# Seeing Through the Plate: Estimating Nutrition from Multiview Food Images

1st Andrés Lomelí

*AI*

*Universidad Panamericana*

Aguascalientes, Mexico

0242956@up.edu.mx

il address or ORCID

*Abstract*—**Accurate nutritional estimation from food images remains a challenging task due to occlusions, lack of depth cues, and the high variability of food presentation. In this paper, we propose a novel, sensor-free, multi-view framework for predicting detailed macronutrient values—calories, carbohydrates, proteins, fats, and weight—from RGB images alone.**

*Index Terms*—**Classification, Carbohydrates, Detector, food**

## I. INTRODUCTION

Accurate dietary assessment plays a crucial role in addressing global health issues such as obesity, diabetes, and cardiovascular diseases. Manual tracking of food intake is time-consuming, error-prone, and highly subjective, often leading to underreporting or misestimation of caloric and nutritional values. With the ubiquity of smartphones and advancements in computer vision, automated food analysis systems based on images have emerged as promising tools to support healthier lifestyles and dietary self-monitoring.

Estimating the nutritional content of a meal from an image, however, is a highly challenging task. Visual information alone is often insufficient to determine the exact type, portion size, or composition of a dish. Different ingredients can appear similar, and the same dish may vary significantly across cultures, restaurants, or individual preparations. Moreover, the volume of food—crucial for calorie estimation—is difficult to infer from a single 2D image without depth cues or scale references.

To overcome these challenges, researchers have developed a wide range of approaches combining deep learning, computer vision, and nutritional databases. Some models classify food types and infer average values from known datasets, while others aim to estimate portion size through geometric reasoning or depth sensors. A more recent and promising direction involves directly regressing caloric and macronutrient values from food images using end-to-end deep learning models.

This paper presents a review of the state of the art in food image-based nutrition estimation, comparing various methodologies and highlighting their strengths and limitations. We also propose a novel multi-view transformer-based framework that leverages multiple RGB images per dish to improve prediction accuracy without requiring depth sensors or manual calibration

## II. STATE OF THE ART

The automatic estimation of nutritional content from food images has become a significant area of research, driven by the growing need for tools that assist individuals in monitoring their dietary intake. This section explores the evolution of methodologies in this domain, focusing on food recognition, portion size estimation, and nutritional analysis, while addressing the challenges and limitations inherent in these approaches.

### A. Food Recognition

Accurate food recognition serves as the foundation for any image-based nutritional estimation system. Early approaches relied on traditional image processing techniques, but the advent of deep learning, particularly Convolutional Neural Networks (CNNs), has significantly enhanced recognition accuracy.

*1) Convolutional Neural Networks (CNNs) for Food Classification:* CNNs have demonstrated remarkable success in image classification tasks, including food recognition. Models such as AlexNet, VGGNet, and ResNet have been employed to classify food items from images. For instance, the work by Ciocca et al. provides a comprehensive overview of deep learning techniques applied to food recognition, emphasizing the effectiveness of CNNs in handling the variability of food appearances [1]. However, food classification presents unique challenges:

- Intra-Class Variability: The same food item can appear differently due to variations in cooking methods, presentation, and lighting conditions.
- Inter-Class Similarity: Visually similar foods may belong to different categories, making differentiation difficult.

To address these challenges, hybrid models combining CNNs with other techniques have been proposed. Min et al. introduced a hybrid deep learning algorithm that integrates multiple CNN architectures to improve food recognition accuracy, demonstrating the potential of ensemble methods in this context [2].

### B. Portion Size and Volume Estimation

Estimating the portion size or volume of food items is crucial for accurate nutritional assessment. Various method-

ologies have been explored, ranging from geometric models to advanced deep learning techniques.

*1) Depth-Based Methods:* Incorporating depth information can enhance volume estimation accuracy. Pouladzadeh et al. developed a system utilizing RGB-D images (combining color and depth data) to estimate food volume, leveraging the additional spatial information provided by depth sensors [3]. Similarly, Lee and Kwon proposed a method that combines color and depth images to estimate food intake amounts, demonstrating improved accuracy over color-only approaches [4].

However, the reliance on depth sensors poses limitations:

- Hardware Requirements: Depth sensors are not commonly available in standard consumer devices, restricting the applicability of these methods.
- Environmental Constraints: Depth sensing can be affected by lighting conditions and the reflective properties of food surfaces.

*2) RGB Image-Based Methods:* To overcome hardware limitations, researchers have explored volume estimation using only RGB images. Shao et al. introduced a framework that reconstructs 3D shapes from monocular images to estimate food portions, employing domain adaptation techniques to enhance accuracy [5]. Similarly, Han et al. proposed DPF-Nutrition, an end-to-end nutrition estimation method that predicts depth maps from monocular images and fuses them with RGB data to improve portion estimation [6].

While these methods eliminate the need for specialized hardware, challenges remain:

- Depth Ambiguity: Inferring depth from a single RGB image is inherently ambiguous and prone to errors.
- Occlusions and Complex Food Arrangements: Overlapping food items and complex presentations can complicate volume estimation.

### C. Nutritional Content Estimation

Beyond recognition and volume estimation, determining the nutritional content—calories and macronutrients—of food items is the ultimate goal. This involves correlating visual information with nutritional databases or directly predicting nutritional values.

*1) Direct Prediction Models:* Thames et al. introduced the Nutrition5k dataset, comprising over 5,000 real-world dishes with detailed nutritional annotations, including calories, fats, proteins, and carbohydrates [7]. Using this dataset, they trained deep learning models to predict nutritional content directly from images, achieving accuracy levels surpassing those of professional nutritionists.

However, direct prediction models face several challenges:

- Dataset Limitations: Large, diverse, and accurately annotated datasets are essential for training robust models, but such datasets are scarce.
- Generalization: Models trained on specific datasets may not generalize well to foods from different cultures or regions.

*2) Multispectral Imaging:* To enhance accuracy, some studies have incorporated multispectral imaging. Lee et al. utilized images captured at various wavelengths, including ultraviolet and near-infrared, to improve food classification and caloric estimation, demonstrating that multispectral data can provide additional discriminative features [8]. Despite its potential, multispectral imaging is limited by:

- Equipment Accessibility: Specialized cameras are required, which are not commonly available to consumers.
- Complexity: Processing multispectral data increases computational complexity and may not be feasible for real-time applications.

### D. In-depth Review: Im2Calories

A seminal work in the field of nutritional estimation from images is Im2Calories by Meyers et al. [9] [14], which presents a comprehensive vision-based framework aimed at building an automated dietary logging tool. The goal of their system is to estimate the caloric content of food using mobile camera images, eliminating manual tracking.

Method Overview The Im2Calories pipeline consists of multiple stages designed to address the complex task of estimating calories from a single RGB image:

1) Food Detection and Segmentation: The system begins by segmenting the food from the background using a deep convolutional neural network. Instance segmentation is used to identify and isolate individual food items on a plate.
2) Food Classification: Each segmented region is classified into a predefined food category using a CNN trained on a food-specific dataset. Fine-grained classification is crucial here, as misclassification can significantly affect calorie estimates.
3) 3D Shape Estimation: A novel aspect of the Im2Calories framework is its depth estimation from RGB. Instead of using dedicated hardware, the authors train a CNN to predict depth maps from single RGB images. These depth maps are then used to infer 3D volume of the food.
4) Calorie Estimation: With known volume and class label, the system queries a food density database to convert volume into mass, and then into calories, using standardized nutritional values.

### E. Challenges and Limitations

Despite advancements, several challenges persist in the field of image-based nutritional estimation:

- Reliance on Specialized Hardware: Methods requiring depth sensors or multispectral cameras are impractical for widespread use due to hardware limitations.
- Data Availability and Quality: The scarcity of large-scale, diverse, and accurately annotated datasets hampers model training and evaluation.
- Generalization Across Diverse Cuisines: Models often struggle to generalize to foods from different cultures, limiting their applicability in

Technical Challenges and Research Implications Meyers et al [9]. explicitly identify numerous sub-problems that their system must address—many of which remain open today:

- Fine-grained recognition: Distinguishing between visually similar variants (e.g., baked vs. fried chicken).
- Open-world classification: Recognizing foods outside the training set.
- Visual attribute detection: Understanding food preparation methods (with/without sauce, grilled/fried).
- Real-time on-device inference: Making the system usable in mobile devices with limited resources.

These problems are not only important for practical deployment but are also of significant research interest across computer vision and human-computer interaction fields.

## III. PROPOSED METHOD

Building upon the open challenges identified in *Im2Calories* [9], we propose a novel, lightweight, and sensor-free framework for detailed nutritional estimation from food images. Our model aims to address key limitations of previous systems by integrating multiple RGB views of a meal, without requiring depth sensors, calibration objects, or manual segmentation. It directly predicts **calories, protein, fat, carbohydrates, and weight** using a transformer-based multi-view architecture.

### A. Motivation and Design Principles

The design of our system is guided by the following observations:

- Single-image depth estimation is unreliable, especially in occluded or cluttered food scenes.
- Mixed dishes (e.g., salads, rice with vegetables) are difficult to classify and segment manually.
- Consumers often take multiple photos of their meals (e.g., from slightly different angles), but existing models ignore this information.
- Macronutrient-level feedback is more informative than calories alone, especially for clinical applications like diabetes or muscle gain.

We propose a **multi-view nutritional estimation pipeline**, inspired by *Im2Calories'* architecture, but extending it in three significant ways:

1) We remove the need for 3D reconstruction or monocular depth prediction.
2) We allow direct regression of macronutrient values.
3) We enable multiview fusion to implicitly learn food geometry and volume.

### B. System Overview

Our proposed system consists of five main components:

*1) Multiview Image Acquisition:* We assume that the user provides **2 to 4 RGB images** of a single meal captured from different angles using a mobile phone camera. No calibration object is required.

*2) Feature Extraction and View Encoding:* Each image is passed through a shared CNN backbone (e.g., EfficientNet or MobileViT) to extract deep visual features. These features retain local texture, color, and shape information necessary for food type and portion understanding.

*3) Cross-View Attention Module (Transformer Encoder):* A **cross-view transformer encoder** performs attention across views, learning how food appearance changes with angle to infer hidden volume and structure—implicitly solving the amodal completion problem described in [9].

*4) Global Nutritional Representation:* The attention-fused representation summarizes the entire dish, encoding spatial geometry and food content. Unlike classification-based approaches, our model is **trained to regress directly** into continuous nutritional values.

*5) Macronutrient Regression Heads:* We include separate fully-connected heads to predict:

- Calories (kcal)
- Carbohydrates (g)
- Proteins (g)
- Fats (g)
- Weight (g)

Each head is trained using Mean Squared Error (MSE) loss.

### C. Training Strategy

We train our model using the **Nutrition5k** dataset [7], which provides multi-angle RGB images and ground-truth nutritional values for over 5,000 meals. During training:

- We randomly sample 2 to 4 views per dish.
- We apply standard augmentations (jitter, crop, flip).
- We normalize ground-truth labels and use MSE for each regression head.

### D. Contribution Summary

- A **multi-view transformer-based system** for food nutrient estimation using RGB images only.
- A cross-view attention module that performs **implicit amodal completion** without 3D supervision.
- Direct regression of **macronutrient values** (calories, carbs, fat, protein, weight).
- A lightweight design that enables **real-time inference on mobile devices**.

## REFERENCES

[1] G. Ciocca, P. Napoletano, 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, May 2017.

[2] W. Min, S. Jiang, J. Sang, H. Wang, and L. Herranz, "Being a supercook: Joint food attributes and multimodal content modeling for recipe retrieval and exploration," IEEE Transactions on Multimedia, vol. 19, no. 5, pp. 1100–1113, May 2017.

[3] P. Pouladzadeh, S. Shirmohammadi, and R. Al-Maghrabi, "Measuring calorie and nutrition from food image," IEEE Transactions on Instrumentation and Measurement, vol. 63, no. 8, pp. 1947–1956, Aug. 2014.

[4] S. Lee and H. Kwon, "A two-stage approach for estimating food intake in mobile diet monitoring," in Proc. IEEE Int. Conf. on Image Processing (ICIP), 2015, pp. 2761–2765.

[5] Y. Shao, T. S. Huang, and Z. Yang, "Food volume estimation using 3D reconstruction from a single-view image," IEEE Transactions on Multimedia, vol. 21, no. 5, pp. 1190–1199, May 2019.

[6] Y. Han, M. Zhu, and S. Wang, "DPF-Nutrition: Depth prediction and fusion for food nutrition estimation," arXiv preprint arXiv:2310.11702, 2023.

[7] M. Thames, A. Spence, T. R. Hayes, and A. Farhadi, "Nutrition5k: Towards automatic nutritional understanding of generic food," in Proc. IEEE Conf. on Computer Vision and Pattern Recognition (CVPR), 2021, pp. 1911–1921.

[8] C. Lee, S. Lin, and H. Wang, "Multispectral imaging for real-time food recognition and caloric estimation," in Proc. ACM Int. Conf. on Multimedia (ACM MM), 2018, pp. 1204–1212.

[9] S. Fang, Y. Zhu, C. Boushey, and E. J. Delp, "Im2Calories: Towards an automated mobile vision food diary," in Proc. IEEE Int. Conf. on Computer Vision (ICCV), 2015, pp. 1233–1241.

[10] S. Myers et al., "Im2Calories: Towards an automated mobile vision food diary," in Proc. ICCV, 2015, pp. 1233–1241.

[11] M. Anthimopoulos, J. Dehais, P. Diem, and S. Mougiakakou, "Computer vision-based carbohydrate estimation for type 1 patients with diabetes using deep learning," Multimedia Tools and Applications, vol. 77, no. 9, pp. 10433–10451, May 2018.

[12] M. Mezgec and B. Koroušić Seljak, "NutriNet: A deep learning food and drink image recognition system for dietary assessment," Nutrients, vol. 9, no. 7, pp. 657, Jul. 2017.

[13] K. Bolliger, R. He, R. Kaehr, and J. Amft, "Deep learning based food image segmentation for dietary assessment," in Proc. Int. Conf. on Smart Health, 2019, pp. 1–14.

[14] A. Meyers, S. Johnston, V. Rathod, A. Korattikara, A. Gorban, N. Silberman, and S. Guadarrama, "Im2Calories: Towards an Automated Mobile Vision Food Diary," in Proc. ICCV, 2015.

[15] H. Kawano and K. Yanai, "Real-time mobile food recognition system," in Proc. IEEE Conf. on Computer Vision and Pattern Recognition Workshops (CVPRW), 2013, pp. 1–7.