10.1109/BigData.2018.8622579.

[10] S. Fang, C. Liu, F. Zhu, E. J. Delp and C. J. Boushey, "Single-View Food Portion Estimation Based on Geometric Models," 2015 IEEE International Symposium on Multimedia (ISM), Miami, FL, USA, 2015, pp. 385-390, doi: 10.1109/ISM.2015.67.

[11] H. -C. Liao, Z. -Y. Lim and H. -W. Lin, "Food intake estimation method using short-range depth camera," 2016 IEEE International Conference on Signal and Image Processing (ICSIP), Beijing, China, 2016, pp. 198-204, doi: 10.1109/SIPROCESS.2016.7888252.

[12] Y. He, C. Xu, N. Khanna, C. J. Boushey and E. J. Delp, "Food image analysis: Segmentation, identification and weight estimation," 2013 IEEE International Conference on Multimedia and Expo (ICME), San Jose, CA, USA, 2013, pp. 1-6, doi: 10.1109/ICME.2013.6607548.

[13] Xiaoning Yu, Yaqian Wang, Jincun Liu, Jia Wang, Dong An, Yaoguang Wei, Non-contact weight estimation system for fish based on instance segmentation, Expert Systems with Applications, Volume 210, 2022, 118403, ISSN 0957-4174, https://doi.org/10.1016/j.eswa.2022.118403.

[14] S. Ren, K. He, R. Girshick and J. Sun, "Faster R-CNN: Towards Real-Time Object Detection with Region Proposal Networks," in IEEE Transactions on Pattern Analysis and Machine Intelligence, vol. 39, no. 6, pp. 1137-1149, 1 June 2017, doi: 10.1109/TPAMI.2016.2577031.

[15] K. He, Kaiming, Xiangyu Zhang, Shaoqing Ren, and Jian Sun. "Deep residual learning for image recognition." In Proceedings of the IEEE conference on computer vision and pattern recognition, pp. 770-778. 2015.

[16] Kingma, Diederik P., and Jimmy Ba. "Adam: A method for stochastic optimization." arXiv preprint arXiv:1412.6980 (2014).*

[17] Agarap, A. F. (2018). Deep learning using rectified linear units (relu). ArXiv Preprint arXiv:1803.08375.

[18] T.-Y. Lin, P. Goyal, R. Girshick, K. He, and P. Dollar, "Focal loss for dense object detection," arXiv preprint arXiv:1708.02002, 2017.

[19] Howard, Andrew & Pang, Ruoming & Adam, Hartwig & Le, Quoc & Sandler, Mark & Chen, Bo & Wang, Weijun & Chen, Liang-Chieh & Tan, Mingxing & Chu, Grace & Vasudevan, Vijay & Zhu, Yukun. (2019). Searching for MobileNetV3. 1314-1324. 10.1109/ICCV.2019.00140.

[20] Krizhevsky, Alex, Ilya Sutskever, and Geoffrey E. Hinton. "ImageNet Classification with Deep Convolutional Neural Networks." Advances in Neural Information Processing Systems 25 (2012): 1097-1105.

[21] Huang, Gao, Zhuang Liu, Laurens van der Maaten, and Kilian Q. Weinberger. "Densely connected convolutional networks." In Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition, pp. 4700-4708. 2017.