# Image Based Food Energy Estimation
# With Depth Domain Adaptation

Gautham Vinod    Zeman Shao    Fengqing Zhu

*School of Electrical and Computer Engineering, Purdue University, West Lafayette, Indiana USA*

*Abstract*—Assessment of dietary intake has primarily relied on self-report instruments, which are prone to measurement errors. Dietary assessment methods have increasingly incorporated technological advances particularly mobile, image based approaches to address some of these limitations and further automation. Mobile, image-based methods can reduce user burden and bias by automatically estimating dietary intake from eating occasion images that are captured by mobile devices. In this paper, we propose an "Energy Density Map" which is a pixel-to-pixel mapping from the RGB image to the energy density of the food. We then incorporate the "Energy Density Map" with an associated depth map that is captured by a depth sensor to estimate the food energy. The proposed method is evaluated on the Nutrition5k dataset. Experimental results show improved results compared to baseline methods with an average error of 13.29 kCal and average percentage error of 13.57% between the ground-truth and the estimated energy of the food.

*Index Terms*—Food Portion Estimation, Energy Density Map, Dietary Analysis, Image-based Food Portion, Depth Map, Depth Domain Adaptation

## I. Introduction

The ubiquity of smartphones and their integration into every facet of life has made it possible to monitor one's health status such as vital signs, exercise and sleep statistics, etc. What and how much a person eats and drinks in a day is another key contributor to one's overall health, yet very difficult to assess [1]. Many existing mobile applications rely on manual user input to capture information about the foods they consumed [2]. This follows the traditional dietary assessment method called the 24 hour Dietary Recalls where the participant either communicates with a dietitian or use a web-based tool to report their food intake in the last 24-hour period [3] which adds user burden and measurement error. For example, study has shown that multiple 24-hour recalls are needed to accurately capture the dietary intake [4]. The development of image-based dietary assessment methods, particularly those using mobile devices, have the potential to reduce user burden and improve reporting bias compared to traditional approaches [5], [6].

Image-based methods have achieved impressive results in many applications for dietary assessment such as food recognition, food segmentation and food portion estimation [7]–[11]. In this paper, we focus on food energy estimation which is challenging even for humans including domain experts (*e.g.*, dietitians) to accurately estimate the energy of foods in an image without known physical references [12], [13]. Food energy estimation is closely related to image-based food portion estimation, which can be classified into three categories [7]:

1) **Single view image methods** where only a single view image of the eating scene is used for food portion estimation.
2) **Multi-view images methods** where multi-view images of the eating scene are captured to extract depth information which is then used for food portion estimation. Usually, the depth is estimated from multiple images of the same scene to recover some of the 3D information of the foods from the 2D images.
3) **Depth based methods** where a depth map is captured by a depth sensor and aligned with the RGB image to estimate food portion size jointly.

Relying only on a single-view image to estimate the food portion is a challenging task since the 3D information is lost during the projection from the real world coordinates to the image coordinates. Multi-view images methods typically increases user burden as the participant needs to capture images of the eating scene from different angles. Furthermore, multi-view images require additional post-processing, such as camera calibration and feature matching after image capture. Depth map captured by the depth sensor represents pixel-wise distance from the object surface to the camera, which can be used to recover the 3D information. However, depth sensor is not commonly available on mobile devices, making it difficult to deploy multi-view based methods on existing mobile devices.

Previously, we proposed a single-view based food portion estimation method by introducing the concept of an "Energy Distribution Map" [11] which is a pixel to pixel mapping of the RGB image of the food to a map of how the energy of the food is distributed. A conditional Generative Adversarial Network (cGAN) [14] is used to learn this mapping from the RGB image to the "Energy Distribution Map." In [12], we further adopted the concept of the "Energy Distribution Map" and combined its features with the features extracted from the RGB image to improve the estimation of food energy. The performance is shown to have improved when using a combination of these features as compared to using the "Energy Distribution Map" alone. The method is evaluated on a real-world dataset collected from a dietary study where the ground-truth energy values are provided by registered dietitians. However, there are only 96 eating occasion images in the dataset. In this paper, we first perform segmentation for each RGB image, and also adopt the "Energy Distribution Map" as the "Energy Density Map" which maps each pixel

in the RGB image to an energy value, to more accurately represent the energy density of the different foods in the image. We replace the RGB features with features extracted from a depth map to jointly perform food energy estimation. Our method is evaluated on a large public food image dataset, Nutrition5K, which consists of 3,493 food images with aligned depth maps captured by a depth sensor [15].

The main contributions of our paper can be summarized as follows:

1) Propose an "Energy Density Map" which maps each pixel in the RGB image to an energy value.
2) Incorporate depth information from a depth sensor with information from the "Energy Density Map" for joint learning of food energy in the image.
3) Evaluate the proposed method on the public Nutrition5k dataset [15], which shows improved performance compared to previous methods using only "Energy Density Map" and the RGB image.

## II. RELATED WORKS

Work related to food portion estimation typically tries to reconstruct 3D data from a 2D image, which is challenging because most of the 3D information is lost during the 3D to 2D projection. Existing methods are either based on extracting geometric information from the image or using a deep learning based approach.

### A. Geometric Based Methods

These methods estimate the food portion/volume from the geometry of objects in the image such as objects with known physical dimensions. For example, the shapes of containers in the image along with a physical reference object in the image such as a colored checkerboard of known dimensions are used as reference to estimate the volume of foods in the image [16]. In [17], the authors use classification and estimation of serving sizes based on certain parameters of the image such as color, number of circles in the image, etc. and then these features are fed into an AdaBoost classifier [18] to estimate the serving sizes of the food. Mobile 3D range is used in [19] where a structured laser grid is projected onto the food and a smartphone is used to capture the resulting image. The laser grid lines are used to create a 3D depth map of the food. Other methods try to estimate the 3D scene using multiple images [20]–[22]. The food volume is estimated in [23] based on a reference object which is a circular dining plate and a user defined 3D model which represents the food image. This 3D model is projected onto the 2D image. The model uses the circular plate and the chosen 3D model to estimate the food volume.

### B. Deep Learning Based Methods

Most deep learning based methods use Neural Networks to estimate the food portion. In [24] a CNN network is used for image classification and another network is used for food attribute estimation using vector space embedding [25]. In [26], a CNN is used for food classification and then the portion

is estimated based on a physical reference object in the image or the mobile camera's sensor measurements. A top and side view image are used in [27] along with the Faster R-CNN network [28] to identify and localize the food and a calibration object from these images. Segmentation is performed next to calculate the food volume based on the calibration object in the image. The concept of energy value per unit mass was used in [15] and [29], which treats portion estimation as a classification task since each food has its own energy per unit mass which is a property inherent to the specific food.

## III. METHOD

### A. Overview of Our Method

Our method is based on the idea that energy of foods in an image can be estimated using the energy density, which is the energy per pixel in the image. We also incorporate the depth map which contains the physical distance information between the camera and the objects in the image. By keeping the distance between the camera and the objects (eating scene in this case) fixed, the depth map provides information about the relative sizes and shapes of the foods in the image. Our method consists of 2 modules - an Energy Density Estimation Module, and a Feature Adaptation Module. The Energy Density Estimation Module consists of a generative model that maps the RGB image $x$ to the Energy Density Map $y$. The Feature Adaptation Module consists of 2 feature extracting networks which extract features from the generated energy density map $y$ and the depth map $z$ separately. The extracted features from the energy density map $y_f$ and the depth map $z_f$ are then normalized and concatenated. The normalized and concatenated features $(\tilde{y}_f, \tilde{z}_f)$ are then passed through a regression network to estimate the energy value of foods in the image. During training, the Energy Density Estimation Module is first trained to generate the Energy Density Map $y$. Once the generative model is trained, it estimates the Energy Density Map for all the images in our dataset which is then used as one of the inputs for the Feature Adaptation Module. The Feature Adaptation Module is trained to use the generated Energy Density Maps and the depth maps to estimate the food energy in the image. However, during the inference, an RGB image is first fed into the Energy Density Estimation Module and the output of this module along with the depth map are served as inputs for the Feature Adaptation Module, *i.e.*, the inference is end-to-end. Figure 1 shows an overview of the proposed method.

### B. Energy Density Map

The Energy Density Map $y$ is a 1-channel image where each pixel value is representative of the energy of the food at that pixel location in the image. During the energy estimation we use a conditional Generative Adversarial Netowrk (cGAN) [30] to generate this Energy Density Map. In order to train this generative model we need to first obtain the "ground-truth" energy density map.

To create the ground-truth energy density map for training, we first segment the food images so that we know the location
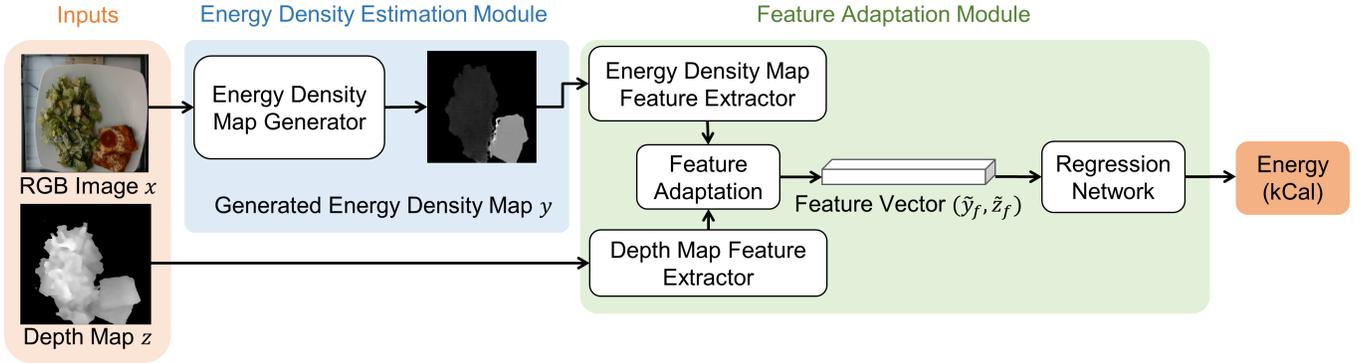
Fig. 1. **Overview of proposed method.** The Energy Density Estimation Module uses a generative model to estimate the Energy Density Map $y$ from the RGB input $x$. The Feature Adaptation Module extracts the features from the Energy Density Map $y$ and the Depth Map $z$, separately and combines them after normalizing the features. Finally, a regression network produces a single output value which is the estimated energy value of the food in the input image.

of each individual food in the image. The segmentation step is necessary because the energy density of each food or ingredient can be different. For example, if we consider a small piece of meat and a small piece of spinach, both of the same size, the ground-truth energy value of the meat should be higher than that of the spinach. We use Google's Seefood mobile food segmenter [31] to perform the food segmentation, which adopts the DeepLab-V3 [32] network architecture with a MobileNet-v2 [33] backbone.

Suppose a food image has $N$ foods or ingredients, then let $k$ denote the $k$-th food item in the image, $k \in [1, N]$. Let $e_k$ denote the energy value of the $k$-th food or ingredient in the image. Let pixel $x(i, j)$ correspond to a pixel in row $i$ and column $j$ of the RGB image $x$. The segmentation network gives us the pixel locations of each food $k$ in the image. Further, if $x(i, j)$ contains a food or ingredient $k$, then the corresponding pixel $y(i, j)$ of the energy density map has a value that is equal to

$$y(i, j) = \frac{e_k}{\text{total number of pixels in food } k}$$

This process is repeated for each food or ingredient in the image. Once all the images in the dataset have a corresponding energy density map, these maps are scaled to the range of $[0, 255]$ over the whole dataset. Therefore, this energy density map can be described as the distribution of a scaled value of the ground-truth food energy over the image.

We used the same training method as described in [12] where the RGB image serves as the input to the generative model which outputs a generated Energy Density Map.

### C. Joint Learning from Energy Density and Depth Map

The energy density map $y$ provides information about how the energy of the foods are distributed over the area occupied by the foods in the image. Previously, we extracted features from the Energy Density Map and the RGB image and then normalized and combined these features. The combination of these features are then used as inputs for a regression network to produce a final estimate of the food portion. In this paper,

the features extracted from the depth map $z$ which contains information about the depth and shape of the food in the image, and hence we use these features instead of the features from the RGB image. We also evaluate the performance of the energy estimation when considering all three features, *i.e.*, energy density, depth and RGB.

The depth map in the Nutrition5k Dataset is a 1 channel image with the same dimensions as the RGB image. Each pixel in the depth map is in the range of $[0, 65, 535]$, which represents a scaled estimate of the distance between the camera and the object in that pixel. The scale factor is $1 \times 10^{-4}$ m, which means that if an object is $10\,\text{cm}$ away from the camera, then the corresponding pixels in the image have a value of $1,000$. We apply post-processing to the depth map from the depth sensor including dilation to smooth the foreground and morphological closing to fill any missing values in the depth map. The energy density map and the depth map are then normalized to $[-1, 1]$ so that the features extracted from them can be concatenated.

### D. Food Energy Estimation

The feature adaption module contains 2 feature extractors - The Energy Density Map feature extractor and the Depth Map feature extractor. Both use the VGG-16 [34] network as their backbone. Previously, VGG-16 was used as the backbone to extract features from the Energy Density Map [12] and Resnet-50 was used as the RGB feature extractor [35]. Since extracting features from a single channel depth map is easier than extracting features from an RGB image, VGG-16 is sufficient in this case. The final fully-connected layers of the VGG-16 networks are removed. In each case, for the depth map and for the energy density map, the output of each VGG-16 feature extractor is a $7 \times 7 \times 512$ feature tensor.

The feature tensors extracted from the Depth map and the Energy Density map are concatenated in the feature adaptation module. Since the features are from different domains, they are normalized before concatenation. We tried different methods of normalization such as Z-score normalization, Layer

Fig. 2. **Erroneous ground-truth data.** (Left) The ground-truth ingredients only contain Cantaloupes and Tomatoes but the image also contains grapes. (Right) Overlapping dishes where the food is not clearly visible because of improper image capture.

Normalization [36] (layer normalization helps in normalizing the distributions of intermediate layers) in the regression network and a combination of Layer Normalization and Group Normalization (group normalization organizes channels into different groups and normalizes the distribution among them) [37]. Finally, these normalized features are passed through the same network used in [12] which consists of 2 fully connected layers with normalization layers in between, which outputs a single value of the estimated energy of the food.

## IV. EXPERIMENTAL RESULTS

### A. Dataset Curation

The Nutrition5k Dataset has 3,493 images that have over-head depth map data. In these 3,493 images, there are many dishes that have a lot of individual ingredients including many that are hard for a segmentation network to detect such as ketchup, mayonnaise, salt and pepper, oil, etc. As we mentioned in the previous section, the generation of the energy density map requires segmentation of the image. One of the limitations of the See-food segmentation network [31] is that the output classes are vague categories such as "Leafy Greens", "Starchy Vegetables", etc. making it difficult to group the ingredients in the dataset into these categories. Therefore, to address these challenges we choose a subset of the Nutrition5k Dataset where the segmentation network can achieve reasonable results. This subset of images consists of images with less than 3 food ingredients. In addition, we also removed images where one ingredient is completely covered by another ingredient and ingredients that do not have consistent appearance in all images such as tofu.

There are a few instances where the ground-truth data provided in the dataset is incorrect such as the images in Figure 2, which are excluded from our subset. We also removed images where the ground-truth energy value of the dish is below 10 kCal since there are only few images in the subset that are in this range of $0 - 10$ kCal. It would be challenge for the network to learn from limited data in this energy range. In the end, our subset contains 909 images. Figure 3 shows the distribution of the ground-truth energy values of the dataset (images with depth map available) and our subset. The energy values in our subset approximately range from 10 kCal to 1,000 kCal while the range of the energy in the original dataset
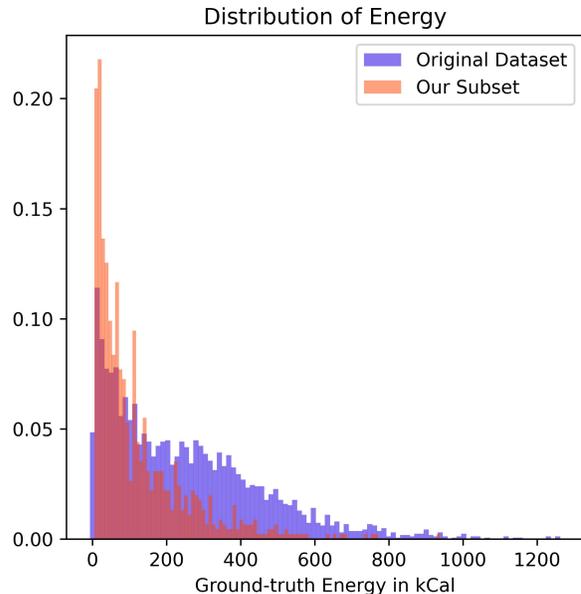


Fig. 3. **Comparison of the distribution of energy in the datasets.** The graph shows that our subset has a higher concentration of images in the low ground-truth energy region as compared to the original dataset. The distribution also shows that the range of energy in our subset covers majority of the energy values in the original dataset.

excluding the outliers (1.4% of the data which is more than 3 standard deviations away from the mean) is around $[0, 800]$ kCal. Excluding the outliers in the original dataset, the range of our subset covers the majority of energy values in the original dataset. Since our subset has a good representation of energy range in the original dataset, it is suitable for evaluating the performance of the proposed method.

### B. Experimental Results on the Nutrition5k Dataset

The feature extractor for both the energy density map and the depth map use the VGG-16 network that is pre-trained on the ImageNet dataset [38]. The feature extractor and the regression network are trained together, *i.e.*, the entire feature adaptation module in Figure 1 is trained together. The network is supervised on the sum of the L1 loss and the Mean Squared Error loss between the estimated and ground-truth energy of the food. The Adam optimizer [39] is used with an initial learning rate of $5e^{-5}$. The output of the network is a scaled version of the estimated energy. We used 300 as the scaling factor for the output of the network, therefore the output of the network is the estimated energy scaled by our scaling factor. This output is scaled back to represent the estimated energy value. The scaling factor was determined experimentally.

The training and testing split is chosen by stratifying the 909 images so that the training data and the testing data have a similar distribution. An 80:20 stratified split is used resulting a total of 731 training images and 175 testing images. The learning rate is reduced by a factor of 0.8 for every 10 epochs and all models were trained for approximately 50 epochs each.

| Method | MAE (kCal) | MAPE (%) |
|--------|-----------|----------|
| $x_f$ | 26.85 | 40.64 |
| $y_f$ | 13.35 | 16.90 |
| $z_f$ | 76.86 | 133.88 |
| $(x_f, y_f)$ | 17.54 | 23.20 |
| $(\tilde{x}_f, \tilde{y}_f)$ | 14.65 | 16.88 |
| $(y_f, z_f)$ | 15.83 | 27.04 |
| $(\tilde{y}_f, \tilde{z}_f)$ | 13.29 | 13.57 |
| $(\tilde{x}_f, \tilde{y}_f, \tilde{z}_f)$ | 12.75 | 16.83 |

We used the Mean Absolute Error (MAE) and Mean Absolute Percentage Error (MAPE) to measure the performance which are defined as

$$\text{MAE} = \frac{1}{N} \sum_{i=1}^{N} |\hat{e} - e|$$

$$\text{MAPE} = \frac{1}{N} \sum_{i=1}^{N} \frac{|\hat{e} - e|}{e} \times 100$$

where $\hat{e}$ is the estimated energy of all foods in the image, $e$ is the ground-truth energy value, and $N$ is the total number of images in the test set. The results in Table I shows the performance of our proposed method. The first three row show the baseline results using the features extracted from the RGB image $x_f$, the energy density map $y_f$, and the depth map $z_f$, respectively. The next two rows show the the combined features from the RGB image $x_f$ and the energy density map $y_f$ with and without feature adaptation. We then show the results of combining the energy density map $y_f$ and the depth map $z_f$ with and without feature adaptation. Finally, we show the result for combined features from the RGB image, energy density map and the depth map.

From the results in Table I, we observed that:

1) From the first three rows, it is clear that features from the Energy Density Map contributes the most to reducing the energy estimation error.
2) The depth map by itself does not have sufficient information to accurately estimate the energy of the food. Since the depth and the energy density maps are in different domains, combining them without normalization degrades the performance compared to using the Energy Density Map alone.
3) The best performance is obtained when the features from the depth map and the energy density map are combined and normalized, resulting in a MAPE of 13.57%. When the RGB, depth and, energy density map features are combined, the MAE is lower than when only the depth map and the Energy Density Map are combined. However, MAPE is a better indicator of

performance in this subset since the low ground-truth energy values of the food push the MAPE to be higher even for a small difference in MAE.

### C. Comparison to Nutrition5k Results

Three portion estimation methods are reported using the Nutrition5k dataset in [15].

1) 2D Direct Prediction - The RGB image is used as an input to a regression network with the Inception v3 [40] network as its backbone.
Reported MAE (kCal) / MAPE: 70.6 / 26.1%
2) Depth as the 4-th channel - The depth map is added as another channel to the RGB input and then this RGB-D image is sampled to form a 3-channel input to the regression network.
Reported MAE (kCal) / MAPE: 47.6 / 18.8%
3) Volume Scalar - The mass of the food is approximated using certain physical approximations and the depth map, and then multiplied with a network that predicts the calories per gram of the food.
Reported MAE (kCal) / MAPE: 41.3 / 16.5%

The mean energy value (124.96 kCal) in our subset is lower than the full dataset (254.94 kCal) and hence we would expect our MAPE to be higher in our subset. This is because an estimate of 20 kCal for a food that has a ground-truth energy of 10 kCal will produce an MAPE of 100% but an MAE of only 10 kCal. However, we see that our proposed method combining the Energy Density Map and the depth map has a lower MAPE of 13.57% than the reported results of 18.8% when depth map is used in [15] .

### V. CONCLUSION

We proposed a method that combines the information from an Energy Density Map and a depth map to estimate energy for the foods in an image, where the Energy Density Map is generated from an RGB input. Our preliminary experiments showed promising results on a subset of the Nutrition5K dataset. Our future work will focus on evaluating our method on the complete Nutrition5K dataset and food images captured in dietary studies.

### REFERENCES

[1] L. M. König, M. V. Emmenis, J. Nurmi, A. Kassavou, and S. Sutton, "Characteristics of Smartphone-based Dietary Assessment Tools: A Systematic Review," *Health Psychology Review*, vol. 0, no. 0, pp. 1–25, 2021.
[2] M. Zečević, D. Mijatović, M. Kos Koklič, V. Žabkar, and P. Gidaković, "User Perspectives of Diet-Tracking Apps: Reviews Content Analysis and Topic Modeling," *Journal of Medical Internet Research*, vol. 23, no. 4, p. e25160, Apr 2021.
[3] A. F. Subar, S. I. Kirkpatrick, B. Mittl, T. P. Zimmerman, F. E. Thompson, C. Bingley, G. Willis, N. G. Islam, T. Baranowski, S. McNutt *et al.*, "The Automated Self-administered 24-Hour Dietary Recall (ASA24): A Resource for Researchers, Clinicians and Educators from the National Cancer Institute," *Journal of the Academy of Nutrition and Dietetics*, vol. 112, no. 8, p. 1134, 2012.
[4] Y. Ma, B. C. Olendzki, S. L. Pagoto, T. G. Hurley, R. P. Magner, I. S. Ockene, K. L. Schneider, P. A. Merriam, and J. R. Hébert, "Number of 24-Hour Diet Recalls Needed to Estimate Energy Intake," *Annals of Epidemiology*, vol. 19, no. 8, pp. 553–559, August 2009.

[5] C. J. Boushey, M. Spoden, F. M. Zhu, E. J. Delp, and D. A. Kerr, "New Mobile Methods for Dietary Assessment: Review of Image-assisted and Image-based Dietary Assessment Methods," *Proceedings of the Nutrition Society*, vol. 76, no. 3, p. 283–294, 2017.

[6] F. Zhu, M. Bosch, I. Woo, S. Kim, C. J. Boushey, D. S. Ebert, and E. J. Delp, "The Use of Mobile Devices in Aiding Dietary Assessment and Evaluation," *IEEE Journal of Selected Topics in Signal Processing*, vol. 4, no. 4, pp. 756–766, 2010.

[7] F. P. W. Lo, Y. Sun, J. Qiu, and B. Lo, "Image-Based Food Classification and Volume Estimation for Dietary Assessment: A Review," *IEEE Journal of Biomedical and Health Informatics*, vol. 24, no. 7, pp. 1926–1939, July 2020.

[8] J. He, R. Mao, Z. Shao, J. L. Wright, D. A. Kerr, C. J. Boushey, and F. Zhu, "An End-to-end Food Image Analysis System," *Electronic Imaging*, vol. 33, no. 8, 2021.

[9] S. K. Yarlagadda, D. M. Montserrat, D. Güera, C. J. Boushey, D. A. Kerr, and F. Zhu, "Saliency-Aware Class-Agnostic Food Image Segmentation," *ACM Transactions on Computing for Healthcare*, vol. 2, no. 3, jul 2021.

[10] R. Mao, J. He, Z. Shao, S. K. Yarlagadda, and F. Zhu, "Visual Aware Hierarchy Based Food Recognition," *Proceedings of the International Conference on Pattern Recognition Workshop*, pp. 571–598, 2021.

[11] S. Fang, Z. Shao, R. Mao, C. Fu, E. J. Delp, F. Zhu, D. A. Kerr, and C. J. Boushey, "Single-View Food Portion Estimation: Learning Image-to-Energy Mappings Using Generative Adversarial Networks," *Proceedings of the IEEE International Conference on Image Processing*, pp. 251–255, 2018.

[12] Z. Shao, S. Fang, R. Mao, J. He, J. L. Wright, D. A. Kerr, C. J. Boushey, and F. Zhu, "Towards Learning Food Portion From Monocular Images With Cross-Domain Feature Adaptation," *Proceedings of the IEEE International Workshop on Multimedia Signal Processing*, pp. 1–6, Oct 2021.

[13] C. D. Lee, J. Chae, T. E. Schap, D. A. Kerr, E. J. Delp, D. S. Ebert, and C. J. Boushey, "Comparison of Known Food Weights with Image-Based Portion-Size Automated Estimation and Adolescents' Self-Reported Portion Size," *Journal of Diabetes Science and Technology*, vol. 6, no. 2, pp. 428–434, 2012, pMID: 22538157.

[14] M. Mirza and S. Osindero, "Conditional Generative Adversarial Nets," *arXiv preprint arXiv:1411.1784*, 2014.

[15] Q. Thames, A. Karpur, W. Norris, F. Xia, L. Panait, T. Weyand, and J. Sim, "Nutrition5k: Towards automatic nutritional understanding of generic food," *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pp. 8903–8911, June 2021.

[16] S. Fang, C. Liu, F. Zhu, E. J. Delp, and C. J. Boushey, "Single-View Food Portion Estimation Based on Geometric Models," *Proceedings of the IEEE International Symposium on Multimedia*, pp. 385–390, 2015.

[17] K. Aizawa, Y. Maruyama, H. Li, and C. Morikawa, "Food Balance Estimation by Using Personal Dietary Tendencies in a Multimedia Food Log," *IEEE Transactions on Multimedia*, vol. 15, no. 8, pp. 2176–2185, 2013.

[18] Y. Freund, R. E. Schapire *et al.*, "Experiments with a New Boosting Algorithm," *Proceedings of the Thirteenth International Conference on International Conference on Machine Learning*, vol. 96, pp. 148–156, 1996.

[19] J. Shang, M. Duong, E. Pepin, X. Zhang, K. Sandara-Rajan, A. Mamishev, and A. Kristal, "A Mobile Structured Light System for Food Volume Estimation," *Proceedings of the IEEE International Conference on Computer Vision Workshops*, pp. 100–101, 2011.

[20] M. Puri, Z. Zhu, Q. Yu, A. Divakaran, and H. Sawhney, "Recognition and Volume Estimation of Food Intake Using a Mobile Device," *Proceedings of the IEEE Workshop on Applications of Computer Vision*, pp. 1–8, 2009.

[21] F. Kong and J. Tan, "DietCam: Automatic Dietary Assessment with Mobile Camera Phones," *Pervasive and Mobile Computing*, vol. 8, no. 1, pp. 147–163, 2012.

[22] J. Dehais, S. Shevchik, P. Diem, and S. G. Mougiakakou, "Food Volume Computation for Self Dietary Assessment Applications," *Proceedings of the IEEE International Conference on BioInformatics and BioEngineering*, pp. 1–4, 2013.

[23] H.-C. Chen, W. Jia, Z. Li, Y.-N. Sun, and M. Sun, "3D/2D Model-to-image Registration for Quantitative Dietary Assessment," *Proceedings of the 38th Annual Northeast Bioengineering Conference*, pp. 95–96, 2012.

[24] R. Yunus, O. Arif, H. Afzal, M. F. Amjad, H. Abbas, H. N. Bokhari, S. T. Haider, N. Zafar, and R. Nawaz, "A Framework to Estimate the Nutritional Value of Food in Real Time Using Deep Learning Techniques," *IEEE Access*, vol. 7, pp. 2643–2652, 2019.

[25] T. Mikolov, K. Chen, G. Corrado, and J. Dean, "Efficient Estimation of Word Representations in Vector Space," *arXiv preprint arXiv:1301.3781*, 2013.

[26] P. Pouladzadeh, P. Kuhad, S. V. B. Peddi, A. Yassine, and S. Shirmohammadi, "Food Calorie Measurement using Deep Learning Neural Network," *Proceedings of the IEEE International Instrumentation and Measurement Technology Conference*, pp. 1–6, 2016.

[27] Y. Liang and J. Li, "Deep Learning-Based Food Calorie Estimation Method in Dietary Assessment," *Computing Research Repository*, vol. abs/1706.04062, 2017.

[28] S. Ren, K. He, R. Girshick, and J. Sun, "Faster R-CNN: Towards Real-Time Object Detection with Region Proposal Networks," *Proceedings of the 28th International Conference on Neural Information Processing Systems - Volume 1*, vol. 28, 2015.

[29] P. Ma, C. P. Lau, N. Yu, A. Li, and J. Sheng, "Application of Deep Learning for Image-based Chinese Market Food Nutrients Estimation," *Food Chemistry*, vol. 373, p. 130994, 2022.

[30] P. Isola, J.-Y. Zhu, T. Zhou, and A. A. Efros, "Image-to-Image Translation with Conditional Adversarial Networks," *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, pp. 1125–1134, July 2017.

[31] Google, "Mobile food segmentation model," https://tfhub.dev/google/seefood/segmenter/mobile_food_segmenter_V1/1.

[32] L.-C. Chen, Y. Zhu, G. Papandreou, F. Schroff, and H. Adam, "Encoder-decoder with Atrous Separable Convolution for Semantic Image Segmentation," *Proceedings of the European Conference on Computer Vision*, pp. 801–818, August 2018.

[33] M. Sandler, A. Howard, M. Zhu, A. Zhmoginov, and L.-C. Chen, "Mobilenetv2: Inverted Residuals and Linear Bottlenecks," *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, pp. 4510–4520, January 2018.

[34] K. Simonyan and A. Zisserman, "Very Deep Convolutional Networks for Large-scale Image Recognition," *arXiv preprint arXiv:1409.1556*, 2014.

[35] K. He, X. Zhang, S. Ren, and J. Sun, "Deep Residual Learning for Image Recognition," *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, pp. 770–778, June 2016.

[36] J. L. Ba, J. R. Kiros, and G. E. Hinton, "Layer Normalization," *arXiv preprint arXiv:1607.06450*, 2016.

[37] Y. Wu and K. He, "Group Normalization," *Proceedings of the European Conference on Computer Vision*, pp. 3–19, 2018.

[38] J. Deng, W. Dong, R. Socher, L.-J. Li, K. Li, and L. Fei-Fei, "ImageNet: A Large-scale Hierarchical Image Database," *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, pp. 248–255, June 2009.

[39] D. P. Kingma and J. Ba, "Adam: A method for Stochastic Optimization," *arXiv preprint arXiv:1412.6980*, 2014.

[40] C. Szegedy, W. Liu, Y. Jia, P. Sermanet, S. Reed, D. Anguelov, D. Erhan, V. Vanhoucke, and A. Rabinovich, "Going Deeper with Convolutions," *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, pp. 1–9, June 2015.