# RecipeGen: A Benchmark for Real-World Recipe Image Generation

Ruoxuan Zhang, Hongxia Xie, Yi Yao, Jian-Yu Jiang-Lin, Bin Wen, Ling Lo, Hong-Han Shuai, Yung-Hui Li, Wen-Huang Cheng

*Abstract*—Recipe image generation is an important challenge in food computing, with applications from culinary education to interactive recipe platforms. However, there is currently no real-world dataset that comprehensively connects recipe goals, sequential steps, and corresponding images. To address this, we introduce RecipeGen, the first real-world goal-step-image benchmark for recipe generation, featuring diverse ingredients, varied recipe steps, multiple cooking styles, and a broad collection of food categories. Data is in https://github.com/zhangdaxia22/RecipeGen.

*Index Terms*—Recipe Dataset, Recipe Image Generation.

## I. INTRODUCTION

GENERATING illustrated instructions has gained growing attention as a means to assist users in comprehending and executing complex, step-by-step tasks. Whether assembling furniture, repairing devices, or following cooking procedures, visual depictions of each step can significantly reduce confusion and errors. One of the compelling applications in this space is recipe generation, where step-by-step illustrations enhance users' understanding of cooking processes. Existing studies in recipe image generation can be categorized into creating step images from the recipe title and ingredients [1], from recipe steps [2], [3], [4] and ingredients, or from a reference dish image [5], [6], [7], [8]. While these approaches have shown promise, they often struggle to generate intermediate steps with clarity and consistency, which leads to confusion and misinterpretation. A primary reason for this limitation lies in existing datasets [9], [10], [11], which typically only include final images of the dish, lacking detailed snapshots that capture essential transitions between steps.

To address these shortcomings, we introduce **RecipeGen**, a real-world goal-step-image benchmark tailored for generating illustrative step-by-step instructions in the cooking domain. We collect recipes spanning diverse cuisines, cooking methods, and food types, guided by 158 targeted keywords to ensure broad coverage. Each sample includes a dish name, ingredient list, detailed cooking steps, and corresponding images for every step. Altogether, our dataset comprises 21,944 recipes, totaling 139,872 images paired with textual descriptions. By retaining comprehensive visual details at each stage of food preparation, RecipeGen lays the groundwork for more accurate and instructive recipe-generation models, serving as a robust resource for the larger domain of generating illustrated instructions.

## II. RELATED WORK

Food is closely related to people's lives, and with the development of human society, diets have diversified. To facilitate the management of human life and health, the field of food computing [12] has emerged. Academically, topics such as food segmentation [13], [14], food recognition [15], [16], [17], food recommendation [18], [19], food reasoning [20] and recipe image generation have gained significant attention.
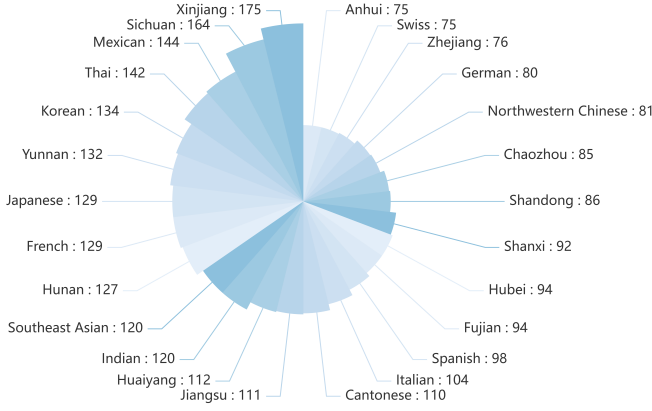
**Food Computing Dataset.** There are several existing datasets for food computing tasks. For instance, Food-101 [21], VireoFood-172 [22], Recipes242k [23]. However, these datasets are primarily used for fine-grained tasks such as food or ingredients classification, nutritional analysis, calorie calculation, etc. In addition to these information, FoodEarth [20] utilizes LLMs to develop a comprehensive dataset for food and nutrition analysis. Similarly, Recipe1M [9] provides recipe instructions linked to each final dish image. However, neither dataset includes images for each individual recipe step. While the recipe section of VGSI [24] offers step images for each instruction, many of these images are presented in a comic style. To address this gap, we propose the RecipeGen Benchmark, which is the first dataset that provides a real-world step-image recipe dataset.
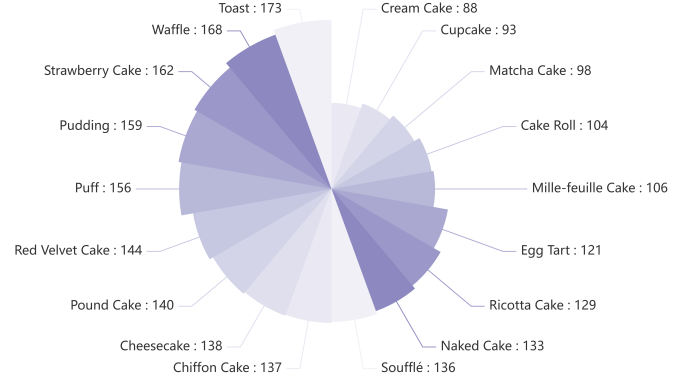
## III. RECIPEGEN BENCHMARK

The diversity of human culinary practices, developed over centuries, has led to a vast array of ingredient combinations and cooking methods, posing challenges for constructing comprehensive recipe step-image datasets. Existing datasets are limited, often constrained to narrow step counts and a restricted range of ingredients and interactions.

To address this gap, we propose **RecipeGen Benchmark (RGB)**, a goal-instruction image dataset designed to evaluate text-to-image models (T2Is) in understanding cooking instructions. RGB includes 21,944 recipes with 139,872 images and steps, encompassing diverse regional dishes and cooking styles. Collected from a large set of user-uploaded, real-world recipes[1], RGB provides a solid foundation for the evaluation of the T2I model. This dataset captures the distinctive ways different culinary traditions prepare similar ingredients and combines unique cooking techniques. For instance, Western recipes may pan-fry beef as a steak, while East Asian cuisine might braise it with root vegetables. Unlike existing datasets focused primarily on Western cuisines, such as VGSI-Recipe
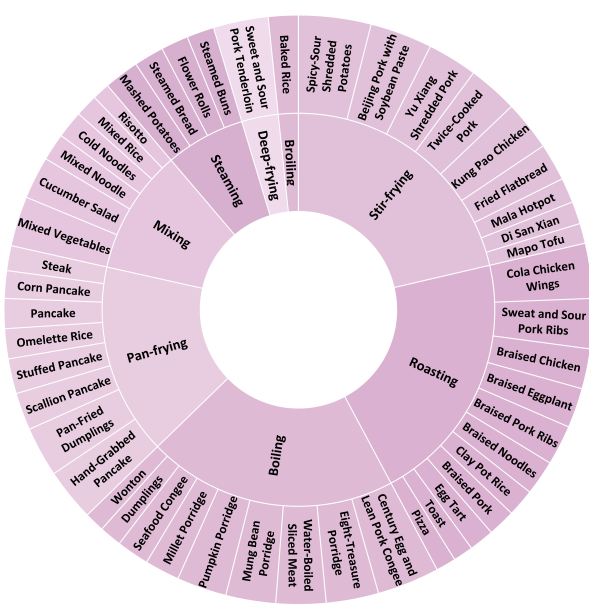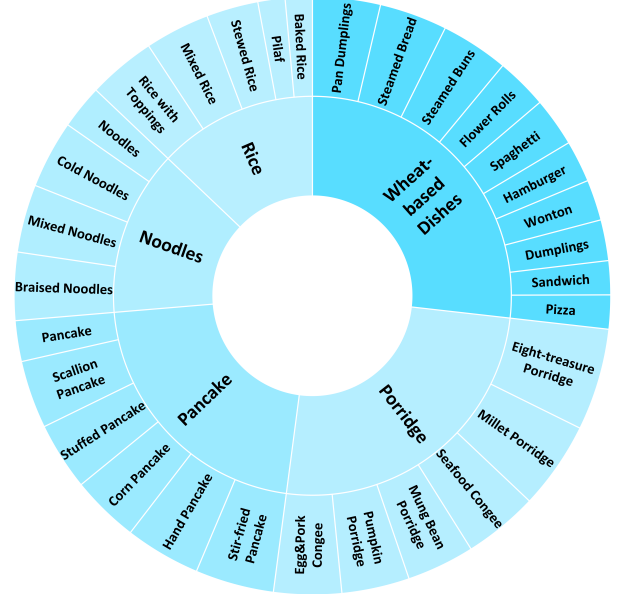
---

[1] www.douguo.com

(a) The distribution of region keywords.



(b) The distribution of dessert keywords.



(c) The distribution of cooking style keywords.



(d) The distribution of staple keywords.

Fig. 1: The distribution of region and dessert keywords in RecipeGen Benchmark.

| DataSet | Categories | Recipes | Steps/Images | Modality | Source | Keywords | | Data Faithfulness | | Step Counts | | |
|---|---|---|---|---|---|---|---|---|---|---|---|---|
| | | | | | | Chinese | Western | GF | SF | 0-6 | 7-10 | >10 |
| VGSI-Recipe | - | 1,157 | 6,417 | Image | wikiHow | ✗ | ✓ | 81.40 | 73.70 | 72.95% | 24.37% | 2.68% |
| YouCook2 | 89 | 2,000 | 15,433 | Video | YouTube | ✗ | ✓ | - | - | 38.50% | 44.15% | 17.35% |
| **RGB (Ours)** | **158** | **21,944** | **139,872** | **Image** | **douguo** | **56+64** | **37+64** | **85.40** | **84.44** | **44.16%** | **53.78%** | **2.06%** |

TABLE I: Comparison among Different Recipe Datasets. In the "keywords" column of the table, it indicates whether the keyword category is Chinese cuisine or Western cuisine. The number before the plus sign represents the total images of keywords that exclusively belong to either Chinese or Western cuisine, while the number after the plus sign represents the count of keywords that are common to both cuisines.

[25], RGB spans a wide range of regional styles, including East and Southeast Asian cuisines, enhancing model understanding of varied preparation techniques.

#### A. Dataset Construction Procedure

To ensure RGB's diversity, we curated 158 keywords across different cuisines (e.g., Cantonese, Mexican), dish types (e.g., main courses, desserts), and cooking techniques (e.g., stir-frying, roasting), gathering 29,026 recipes.

**Quality Control.** The process is shown in Fig. 4. To maintain consistency and clarity, we implement a quality control process to address common issues like incomplete instructions or overly detailed steps. This process involved: (1) Removing recipes with low quality (including lack images or steps, mismatched numbers of images and steps, vague directions [2]). In this step, we delete 4,978 recipes; (2) Using GPT-4o [26]

---

[2]For example, some steps in the recipe only contain phrases like "see as images" without providing clear or actionable instructions.

Fig. 2: Some examples in RecipeGen. Each sample contains its goal, steps, and the images corresponding to each step.

to clean and refine recipes by eliminating irrelevant content [3], merging redundant steps, and summarizing for coherence. And then, we remove recipes where the LLM's hallucination issues caused a single step to be split into two, as well as recipes with incorrect output formats. A total of 2,104 recipes were deleted. After Quality Control, the number of recipe is 21,944 and the number of images is 139,872; (3) Computing data faithfulness and using human check to insure the quality. The result of metrics is in Tab. I.

**Prompts for GPT-4o.** To ensure Quality Control, we employ GPT-4o to combine adjacent simple steps, generate captions to prevent over-merging by GPT-4o and translate the output into English. We establish clear principles for GPT-4o and outline a step-by-step process to guide it in producing accurate and appropriate results. The prompt is in Fig. 3.

**Human Check.** In addition to using the step faithfulness and goal faithfulness metrics to evaluate the quality of the dataset, we also conduct a human review. Specifically, we ask annotators to verify whether the recipe matched its dish name, whether the steps corresponded to each step image, and whether the images within the same recipe showed continuity. We randomly select 50 recipes for this review. Among these samples, we found one instance where the dish name was mislabeled as "dumplings" instead of "steamed buns," and another instance where all the steps lacked specific ingredients. Instead, the textual steps were generic, such as "as shown in the image" or "ready to serve".

### B. Keywords used in RecipeGen Benchmark Collection Process

Before gathering recipe data, we first analyze the key characteristics of dishes, based on the idea that a dish is defined by its primary ingredients, cooking techniques, and preparation methods. To ensure a comprehensive collection of recipes, we select a set of 158 keywords encompassing cuisines from various regions, a wide range of cooking methods (e.g., mixing, frying, stir-frying), and categories such as staples, desserts, and ingredient-specific recipes. Additionally, we scrape 996 unique new dishes created by cooking enthusiasts from "jingxuan" section of the website to test the robustness of the model. Fig. 1 shows some distribution of keywords and the full list is provided in Tab. II. Fig. 2 shows some recipe examples from RecipeGen.

### C. Dataset Statistics

Tab. I presents statistical details of our RGB in comparison to existing benchmarks. As shown in Tab. I, existing datasets for recipe image generation are limited in the amount of data available, rendering them insufficient for training T2I models. Moreover, the **VGSI-Recipe** [25] [4] dataset is predominantly Western-focused, offering limited diversity in terms of cooking regions and step counts. Additionally, the dataset includes comic-style images, which diminishes its effectiveness for training models aimed at generating realistic images. On the

---

[3]For instance, some users include exclamatory phrases such as "This tastes amazing!" or suggestions like "Using ingredient A instead of B works better."

[4]VGSI does not categorize the recipe data, so we search the recipe data in its JSON structure under the "method" key using the keyword "cook".

**Quality Control Prompt For GPT-4o**

```
You are now a strict expert in formatting and outputting recipes. Please read the following recipe and check whether
any steps need to be merged.
The rules are as follows:
1.Please merge adjacent simple recipe steps.
2.Provide the recipe in both Chinese and its English translation. Even if there are no modifications or merges,
output both the Chinese recipe and the English version.
3.Generate a summary phrase for each step in the recipe.
4.During the preparation phase, especially at the beginning, if two steps involve handling different ingredients,
keep these steps separate and do not merge them. Ensure the preparation steps for different ingredients are
independent.
5.Avoid merging wherever possible. If merging is necessary, it should only involve adjacent steps and not skip over
or reorder them. The relative sequence of the original steps must remain unchanged.
6.Do not split steps. Do not divide one original step into two or more new steps. The number of steps in the output
should not exceed the number of steps in the original recipe.
7.For the final steps of the recipe, if they are unrelated to the cooking process, you may omit them. Also,
exclamatory sentences within the steps can be removed.
8.Note: Output the recipe in both Chinese and English, with Chinese first, followed by the English translation. The
English version should be clear, easy to understand, and correspond to the Chinese steps one-to-one. Do not split
steps or create new steps from a single original step. Only merge steps when necessary, following the rules.
9.The output must follow this format:
"Summary of the step: Content of the step – Original step start number – Original step end number"
Step-by-Step Process:
1. Merge steps and generate captions based on the above principles.
2. Add the corresponding original recipe's start and end step numbers at the end of each new step.
3. Translate into English.
4. Maintain continuity between steps.
```

Fig. 3: Quality Control Prompt for GPT-4o. We utilize GPT-4o to merge steps and generate captions.

| Type | Keywords |
|---|---|
| Regions | Mexican cuisine, Japanese cuisine, Korean cuisine, Indian cuisine, Swiss cuisine, Yunnan cuisine, Su cuisine, Anhui cuisine, Sichuan cuisine, Northeast Chinese cuisine, Zhejiang cuisine, Hubei cuisine, Shandong cuisine, Fujian cuisine, Northwest Chinese cuisine, Shanxi cuisine, Hunan cuisine, Cantonese cuisine, Huaiyang cuisine, Chaozhou cuisine, Xinjiang cuisine, French cuisine, Italian cuisine, German cuisine, Southeast Asian cuisine, Spanish cuisine, Thai cuisine, Chinese cuisine. |
| Dishes | Kung Pao Chicken, Di San Xian, Mapo Tofu, Claypot Rice, Omurice, Boiled Pork Slices, Mashed Potatoes, Braised Pork Belly, Sweet and Sour Pork Tenderloin, Braised Eggplant, Twice-Cooked Pork, Cold Cucumber Salad, Sweet and Sour Pork Ribs, Shredded Pork in Garlic Sauce, Braised Pork Ribs, Spicy Hot Pot, Yellow Braised Chicken, Spicy and Sour Shredded Potatoes, Shredded Pork with Peking Sauce, Steak, Cola Chicken Wings. |
| Staples | Noodles, Cold Noodles, Mixed Noodles, Braised Noodles, Steamed Bun, Dumplings, Steamed Bread, Wonton, Flower Rolls, Pan-Fried Dumplings, Scallion Pancake Roll, Pancake, Scallion Pancake, Stuffed Pancake, Corn Pancake, Stir-Fried Pancakes, Century Egg and Pork Congee, Millet Congee, Seafood Congee, Pumpkin Congee, Mung Bean Congee, Eight Treasure Congee, Rice with Toppings, Mixed Rice, Baked Rice, Braised Rice, Pilaf, Hamburger, Spaghetti, Pizza, Sandwich, Fried Rice, Braised Rice. |
| Ingredients | Cauliflower, Water Spinach, Asparagus, Lettuce Stem, Potato, Celery, Baby Bok Choy, Fennel, Chinese Toon, Rapeseed Greens, Celtuce, Napa Cabbage, Chives, Broccoli, Mustard Greens, Chinese Cabbage, Lettuce, Eggplant, Seaweed, Yellow Chives, Chinese Green Cabbage, Ice Plant, Garlic Sprouts, Vegetable Shoots, Spring Bamboo Shoots, Broccoli, Crown Daisy, Tomato, Spinach, Bitter Chrysanthemum, Cabbage, Enoki Mushroom, King Oyster Mushroom, Shiitake Mushroom, Wood Ear Mushroom, White Shimeji Mushroom, Tremella, Chinese Yam, Carrot, Water Bamboo, White Radish, Sweet Potato, Taro, Lotus Root, Winter Melon, Cucumber, Pumpkin, Loofah, Squash. |
| Dessert | Toast, Egg Tart, Pudding, Puff, Waffle, Cheesecake, Cream Cake, Layer Cake, Cupcake, Chiffon Cake, Cake Roll, Soufflé, Cream Cheese Cheesecake, Naked Cake, Matcha Cake, Red Velvet Cake, Pound Cake, Strawberry Cake. |
| Others | Breakfast, Cold Dishes, Lunch, Hot Dishes, Midnight Snack, Home-cooked Dishes, Dinner, Side Dishes, Jingxuan. |

TABLE II: Keywords used in Data Collection. We categorize these 158 keywords into six main groups, including regions and countries, various ingredients, classic dishes, staple foods, and desserts.

other hand, while the **YouCook2** [24] dataset features a diverse range of cooking styles, it is composed of video data. The data format often obscures essential visual cues due to the position of the chef or suboptimal camera angles, compromising instructional clarity when used for training image-based recipe models. In contrast, RGB offers 21,944 recipes with 139,872 image-step pairs, providing a substantial variety of detailed instructions. Additionally, the images focus solely on the food itself, as they are contributed by amateur food enthusiasts rather than being professionally staged to capture the chef. Tab. I further illustrates that most RGB recipes contain over six steps, enhancing the instructional richness of the dataset.

Furthermore, to validate quality of the datasets, we assessed goal faithfulness (GF) and step faithfulness (SF) to evaluate the alignment between recipe steps and intended outcomes. GF computes the CLIP score between the caption of the final step and the final image of the recipe. SF computes the CLIP score between each step and its respective image. As indicated in Tab. I, RGB achieves significantly higher GF and SF scores than VGSI-Recipe, highlighting its superior adherence to recipe objectives and procedural accuracy.

Overall, our proposed RecipeGen Benchmark (RGB) has four notable features:

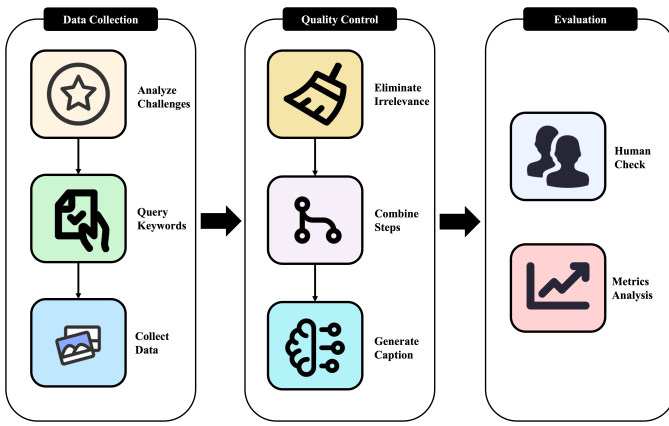- **Broad Step Distribution**: RGB includes recipes with 2

Fig. 4: Dataset Construction Procedure. We first analyze the characteristics of the dishes and select 158 keywords. Subsequently, we utilize GPT-4o to perform quality control by omitting irrelevant steps, merging adjacent simple actions, and generating captions. Finally, we calculate metrics and conduct human checks to ensure the usability of the dataset.

to 25 steps, with an average of 6.4 steps and 55.84% of recipes containing more than 6 steps, supporting the modeling of long-range semantic relationships.

- **Ingredient Diversity**: RGB provides an average of 9 ingredients per recipe (including seasonings), capturing a rich variety of ingredients and interactions.
- **Variety of Cooking Styles**: The dataset encompasses a wide range of cooking styles enhanced by numerous unique keywords, making it versatile across different culinary processes.
- **Real-World Representativeness**: Collected from various users, RGB consists entirely of real-world recipes and closely resembles the types of instructions that users might upload in practical scenarios. Quality control using GPT-4o ensures that steps are trustworthy and concise, reflecting actual cooking practices.

## IV. CONCLUSION

In this work, we introduced RecipeGen, the first real-world goal-step-image benchmark for recipe image generation, addressing the lack of comprehensive datasets in food computing. We believe that our benchmark can foster further advancements in food computing, particularly in real-world, interactive culinary applications.

## REFERENCES

[1] H. H. Lee, K. Shu, P. Achananuparp, P. K. Prasetyo, Y. Liu, E.-P. Lim, and L. R. Varshney, "Recipegpt: Generative pre-training based cooking recipe generation and evaluation system," in *Companion Proceedings of the Web Conference 2020*, 2020, pp. 181–184.

[2] S. Pan, L. Dai, X. Hou, H. Li, and B. Sheng, "Chefgan: Food image generation from recipes," in *Proceedings of the 28th ACM International Conference on Multimedia*, 2020, pp. 4244–4252.

[3] Z. Liu, K. Niu, and Z. He, "Ml-cookgan: Multi-label generative adversarial network for food image generation," *ACM Transactions on Multimedia Computing, Communications and Applications*, vol. 19, no. 2s, pp. 1–21, 2023.

[4] F. Han, R. Guerrero, and V. Pavlovic, "Cookgan: Meal image synthesis from ingredients," in *Proceedings of the IEEE/CVF Winter Conference on Applications of Computer Vision*, 2020, pp. 1450–1458.

[5] P. Chhikara, D. Chaurasia, Y. Jiang, O. Masur, and F. Ilievski, "Fire: Food image to recipe generation," in *Proceedings of the IEEE/CVF Winter Conference on Applications of Computer Vision*, 2024, pp. 8184–8194.

[6] A. Dosovitskiy, "An image is worth 16x16 words: Transformers for image recognition at scale," *arXiv preprint arXiv:2010.11929*, 2020.

[7] H. Wang, G. Lin, S. C. Hoi, and C. Miao, "Structure-aware generation network for recipe generation from images," in *Proceedings of the European Conference on Computer Vision*. Springer, 2020, pp. 359–374.

[8] ——, "Learning structural representations for recipe generation and food retrieval," *IEEE Transactions on Pattern Analysis and Machine Intelligence*, vol. 45, no. 3, pp. 3363–3377, 2022.

[9] J. Marın, A. Biswas, F. Ofli, N. Hynes, A. Salvador, Y. Aytar, I. Weber, and A. Torralba, "Recipe1m+: A dataset for learning cross-modal embeddings for cooking recipes and food images," *IEEE Transactions on Pattern Analysis and Machine Intelligence*, vol. 43, no. 1, pp. 187–203, 2021.

[10] M. Bień, M. Gilski, M. Maciejewska, W. Taisner, D. Wisniewski, and A. Lawrynowicz, "RecipeNLG: A cooking recipes dataset for semi-structured text generation," in *Proceedings of the 13th International Conference on Natural Language Generation*, B. Davis, Y. Graham, J. Kelleher, and Y. Sripada, Eds. Dublin, Ireland: Association for Computational Linguistics, Dec. 2020, pp. 22–28. [Online]. Available: https://aclanthology.org/2020.inlg-1.4

[11] D. Batra, N. Diwan, U. Upadhyay, J. S. Kalra, T. Sharma, A. K. Sharma, D. Khanna, J. S. Marwah, S. Kalathil, N. Singh *et al.*, "Recipedb: a resource for exploring recipes," *Database*, vol. 2020, p. baaa077, 2020.

[12] W. Min, S. Jiang, L. Liu, Y. Rui, and R. Jain, "A survey on food computing," *ACM Computing Surveys (CSUR)*, vol. 52, no. 5, pp. 1–36, 2019.

[13] X. Lan, J. Lyu, H. Jiang, K. Dong, Z. Niu, Y. Zhang, and J. Xue, "Foodsam: Any food segmentation," *IEEE Transactions on Multimedia*, 2023.

[14] Y. Yin, H. Qi, B. Zhu, J. Chen, Y.-G. Jiang, and C.-W. Ngo, "Foodlmm: A versatile food assistant using large multi-modal model," *arXiv preprint arXiv:2312.14991*, 2023.

[15] W. Min, Z. Wang, Y. Liu, M. Luo, L. Kang, X. Wei, X. Wei, and S. Jiang, "Large scale visual food recognition," *IEEE Transactions on Pattern Analysis and Machine Intelligence*, vol. 45, no. 8, pp. 9932–9949, 2023.

[16] S. Jiang, W. Min, L. Liu, and Z. Luo, "Multi-scale multi-view deep feature aggregation for food recognition," *IEEE Transactions on Image Processing*, vol. 29, pp. 265–276, 2020.

[17] Y. Zhang, L. Deng, H. Zhu, W. Wang, Z. Ren, Q. Zhou, S. Lu, S. Sun, Z. Zhu, J. M. Gorriz *et al.*, "Deep learning in food category recognition," *Information Fusion*, vol. 98, p. 101859, 2023.

[18] W. Min, S. Jiang, and R. Jain, "Food recommendation: Framework, existing solutions, and challenges," *IEEE Transactions on Multimedia*, vol. 22, no. 10, pp. 2659–2671, 2019.

[19] W. Wang, L.-Y. Duan, H. Jiang, P. Jing, X. Song, and L. Nie, "Market2dish: health-aware food recommendation," *ACM Transactions on Multimedia Computing, Communications, and Applications (TOMM)*, vol. 17, no. 1, pp. 1–19, 2021.

[20] P. Zhou, W. Min, C. Fu, Y. Jin, M. Huang, X. Li, S. Mei, and S. Jiang, "Foodsky: A food-oriented large language model that passes the chef and dietetic examination," *arXiv preprint arXiv:2406.10261*, 2024.

[21] L. Bossard, M. Guillaumin, and L. Van Gool, "Food-101–mining discriminative components with random forests," in *Proceedings of the European Conference on Computer Vision*. Springer, 2014, pp. 446–461.

[22] J. Chen and C.-W. Ngo, "Deep-based ingredient recognition for cooking recipe retrieval," in *Proceedings of the 24th ACM international conference on Multimedia*, 2016, pp. 32–41.

[23] M. Rokicki, C. Trattner, and E. Herder, "The impact of recipe features, social cues and demographics on estimating the healthiness of online recipes," in *Proceedings of the international AAAI conference on web and social media*, vol. 12, no. 1, 2018.

[24] L. Zhou, C. Xu, and J. Corso, "Towards automatic learning of procedures from web instructional videos," in *Proceedings of the AAAI Conference on Artificial Intelligence*, vol. 32, no. 1, 2018.

[25] Y. Yang, A. Panagopoulou, Q. Lyu, L. Zhang, M. Yatskar, and C. Callison-Burch, "Visual goal-step inference using wikihow," *arXiv preprint arXiv:2104.05845*, 2021.

[26] J. Achiam, S. Adler, S. Agarwal, L. Ahmad, I. Akkaya, F. L. Aleman, D. Almeida, J. Altenschmidt, S. Altman, S. Anadkat *et al.*, "Gpt-4 technical report," *arXiv preprint arXiv:2303.08774*, 2023.