

Deep Learning based Food Instance Segmentation using Synthetic Data

Deokhwan Park¹, Joosoon Lee¹, Junseok Lee¹ and Kyoobin Lee¹

Abstract— In the process of intelligently segmenting foods in images using deep neural networks for diet management, data collection and labeling for network training are very important but labor-intensive tasks. In order to solve the difficulties of data collection and annotations, this paper proposes a food segmentation method applicable to real-world through synthetic data. To perform food segmentation on healthcare robot systems, such as meal assistance robot arm, we generate synthetic data using the open-source 3D graphics software Blender placing multiple objects on meal plate and train Mask R-CNN for instance segmentation. Also, we build a data collection system and verify our segmentation model on real-world food data. As a result, on our real-world dataset, the model trained only synthetic data is available to segment food instances that are not trained with 52.2% mask AP@all, and improve performance by +6.4%p after fine-tuning comparing to the model trained from scratch. In addition, we also confirm the possibility and performance improvement on the public dataset for fair analysis. Our code and pre-trained weights are available online at: <https://github.com/gist-ailab/Food-Instance-Segmentation>

I. INTRODUCTION

Some experts predict that 38 percent of adults in the world will be overweight and 20 percent obese by 2030 if the trend continues [1]. Due to the increasing obesity rate, the importance of diet management and balanced nutrition intake has recently increased. In particular, services are gradually being developed to automatically calculate and record kinds of food and calories through photos of food to be consumed. The most important technology in this service is food recognition and can be widely used in a variety of service robots, including meal assistance robots, serving robots, and cooking robots.

Because of increasing importance of food-aware tasks, many researchers are working hard on the production of food-aware datasets and the development of food recognition. There are three methods to recognize food: food classification, food detection, and food segmentation. Food classification is a task that matches the type of food in an image through a single image, and many public datasets are also released because it is relatively easier than other tasks during the data collection and labeling phase. However, in order to determine a more accurate food intake in the diet management service, it is necessary to pinpoint the real food

portion within the image. Therefore, food segmentation are more useful in this service than food classification and food detection which provides information of the food types and the position by expressing in a bounding box. Nevertheless, there are three difficulties in food segmentation. First, as shown in Table I, released public datasets available for food segmentation are very scarce compared to food classification public datasets, and most datasets are not publicly available even if released. Second, when producing a dataset personally, there are a tremendous variety of food types and it takes a huge labor cost to labeling. Third, food segmentation is still a challenging task because the variations in shape, volume, texture, color, and composition of food are too large.

To address the presented difficulties, we employed two methods. First, We introduced deep neural network for instance food segmentation. In the early works of food segmentation, multiple food items were recognized mainly through image processing techniques: Normalized Cuts [2], Deformable Part Model [3], RANSAC [4], JSEG [3], [5], Grab Cut [6], and Random Forest [7]. In those cases, the sophistication of technique is more important than the acquisition of datasets. Lately, with the introduction of deep learning, deep neural network has eliminated the hassle of image processing by finding food features in the image on its own. There is a study that simultaneously localizes and recognizes foods in images using Fast R-CNN [8]. Moreover, there is CNN-based food segmentation using pixel-wise annotation-free data through saliency map estimation [9]. However, most relevant studies do not distinguish the same type of food in different locations as semantic segmentation, and it is most important to provide sufficient data to allow itself to learn. In that sense, secondly, we generated synthetic data and train these to apply food segmentation in real-world environments, called Sim-to-Real technique. The Sim-to-Real is an efficient technique that is already being studied in robot-related tasks, such as robot simulation [10] and robot control [11], etc. Also, it has the advantage of overcoming environmental simulations or data that are difficult to implement in real-world environments. Using this application, segmentation masks were easily obtained by randomly placing plates and multiple objects in a virtual environment to create a synthetic data describing the situation in which food was contained on the plate in a real world. Using this synthetic data and the Mask R-CNN [12] model, which is most commonly used in segmentation tasks. We conduct a class-agnostic food instance segmentation that recognizes that food types are not classified (only classify background and food) but are different. Furthermore, we found the following effects:

¹School of Integrated Technology (SIT), Gwangju Institute of Science and Technology (GIST), Republic of Korea. joosoon1111@gist.ac.kr

© 2021 IEEE. Personal use of this material is permitted. Permission from IEEE must be obtained for all other uses, in any current or future media, including reprinting/republishing this material for advertising or promotional purposes, creating new collective works, for resale or redistribution to servers or lists, or reuse of any copyrighted component of this work in other works.

TABLE I
LIST OF FOOD DATASETS

Name	Task	Reference
Food50	Classification	[14]
PFID	Classification	[15]
TADA	Classification	[16]
Food85	Classification	[17]
Food50Chen	Classification	[18]
UEC FOOD-100	Detection	[19]
Food-101	Classification	[20]
UEC FOOD-256	Detection	[21]
UNICT-FD1200	Classification	[22]
VIREO	Classification	[23]
Food524DB	Classification	[24]
Food475DB	Classification	[25]
MAFood-121	Classification	[26]
ISIA Food-200	Classification	[27]
FoodX-251	Classification	[28]
ChineseFoodNet	Classification	[29]
UNIMIB2015	Classification and Leftover	[30]
Food201-Segmented	Segmentation	[31]
UNIMIB2016	Classification and Segmentation	[13]
SUECFood	Segmentation	[32]
Food50Seg	Segmentation	[33]

- Unseen food instance segmentation of first-time tableware and first-time food is possible in real-world environments through random object creation of synthetic data and deep learning
- Food segmentation is sufficiently possible in real-world environments when learning using synthetic data only
- After learning with synthetic data, fine-tuning with real-world data improves performance
- By distinguishing the same food in different locations, it can be used efficiently in robot fields, such as food picking, which can be utilized later

This paper is divided into 4 sections including this introduction section. In section 2, the data production process of synthetic and real-world data, models and parameters used in learning, and evaluation metrics used in performance comparisons are described. In section 3, performance comparison results were described according to the combination of learning data: synthetic data, our real-world data we collected, and public data called UNIMIB2016 [13]. Finally, in section 4, the conclusions are given.

II. METHODS

We propose a unseen food segmentation method that enables segmentation of untrained foods from real-world images. We used a deep neural network and constructed a synthetic dataset and a real-world dataset for unseen food segmentation using deep learning. For training deep neural network, data is the most important factor. In reality, however, it is quite challenging to build an appropriate dataset for every task. Therefore, we used Sim-to-Real, which learns deep neural networks using synthetic data and applies them

to real-world. If we get real food data, a lot of time and expense is needed for data collection and annotation. So, we generated synthetic data using Blender simulator, a computer graphics software, to conserve resource. Also, we collected real-world data by building our food image acquisition system for verification of unseen real-world food segmentation.

A. Dataset

Synthetic Dataset The use of synthetic data for training and testing deep neural networks has gained in popularity in recent years [34], [35]. The Blender, a 3D computer graphics production software capable of realistic rendering [36], is often used to create synthetic data. Realistic high-quality synthetic data is required for deep neural networks to show high performance for real situations. Therefore, we generated the synthetic data using Blender.

In general, food is usually served in bowls and plates. Especially, meal tray is usually used in hospital and school. We actively introduce domain randomization to ensure that the distribution of synthetic data includes real-world data. So we created a random texture on the meal tray to recognize a variety of plates robustly, and created a variety of background textures and distractors around meal tray to be robust against environmental changes. In addition lighting conditions, such as the number, position, and intensity of light points in the virtual environment, also changed during data generation phase. To express food in synthetic data, various kinds of primitives were grouped together and placed on a plate, that resemble food with various colors and textures on the meal tray so that the network can recognize various foods robustly. Therefore, we placed meal tray modeled using the blender and generated objects of various sizes and shapes on the meal tray in a virtual simulation space as shown in Fig 1. We then generated synthetic data by capturing it at various angles and locations of camera. We created 28,839 synthetic data, including RGB images and mask images, for unseen food segmentation. As shown in Fig 2 as the examples of dataset, textures on the plate and background are in a complex form of mixed colors and patterns. Food-like objects located in the food tray and distractors outsider of meal tray composed of clustered primitives also have diverse colors. However, in the mask images, only the objects expressing food are projected as instance for segmentation, while the distractors are expressed as background.

Real-world Dataset We built a real-world food dataset for 50 kinds of Korean food. We selected 50 kinds of Korean food (rice, kimchi, fried egg, etc) through consultation and investigation by experts of hospital and institution, and collected dataset. Each meal tray was assembled with five food items shown in Fig 3. We built a real-world food dataset using the food acquisition system that captures images from various angles, heights, and lights as shown in fig 3. We generated data from various backgrounds to verify the robustness of the network even in environmental changes. We make a real-world food dataset by annotating 229 images acquired

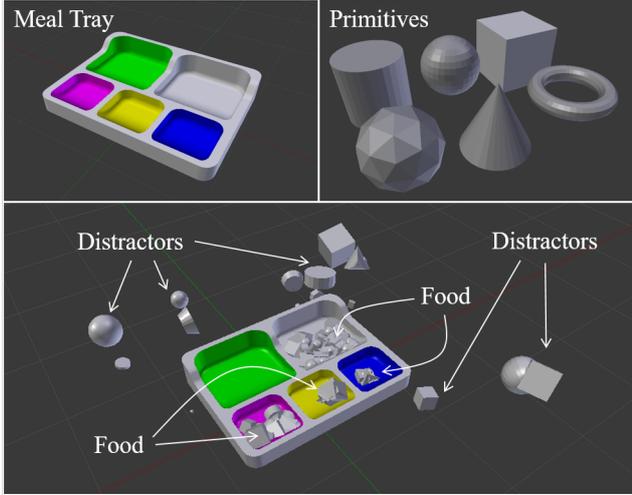


Fig. 1. Examples of synthetic dataset

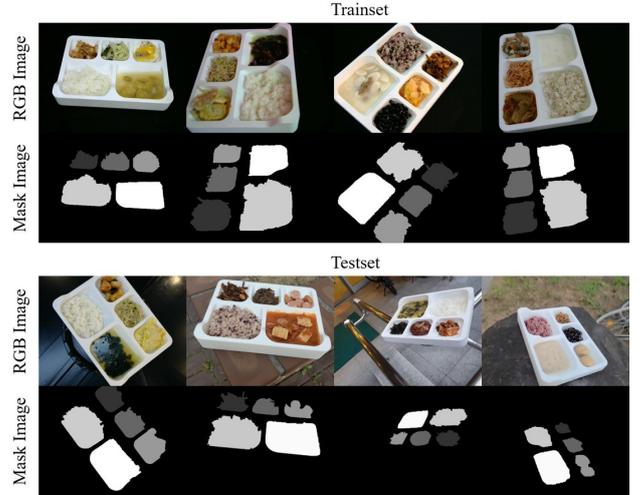


Fig. 4. Examples of real-world dataset.

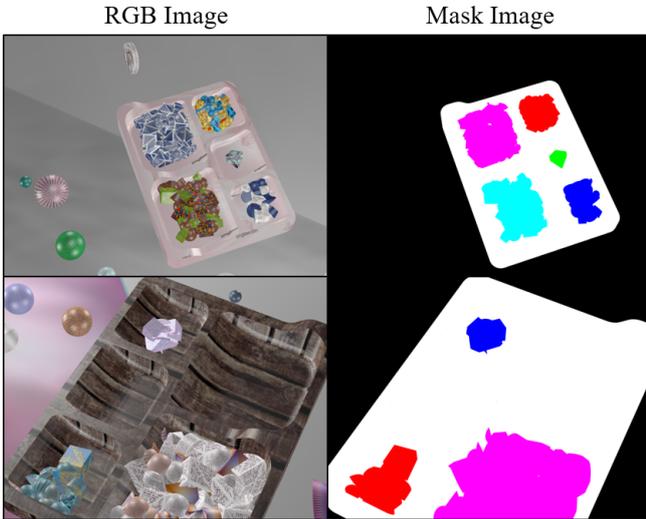


Fig. 2. Examples of synthetic dataset

through the food image acquisition system. The examples of dataset is shown in Fig 4.



Fig. 3. (Left) Examples of Korean food (Right) Data acquisition system

B. Deep Neural Network

We used Mask R-CNN [12] that widely used in the instance segmentation. Mask R-CNN [12] is an extension of the Fast RCNN [37], an algorithm used for object detection. The overall network architecture are shown in Fig . As shown

in the Figure 5, Mask R-CNN [12] consists of Backbone network, region proposal network(RPN), feature pyramid network(FPN), RoIAlign, and classifier. Mask-RCNN [12] is built on a backbone convolutional neural network architecture for feature extraction. Backbone network used a feature pyramid network based on a ResNet-50. In feature pyramid network, the features of various layers are considered together in a pyramid-shaped manner, it gives rich semantic information compared to single networks that use only the last feature. Region proposal network is a network that scans images by sliding window and finds areas containing objects. We refer to the area that RPN searches as anchors and use RPN predictions to select higher anchors that are likely to contain objects and refine their location and size. On the last stage, Region proposal network uses the proposed ROI to perform class preference, bounding-box regression, and mask preference. We give data and ground truth of food image as input to the network and we get output the instances of segmentation.

C. Training Details

We trained MASK R-CNN [12] model implemented in Py-Torch [38] with stochastic gradient descent(SGD) optimizer configured with learning rate of 0.0001, weight decay of 0.00005 and batch size of 8 on Titan RTX (24GB) GPU. We trained model on three types, first training with only synthetic dataset, second training only real-world dataset, the last fine-tuning with real-world dataset after pre-training on synthetic dataset. During fine-tuning the model, the model trained with only synthetic data first, and then only real dataset is used to fine-tune the pre-trained model.

D. Evaluation Metrics

For performance evaluation for unseen food segmentation, we utilize the same metric of COCO dataset [39], one of the most popular criteria of instance segmentation. The Intersection over Union (IoU), also known as the Jacquard Index, is a simple and highly effective rating metric that

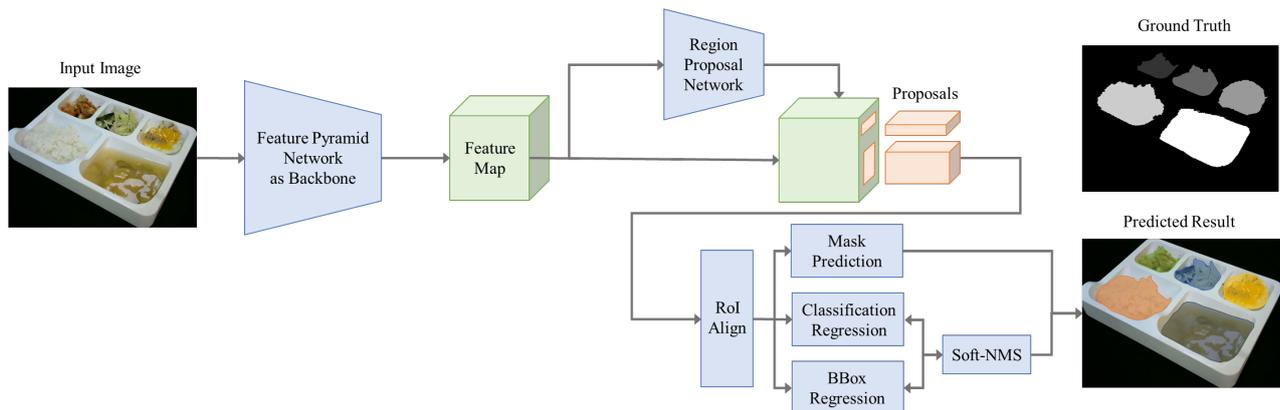


Fig. 5. The architecture of Mask R-CNN using for food instance segmentation.

calculates the overlapping area between the predicted and ground truth divisions: $IoU = \text{area of overlap} / \text{area of union}$. The proposed outputs of segmentation are post-processed with non-max suppression by the threshold of 0.5 for IoU.

The mean Average Precision (mAP) is used as an evaluation metric of the performance of the instance segmentation. Precision and recall are required to calculate the mAP. Precision means the true positive ratio of predicted results which can be calculated by adding true positive and false positive: $Precision = \text{true positive} / (\text{true positive} + \text{false positive})$. Recall means the true positive ratio of all ground truths which can be calculated by adding true positive and false negative: $Recall = \text{true positive} / (\text{true positive} + \text{false negative})$. Therefore, a high Recall value means that deep neural network recorded a high proportion of the predicted results among ground truths.

This results in mean Average Precision (mAP) being obtained through the Recall and Precision values. The main metric for evaluation is mean Average Precision (mAP), which is calculated by averaging the precisions under Intersection over Union (IoU) thresholds from 0.50 to 0.95 at the step of 0.05.

III. EXPERIMENT AND RESULTS

We experimented that training the MASK R-CNN [12] model on synthetic dataset and evaluation on our real-world dataset to verify the performance of unseen food segmentation. Furthermore, we conducted an experiment using a public dataset to verify the generalized performance of the algorithm. In all the experiments, our trained model segments food instances, which are category-agnostic and only certain to be food as a single category.

A. Result on our dataset

We categorized our real-world dataset into three types: easy, medium, and hard, based on background diversity within the image. Easy samples have a completely black background, medium samples have a black background with light reflection and hard samples have a wide variety of backgrounds. We have 73 easy samples, 61 medium samples, and 95 hard samples. Easy samples were used for training and medium and hard samples were used for testing.

TABLE II
SEGMENTATION EVALUATION RESULTS OF MASK AP AND BOX AP FOR EACH DATASET.

Test Sets	Metric	Synthetic+Real ¹	Synthetic Only	Real Only
Our test set	BBOX ²	-	80.0	-
	SEG ²	-	87.9	-
Our test set	BBOX	76.1	51.4	65.6
	SEG	79.0	52.2	72.6
UNIMIB	BBOX	80.6	35.7	79.3
	SEG	82.7	32.9	81.7

¹Synthetic+Real means pre-training with synthetic data and then fine-tuning with real-world data.

²BBOX means box AP@all and SEG means mask AP@all as defined in COCO dataset [39].

The experimental results can be found in Table II. The two columns of Table II (headed as *Synthetic Only* and *Real Only*) demonstrate the performance of models that trained only synthetic data and real data from-scratch, respectively. The column of *Syn+Real* shows the performance of the model fine-tuned on real-world data after pre-training on synthetic data. The real-world data utilized on each training phase, are our dataset and public dataset UNIMIB2016 [13], headed on each rows. Sim-to-Real can show good performance by training network using similar synthetic data to real-world reported on the column of *Synthetic Only*. Our network trains with only synthetic data and shows 52.2% in terms of mAP as a result of evaluating with real-world data. The result suggested that the network learned by using only synthetic data via Sim-to-Real to become unseen food segmentation for real-world data. Furthermore, we confirm that the performance increased by about 8.8% when the model was fine-tuned with real data compared to learning with real-world data from scratch. As shown in Figure 6, the model trained with synthetic data only tends not to recognize watery foods such as soup. This seems unresponsive due to the lack of liquid modeling in training synthetic data, but it is simply overcome by fine-tuning with real data. Also the fine-tuned model shows the advantage of robustness not mistaking in the background compared to the model trained with real data only.

B. Result on public dataset

The UNIMIB2016 [13] has been collected in a real canteen environment. The images contain food on their plates and are also placed outside their plates. In some cases, there are several foods on a plate. The UNIMIB2016 is a dataset for food instance segmentation that captures food from the top view. The UNIMIB2016 [13] is composed of 1,010 tray images with multiple foods and containing 73 food categories. The 1,010 tray images are split into a training set and a test set to contain about 70% and 30% of each food instance, resulted in 650 tray image training sets and 360 image test sets. Although the UNIMIB2016 [13] contains the categories of each food, we utilize all data as single category, food, for comparison with our unseen food segmentation performance.

We conducted experiments using synthetic data, UNIMIB2016 [13] as real-world data, fine-tuning with real-world data after pre-training on synthetic data, and the results can be seen through Table II. When the network was trained with only the synthetic data, mAP was 32.9. Because some data of UNIMIB2016 [13] dataset is several food closely attached on a same plate, Although the network did not train with foods in the UNIMIB2016 [13], network can implement food instance segmentation as shown in Figure 6. Unlike synthetic data, because some data in the UNIMIB2016 [13] dataset multiple foods are clustered together on the same plate, the model trained on synthetic data tends to recognize foods on a single plate as one instance. Despite, using real-world data shows better results than using the synthetic data, in the case of training with fine-tuning with real-world data after pre-training on synthetic dataset, the highest result was obtained with 82.7% in terms of mAP. As a result, training on synthetic dataset is applicable to real-world data via Sim-to-Real and also takes a roll of general feature extraction that is more appropriate for fine-tuning as task-specific adaption.

IV. CONCLUSIONS

In this paper, we demonstrate the possibility of food instance segmentation that have never been seen in real-world environment through synthetic data generation and training of Mask R-CNN [12] model. On our real-world dataset, food instances can be segmented sufficiently with a performance of 52.2% as using a network learned from only synthetic data. Also, when fine-tuning a model learned from only synthetic data with real-world data, +6.4%p performance is improved better than the model trained from scratch. Experiments on public dataset(UNIMIB 2016 [13]) show that it is sufficient to segment food, even if it is not the same meal tray. Since this work can distinguish between different food instances but cannot recognize the type of food, it is also remaining challenge to expand intelligence for recognition of food categories. We suggest a study as our future work, transferring knowledge from classification intelligence that can be implemented with relatively easy to collect data to recognize the category of mask instance in our food instance segmentation models.



Fig. 6. Inference examples of segmentation results

REFERENCES

- [1] T. Kelly, W. Yang, C.-S. Chen, K. Reynolds, and J. He, "Global burden of obesity in 2005 and projections to 2030," *International journal of obesity*, vol. 32, no. 9, pp. 1431–1437, 2008.
- [2] F. Zhu, M. Bosch, N. Khanna, C. J. Boushey, and E. J. Delp, "Multiple hypotheses image segmentation and classification with application to dietary assessment," *IEEE journal of biomedical and health informatics*, vol. 19, no. 1, pp. 377–388, 2014.
- [3] Y. Matsuda, H. Hoashi, and K. Yanai, "Recognition of multiple-food images by detecting candidate regions," in *2012 IEEE International Conference on Multimedia and Expo*, pp. 25–30, IEEE, 2012.
- [4] M. Anthimopoulos, J. Dehais, P. Diem, and S. Mougiakakou, "Segmentation and recognition of multi-food meal images for carbohydrate counting," in *13th IEEE International Conference on Bioinformatics and BioEngineering*, pp. 1–4, IEEE, 2013.
- [5] G. Ciocca, P. Napolitano, and R. Schettini, "Food recognition: a new dataset, experiments, and results," *IEEE journal of biomedical and health informatics*, vol. 21, no. 3, pp. 588–598, 2016.
- [6] S. Fang, C. Liu, K. Tahboub, F. Zhu, E. J. Delp, and C. J. Boushey, "ctada: The design of a crowdsourcing tool for online food image identification and segmentation," in *2018 IEEE Southwest Symposium on image analysis and interpretation (SSIAI)*, pp. 25–28, IEEE, 2018.
- [7] S. Inunganbi, A. Seal, and P. Khanna, "Classification of food images through interactive image segmentation," in *Asian Conference on Intelligent Information and Database Systems*, pp. 519–528, Springer, 2018.
- [8] W. Shimoda and K. Yanai, "Cnn-based food image segmentation without pixel-wise annotation," in *International Conference on Image Analysis and Processing*, pp. 449–457, Springer, 2015.
- [9] M. Bolanos and P. Radeva, "Simultaneous food localization and recognition," in *2016 23rd International Conference on Pattern Recognition (ICPR)*, pp. 3140–3145, IEEE, 2016.
- [10] F. Golemo, A. A. Taiga, A. Courville, and P.-Y. Oudeyer, "Sim-to-real transfer with neural-augmented robot simulation," in *Conference on Robot Learning*, pp. 817–828, PMLR, 2018.
- [11] X. B. Peng, M. Andrychowicz, W. Zaremba, and P. Abbeel, "Sim-to-real transfer of robotic control with dynamics randomization," in *2018 IEEE international conference on robotics and automation (ICRA)*, pp. 3803–3810, IEEE, 2018.

- [12] K. He, G. Gkioxari, P. Dollár, and R. Girshick, "Mask r-cnn," in *2017 IEEE International Conference on Computer Vision (ICCV)*, pp. 2980–2988, 2017.
- [13] G. Ciocca, P. Napolitano, and R. Schettini, "Food recognition: a new dataset, experiments, and results," *IEEE journal of biomedical and health informatics*, vol. 21, no. 3, pp. 588–598, 2016.
- [14] T. Joutou and K. Yanai, "A food image recognition system with multiple kernel learning," in *Proceedings of the 16th IEEE International Conference on Image Processing, ICIP'09*, p. 285–288, IEEE Press, 2009.
- [15] M. Chen, K. Dhingra, W. Wu, L. Yang, R. Sukthankar, and J. Yang, "Pfid: Pittsburgh fast-food image dataset," in *2009 16th IEEE International Conference on Image Processing (ICIP)*, pp. 289–292, IEEE, 2009.
- [16] A. Mariappan, M. Bosch, F. Zhu, C. J. Boushey, D. A. Kerr, D. S. Ebert, and E. J. Delp, "Personal dietary assessment using mobile devices," in *Computational Imaging VII*, vol. 7246, p. 72460Z, International Society for Optics and Photonics, 2009.
- [17] H. Hoashi, T. Joutou, and K. Yanai, "Image recognition of 85 food categories by feature fusion," in *2010 IEEE International Symposium on Multimedia*, pp. 296–301, IEEE, 2010.
- [18] M.-Y. Chen, Y.-H. Yang, C.-J. Ho, S.-H. Wang, S.-M. Liu, E. Chang, C.-H. Yeh, and M. Ouhyoung, "Automatic chinese food identification and quantity estimation," in *SIGGRAPH Asia 2012 Technical Briefs*, pp. 1–4, 2012.
- [19] Y. Matsuda, H. Hoashi, and K. Yanai, "Recognition of multiple-food images by detecting candidate regions," in *2012 IEEE International Conference on Multimedia and Expo*, pp. 25–30, IEEE, 2012.
- [20] L. Bossard, M. Guillaumin, and L. Van Gool, "Food-101—mining discriminative components with random forests," in *European conference on computer vision*, pp. 446–461, Springer, 2014.
- [21] Y. Kawano and K. Yanai, "Automatic expansion of a food image dataset leveraging existing categories with domain adaptation," in *European Conference on Computer Vision*, pp. 3–17, Springer, 2014.
- [22] G. M. Farinella, D. Allegra, M. Moltisanti, F. Stanco, and S. Battiato, "Retrieval and classification of food images," *Computers in biology and medicine*, vol. 77, pp. 23–39, 2016.
- [23] J. Chen and C.-W. Ngo, "Deep-based ingredient recognition for cooking recipe retrieval," in *Proceedings of the 24th ACM international conference on Multimedia*, pp. 32–41, 2016.
- [24] G. Ciocca, P. Napolitano, and R. Schettini, "Learning cnn-based features for retrieval of food images," in *International Conference on Image Analysis and Processing*, pp. 426–434, Springer, 2017.
- [25] G. Ciocca, P. Napolitano, and R. Schettini, "Cnn-based features for retrieval and classification of food images," *Computer Vision and Image Understanding*, vol. 176, pp. 70–77, 2018.
- [26] E. Aguilar, M. Bolaños, and P. Radeva, "Regularized uncertainty-based multi-task learning model for food analysis," *Journal of Visual Communication and Image Representation*, vol. 60, pp. 360–370, 2019.
- [27] W. Min, L. Liu, Z. Luo, and S. Jiang, "Ingredient-guided cascaded multi-attention network for food recognition," in *Proceedings of the 27th ACM International Conference on Multimedia*, pp. 1331–1339, 2019.
- [28] P. Kaur, K. Sikka, W. Wang, S. Belongie, and A. Divakaran, "Foodx-251: a dataset for fine-grained food classification," *arXiv preprint arXiv:1907.06167*, 2019.
- [29] X. Chen, Y. Zhu, H. Zhou, L. Diao, and D. Wang, "Chinesefoodnet: A large-scale image dataset for chinese food recognition," *arXiv preprint arXiv:1705.02743*, 2017.
- [30] G. Ciocca, P. Napolitano, and R. Schettini, "Food recognition and left-over estimation for daily diet monitoring," in *International Conference on Image Analysis and Processing*, pp. 334–341, Springer, 2015.
- [31] A. Meyers, N. Johnston, V. Rathod, A. Korattikara, A. Gorban, N. Silberman, S. Guadarrama, G. Papandreou, J. Huang, and K. P. Murphy, "Im2calories: towards an automated mobile vision food diary," in *Proceedings of the IEEE International Conference on Computer Vision*, pp. 1233–1241, 2015.
- [32] J. Gao, W. Tan, L. Ma, Y. Wang, and W. Tang, "Musefood: Multi-sensor-based food volume estimation on smartphones," in *2019 IEEE SmartWorld, Ubiquitous Intelligence & Computing, Advanced & Trusted Computing, Scalable Computing & Communications, Cloud & Big Data Computing, Internet of People and Smart City Innovation (SmartWorld/SCALCOM/UIC/ATC/CBDCOM/IOP/SCI)*, pp. 899–906, IEEE, 2019.
- [33] S. Aslan, G. Ciocca, D. Mazzini, and R. Schettini, "Benchmarking algorithms for food localization and semantic segmentation," *International Journal of Machine Learning and Cybernetics*, vol. 11, no. 12, pp. 2827–2847, 2020.
- [34] S. Back, J. Kim, R. Kang, S. Choi, and K. Lee, "Segmenting unseen industrial components in a heavy clutter using rgb-d fusion and synthetic data," in *2020 IEEE International Conference on Image Processing (ICIP)*, pp. 828–832, IEEE, 2020.
- [35] M. Danielczuk, M. Matl, S. Gupta, A. Li, A. Lee, J. Mahler, and K. Goldberg, "Segmenting unknown 3d objects from real depth images using mask r-cnn trained on synthetic data," in *2019 International Conference on Robotics and Automation (ICRA)*, pp. 7283–7290, IEEE, 2019.
- [36] Blender Online Community, *Blender - a 3D modelling and rendering package*. Blender Foundation, Blender Institute, Amsterdam, Mon 08/06/2018.
- [37] R. Girshick, "Fast r-cnn," in *Proceedings of the IEEE international conference on computer vision*, pp. 1440–1448, 2015.
- [38] A. Paszke, S. Gross, S. Chintala, G. Chanan, E. Yang, Z. DeVito, Z. Lin, A. Desmaison, L. Antiga, and A. Lerer, "Automatic differentiation in pytorch," in *NIPS-W*, 2017.
- [39] T.-Y. Lin, M. Maire, S. Belongie, J. Hays, P. Perona, D. Ramanan, P. Dollár, and C. L. Zitnick, "Microsoft coco: Common objects in context," in *European conference on computer vision*, pp. 740–755, Springer, 2014.