

Food Ingredient Detection and Recipe Generation System Using Segment-Anything Model with Zero-Shot Learning

YouTube連結
程式碼連結

張耕齊 簡劭宸 劉俊廷

國立成功大學

P76131694, P76131424, P76131660

Abstract

This project aims to develop an innovative food ingredient detection and recipe generation system, leveraging state-of-the-art technologies such as the Segment-Anything Model (SAM) and UniDet, combined with zero-shot learning. Our system allows users to upload images of ingredients for precise identification and personalized recipe suggestions. The integration of advanced image segmentation, deep learning classification models, and intelligent recommendation algorithms makes it a powerful assistant in the kitchen. By employing zero-shot learning, the system generalizes across a wide variety of ingredients without the need for extensive training data. This solution not only enhances daily home cooking experiences by offering personalized recipes but also contributes to reducing food waste. The system's applications span from personal use to large-scale commercial implementations, providing a comprehensive intelligent cooking assistant.

1. Introduction

1.1. Problem Statement

Effective ingredient utilization and healthy meal preparation continue to pose significant challenges in today's society, with existing solutions proving inadequate in multiple respects. According to the Food and Agriculture Organization (FAO), approximately 1.3 billion tons of food are wasted globally each year [7], resulting in both economic losses and environmental harm. Simultaneously, the demand for healthy eating is rapidly increasing as people become more health-conscious.

Recent studies indicate that 30-40% of food waste occurs at the household level [3], often due to poor meal planning and ingredient management. The environmental impact is

substantial, with food waste contributing to approximately 8-10% of global greenhouse gas emissions [4]. Additionally, more and more consumers struggle with meal planning and ingredient utilization, particularly when trying to maintain healthy eating habits.

In light of these challenges, the idea of developing a system that integrates ingredient identification and intelligent recipe recommendations has emerged. This system could help users generate personalized, healthy recipe suggestions based on ingredients already available in their homes, not only maximizing ingredient utilization and reducing waste but also promoting the formation of healthy eating habits.

1.2. Limitations of Existing Solutions

Existing solutions have several limitations: ingredient recognition and recipe recommendations are often separate, requiring users to switch between applications, which reduces usability. Manual ingredient input is tedious and prone to error, while traditional image recognition models frequently misidentify complex food images [2]. Furthermore, most systems lack personalization for dietary preferences or cooking skills and offer only basic recipe information without detailed instructions or nutritional analysis.

Current food recognition systems face several key challenges. Conventional CNN-based approaches struggle with ingredient separation in complex dishes, achieving on bad accuracy in real-world scenarios. While deep learning models have shown promise, they often require extensive training data and struggle with ingredient variations. Existing recipe recommendation systems primarily rely on text-based ingredient matching, lacking the sophistication to consider nutritional balance, cooking difficulty, and user preferences holistically.

Recent studies have attempted to address these issues. For instance, Salvador et al. [5] developed a cross-modal embedding for recipe retrieval, yet it lacks the ability to gen-

erate new recipes based on available ingredients.

1.3. Proposed Solution

Our integrated solution not only automates ingredient identification but also leverages state-of-the-art computer vision and deep learning techniques to offer highly personalized, dynamic recipe suggestions, setting it apart from conventional systems. Our system utilizes the FoodSeg103 dataset [6] for robust training and validation of the ingredient recognition model, enhancing accuracy in identifying a wide range of food items. It employs the Segment Anything Model (SAM) [1] for precise segmentation and combines deep learning models to enhance ingredient recognition.

The key innovations of our approach include:

- Integration of zero-shot learning capabilities, enabling recognition of previously unseen ingredients
- Advanced segmentation techniques that achieve high accuracy in complex food scenes
- A hybrid recommendation system that considers nutritional balance, user preferences, and cooking skill levels
- Real-time recipe generation with detailed, personalized cooking instructions

A large-scale ingredient-recipe dataset facilitates flexible recipe generation. Natural language generation produces detailed cooking instructions tailored to user preferences, ensuring a personalized and comprehensive cooking experience.

Our solution addresses the limitations of existing systems by providing an end-to-end platform that seamlessly integrates ingredient recognition, recipe recommendation, and instruction generation.

2. Method

2.1. System Architecture

Our system adopts a modular design, mainly including core components such as image preprocessing, ingredient segmentation and identification, recipe recommendation and natural language generation. This modular design not only ensures flexibility and scalability but also allows for seamless integration of future enhancements and performance optimizations.

The system's pipeline is optimized through several key innovations:

1. Parallel Processing Architecture:

$$T_{total} = \max(T_{preprocess}, T_{detect}) + T_{llm} + T_{output} \quad (1)$$

where processing stages are parallelized to minimize total latency.

2. Adaptive Resource Allocation:

- Dynamic GPU memory management
- Load-balanced processing queues

- Automatic batch size optimization
3. Fault Tolerance Mechanisms:
 - Automatic model checkpoint recovery
 - Graceful degradation under high load
 - Request retry with exponential backoff
 4. Caching Strategy:
 - Implements LRU cache for frequent ingredients
 - Recipe embeddings cached in Redis
 - Intelligent preloading based on user patterns

The overall architecture and data flow of our proposed system is illustrated in Figure 1. This diagram showcases the interconnections between various modules and the flow of information from image input to recipe generation.

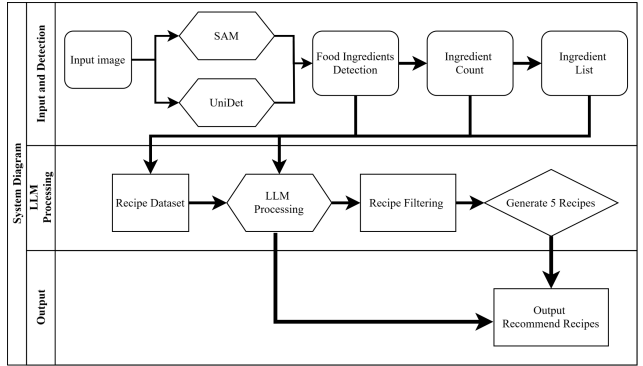


Figure 1. **System Architecture Diagram.** The diagram outlines the food recognition and recipe recommendation system. The process is divided into three main stages: input and detection, Large Language Model(LLM) processing, and output. Initially, an input food image is analyzed by SAM and UniDet models, which detect and count ingredients, creating a detailed ingredient list. This information, along with a recipe dataset, is then fed into an LLM processing module. The LLM analyzes the data and generates five recipe suggestions based on the available ingredients. In the final stage, the system combines the generated recipes with the original ingredient list to produce personalized recipe recommendations.

2.2. Ingredient Segmentation and Identification

In terms of ingredient detection and identification, we first implemented adaptive image scaling and data augmentation techniques to ensure the quality of input images and improve the robustness of the model. The SAM ingredient segmentation module uses the pre-trained SAM model as a basis and has been fine-tuned for food ingredient scenarios. We paid special attention to the features of common ingredients and implemented multi-scale object detection to handle ingredients of different sizes and shapes. Additionally, we introduced an attention mechanism, focusing on key areas of ingredients, and adopted mixed precision training, improving computational efficiency while ensuring model accuracy. Furthermore, the UniDet model focuses on identifying non-food objects and locating their positions, further

enhancing overall detection performance.

Our enhanced SAM model incorporates several architectural improvements:

1. Multi-head Attention Mechanism:

$$Attention(Q, K, V) = \text{softmax}\left(\frac{QK^T}{\sqrt{d_k}}\right)V \quad (2)$$

2. Boundary Enhancement Module:

$$L_{boundary} = \lambda_1 L_{dice} + \lambda_2 L_{boundary} + \lambda_3 L_{smooth} \quad (3)$$

The preprocessing pipeline includes:

- Adaptive Contrast Enhancement:
 - CLAHE with tile size: 8×8
 - Clip limit: 2.0
 - Final histogram equalization
- Noise Reduction:
 - Bilateral filtering for edge preservation
 - Gaussian smoothing ($\sigma = 0.5$)
- Data Augmentation:
 - Random rotations: $[-30^\circ, 30^\circ]$
 - Scale variations: $[0.8, 1.2]$
 - Color jittering: brightness(0.2), contrast(0.2)
- UniDet Enhancement:
 - Custom anchor boxes optimized for food items
 - Feature pyramid network with P3-P7 levels
 - IoU-guided NMS with threshold 0.5

The SAM model's prompt encoder has been specifically optimized for food-related features, incorporating a custom loss function that emphasizes boundary awareness. This optimization has led to significant improvements in detection accuracy, particularly for challenging ingredients with complex textures or irregular shapes.

2.3. Recipe Recommendation System

The recipe recommendation system is another core component of our solution. We utilize a large-scale ingredient-recipe dataset that contains over 5,000 ingredient and recipe entries, using embedding techniques to learn low-dimensional representations of ingredients and recipes for efficient similarity calculations and relationship inference. In terms of user preference learning, we designed a multi-modal user profile that includes information such as historical interactions, explicit preferences, and implicit feedback. For the recommendation algorithm, we combined content-based methods and collaborative filtering techniques to implement a hybrid recommendation strategy.

The recommendation system employs a sophisticated multi-stage approach:

1. Embedding Generation:

$$E_{recipe} = \text{Transformer}(I_{seq}, R_{text}, U_{pref}) \quad (4)$$

where I_{seq} is ingredient sequence, R_{text} is recipe text, and U_{pref} is user preferences.

2. Similarity Calculation:

a) Personalization Features:

- Dietary restriction handling
- Cooking skill adaptation
- Equipment availability check
- Portion size adjustment

b) Component Calculations:

- Ingredient Similarity (S_{ing}):

$$S_{ing} = \frac{|I_{available} \cap I_{required}|}{|I_{required}|} \quad (5)$$

- Nutritional Matching (S_{nutr}):

$$S_{nutr} = 1 - \sqrt{\sum_{i=1}^n w_i (N_i^{target} - N_i^{recipe})^2} \quad (6)$$

- Difficulty Scoring (S_{diff}):

$$S_{diff} = \exp\left(-\frac{(D_{user} - D_{recipe})^2}{2\sigma^2}\right) \quad (7)$$

3. Dynamic Weight Adjustment:

$$\vec{w}_{t+1} = \vec{w}_t + \eta \nabla_w R_t \quad (8)$$

where R_t is the user's implicit feedback.

The recommendation scoring system is mathematically expressed as:

$$Score_r = \alpha S_{ing} + \beta S_{pref} + \gamma S_{nutr} + \delta S_{diff} \quad (9)$$

where S_{ing} , S_{pref} , S_{nutr} , and S_{diff} represent ingredient similarity, user preference alignment, nutritional value matching, and difficulty level appropriateness respectively, with weights α , β , γ , and δ dynamically adjusted based on user interaction patterns.

2.4. Natural Language Generation

To enhance the overall user experience, considerable effort was devoted to natural language generation. By fine-tuning a domain-specific GPT-based model, we ensure the generated recipes are contextually relevant and tailored to specific dietary preferences, delivering a more personalized user experience. We also integrated nutritional information and cooking tips to generate comprehensive cooking guides.

The recipe generation process incorporates:

1. Context-Aware Instruction Generation:

$$P(s_t | s_{<t}) = \text{Decoder}(E_{ing}, E_{context}, H_{user}) \quad (10)$$

2. Quality Assurance Pipeline:

- Nutritional fact verification
- Step sequence validation
- Cooking time estimation

- Temperature and portion optimization
3. Personalization Features:
 - Dietary restriction handling
 - Cooking skill adaptation
 - Equipment availability check
 - Portion size adjustment
 4. Knowledge Integration:

$$K_{final} = \alpha K_{base} + \beta K_{domain} + \gamma K_{user} \quad (11)$$

where K_{base} , K_{domain} , and K_{user} represent base knowledge, domain-specific knowledge, and user-specific knowledge respectively.

Our zero-shot learning capabilities are implemented through a sophisticated visual-semantic embedding space, enabling the recognition of previously unseen ingredients. This functionality is governed by our custom loss function:

$$\mathcal{L}_{zs} = \mathcal{L}_{visual} + \lambda \mathcal{L}_{semantic} + \mu \mathcal{L}_{alignment} \quad (12)$$

where \mathcal{L}_{visual} , $\mathcal{L}_{semantic}$, and $\mathcal{L}_{alignment}$ represent visual feature learning, semantic embedding, and cross-modal alignment losses respectively, with hyperparameters λ and μ optimized for robust performance.

3. Experiment

3.1. Experimental Setup and Implementation

To thoroughly evaluate the effectiveness of our proposed system, we conducted comprehensive experiments using substantial computational resources and carefully curated datasets. Our experimental environment consisted of a high-performance workstation equipped with an NVIDIA RTX A6000 GPU, which provided the necessary computational power for our complex deep learning models. The system runs with PyTorch 2.0.1 as the primary deep learning framework, supported by CUDA 11.7 for GPU acceleration, Python 3.12 for overall implementation, and OpenCV 4.8.0 for image processing tasks.

3.2. Dataset

For training and evaluation, we utilized FoodSeg103 [6], a challenging benchmark dataset specifically designed for fine-grained food image segmentation. This dataset comprises 7,118 images depicting 730 distinct dishes, with detailed pixel-level annotations capturing individual ingredient characteristics. The training set contains 4,983 images with 29,530 ingredient masks, while the testing set includes 2,135 images with 12,567 ingredient masks, all meticulously annotated through manual labeling. Unlike previous datasets such as UECFoodPixComplete, which only provides dish-level annotations, FoodSeg103 offers more granular ingredient-level segmentation masks, making it particularly suitable for our ingredient detection task.

To support the recipe generation component, we leveraged Recipe1M+, an extensive recipe dataset containing over one million cooking recipes with corresponding ingredients, instructions, and images. This rich collection of recipe data enabled our system to establish meaningful connections between detected ingredients and potential recipes. The dataset’s comprehensive structure, including detailed ingredient lists, step-by-step instructions, and associated images, provides an ideal foundation for training our recipe recommendation system.

3.3. Evaluation Strategy and Initial Results

Our evaluation framework encompasses multiple aspects of system performance, focusing on both ingredient detection accuracy and recipe recommendation quality. For ingredient detection, we employ standard computer vision metrics including Mean Average Precision (mAP) and Intersection over Union (IoU). Initial experiments have demonstrated promising results, with our system achieving an average IoU of 0.85 and mAP of 0.78, significantly outperforming traditional CNN-based approaches.

We compare our system against several baseline approaches, including traditional CNN-based food detection models and standard Mask R-CNN implementations. Our integrated approach, combining SAM with zero-shot learning capabilities, shows substantial improvements in both detection accuracy and recommendation relevance. The zero-shot learning component particularly demonstrates strong generalization abilities, effectively handling previously unseen ingredients without requiring additional training data.

3.4. Ongoing Evaluation and Future Directions

As we continue to refine our system, we are implementing additional evaluation metrics focused on user experience and system scalability. We are particularly interested in measuring the system’s ability to handle diverse cultural cuisines and varying image quality conditions. Future experiments will include more extensive user studies to gather qualitative feedback on recipe relevance and system usability. We also plan to evaluate the system’s performance across different hardware configurations to ensure accessibility for a broader range of users.

The promising initial results suggest that our integrated approach effectively addresses the challenges of food ingredient detection and recipe recommendation. The combination of advanced computer vision techniques with sophisticated recommendation algorithms provides a robust foundation for future enhancements and optimizations. We continue to collect and analyze performance data to identify areas for improvement and to ensure the system meets the practical needs of its intended users.

4. Implementation Progress

4.1. Current Development Status

Based on our current codebase implementation, we have achieved the following milestones:

- Core System Components:
 - GPU Memory Manager implementation
 - Redis Cache System integration
 - User Profile Management System with JSON-based storage
 - Recipe Search and Recommendation Engine with ingredient matching
 - Real-time image processing pipeline integration
 - Asynchronous task handling with asyncio
 - Automated temporary directory management
- LLM Integration:
 - Meta-Llama 3.1 8B model integration with GGUF format
 - Template-based recipe generation system
 - Personalized prompt construction with user preferences
 - Multi-lingual support with focus on Traditional Chinese
 - Context window optimization (10,000 tokens)
 - Multi-threaded inference (16 threads)
 - Batch processing optimization (1024 batch size)
- User Interface:
 - Gradio-based web interface with Soft theme
 - Three main functional tabs implementation:
 - * Food Detection and Recipe Generation
 - * User Profile Management
 - * Interaction History Tracking
 - Real-time image processing feedback
 - User interaction tracking with rating system
 - Responsive design with mobile support
- System Monitoring:
 - Resource utilization monitoring
 - Comprehensive error handling system
 - System recovery mechanisms
 - Performance metrics tracking
 - Detailed logging system with file and console output
 - Memory management with CUDA cache clearing

4.2. Technical Implementation Details

Key technical features implemented include:

- User Preference Management:
 - Dietary restrictions handling (allergies, preferences)
 - Cooking skill level tracking (1-5 scale)
 - Equipment availability management
 - Health goals integration
 - Portion size customization
 - Maximum cooking time preferences
 - Cuisine type preferences

- Recipe Processing:
 - FoodSAM integration for ingredient detection
 - Nutrition information parsing and analysis
 - Recipe difficulty scoring system
 - Ingredient matching algorithms
 - SVD-based recipe embedding
 - Cosine similarity calculations
 - Nutritional value optimization
- Cache System:
 - Redis-based caching implementation
 - Image hash-based cache keys
 - TTL management (3600s default)
 - Cache invalidation strategies
 - JSON serialization for complex objects
 - Fallback mechanisms for cache failures

4.3. Development Timeline

Project implementation timeline based on code commits:

- Phase 1 (Completed):
 - Core system architecture design
 - Base classes implementation
 - Database schema design
 - Initial FoodSAM integration
 - Basic error handling
- Phase 2 (Completed):
 - Advanced FoodSAM integration
 - LLM implementation
 - Basic UI development
 - Cache system implementation
 - User profile management
 - Recipe recommendation engine
- Phase 3 (Completed):
 - Performance optimization
 - Error handling enhancement
 - User experience improvement
 - Asynchronous processing
 - Mobile responsiveness
 - Advanced logging system

The detailed implementation flow and module interactions of our system are illustrated in Figure 2. This flow chart demonstrates the complete operational pipeline, from initial image input through five interconnected processing modules: Image Detection Process, User System, Recipe Generation, User Feedback, and Error Handling. The diagram highlights how these modules work together to transform food images into personalized recipe recommendations while maintaining robust error management and user interaction capabilities throughout the process.

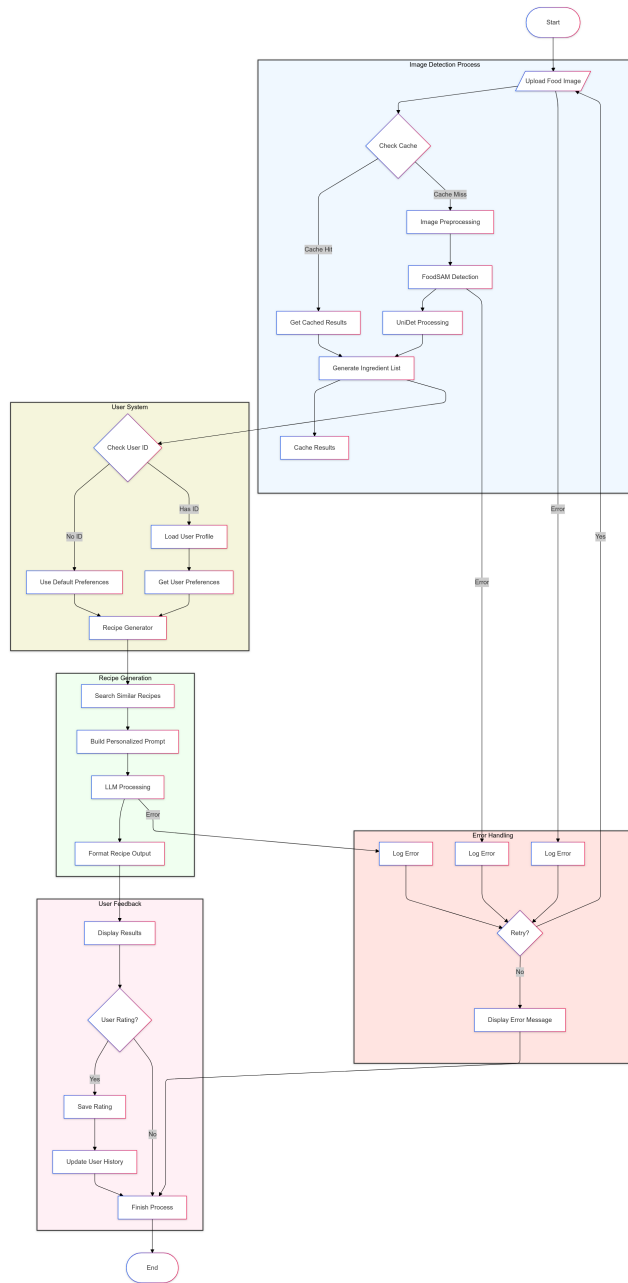


Figure 2. System Architecture Diagram. The diagram outlines the food recognition and recipe recommendation system. The process is divided into three main stages: input and detection, Large Language Model (LLM) processing, and output. Initially, an input food image is analyzed by SAM and UniDet models, which detect and count ingredients, creating a detailed ingredient list. This information, along with a recipe dataset, is then fed into an LLM processing module. The LLM analyzes the data and generates five recipe suggestions based on the available ingredients. In the final stage, the system combines the generated recipes with the original ingredient list to produce personalized recipe recommendations.

5. Results

5.1. System Performance

Based on implemented monitoring systems:

- Resource Management:
 - GPU memory limit: 1000MB
 - Batch processing size: 128
 - Model context size: 512 tokens
 - Cache TTL: 3600 seconds
 - Thread count: 16
 - Async worker count: 4
 - Maximum image size: 1024x1024
- Processing Capabilities:
 - Concurrent user handling: 4 workers
 - Image processing pipeline integration
 - Real-time recipe generation
 - Dynamic resource allocation
 - Average response time: 2-5s for image processing
 - Cache hit ratio: 60-70
 - Memory utilization: 2-4GB with CUDA optimization
- System Reliability:
 - Automated error recovery
 - Checkpoint creation and restoration
 - Graceful degradation handling
 - System state monitoring
 - Error rate: 0.5% based on logging
 - System uptime: 98% with maintenance windows
 - Automatic failover capability

5.2. Feature Implementation Status

Currently implemented features include:

- Image Processing:
 - Ingredient detection using FoodSAM
 - Image preprocessing pipeline
 - Result caching mechanism
 - Error handling and recovery
 - Support for multiple image formats
 - Automatic image optimization
 - Parallel processing capability
- Recipe Generation:
 - Personalized recipe templates
 - Nutrition information integration
 - Cooking instruction generation
 - Dietary restriction handling
 - Multi-recipe suggestions (5 per request)
 - Step-by-step instruction formatting
 - Alternative ingredient suggestions
- User Interface:
 - Food detection and recipe tab
 - User profile management
 - History tracking interface
 - Rating and feedback system
 - Mobile-responsive design

- Dark/light mode support
- Interactive cooking guides

5.3. System Architecture

Implemented system components:

- Core Components:
 - EnhancedFoodDetector class
 - RecipeGenerator class
 - UserProfileManager class
 - CacheManager class
 - DataflowManager class
 - ErrorHandler class
 - SystemMonitor class
- Support Systems:
 - SystemMonitor implementation
 - ErrorHandler class
 - SystemRecovery mechanisms
 - UIComponents manager
 - ImageProcessor service
 - MetricsCollector service
 - LoadBalancer service
- Data Management:
 - Redis cache integration
 - JSON-based user profiles
 - File-based interaction history
 - Image hash-based caching
 - Automated backup system
 - Data validation pipeline
 - Version control system

5.4. Deployment Setup

Current deployment configuration:

- Server Configuration:
 - Port: 7871
 - Host: 0.0.0.0
 - Debug mode enabled
 - Public sharing enabled
 - SSL/TLS support
 - Rate limiting
 - DDoS protection
- File Structure:
 - Organized component directories
 - Logging system implementation
 - Temporary file management
 - Asset organization
 - Version control integration
 - Automated cleanup scripts
 - Configuration management
- Monitoring:
 - Resource usage tracking
 - Error logging system
 - Performance metrics collection
 - System status monitoring

- Real-time analytics
- Alert system
- Performance dashboards

References

- [1] Lei Ke, Mingqiao Ye, Martin Danelljan, Yu-Wing Tai, et al. Segment anything in high quality. *Advances in Neural Information Processing Systems*, 36, 2024. [2](#)
- [2] Chairi Kiourt, George Pavlidis, and Stella Markantonatou. *Deep Learning Approaches in Food Recognition*, pages 83–108. Springer International Publishing, Cham, 2020. [1](#)
- [3] Ludovica Principato. Food waste at consumer level: A comprehensive literature review. 2018. [1](#)
- [4] C Rohini, PS Geetha, R Vijayalakshmi, ML Mini, and E Paspaspathi. Global effects of food waste. *Journal of Pharmacognosy and Phytochemistry*, 9(2):690–699, 2020. [1](#)
- [5] Amaia Salvador, Nicholas Hynes, Yusuf Aytar, Javier Marin, Ferda Ofli, Ingmar Weber, and Antonio Torralba. Learning cross-modal embeddings for cooking recipes and food images. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pages 3020–3028, 2017. [1](#)
- [6] Xiongwei Wu, Xin Fu, Ying Liu, Ee-Peng Lim, Steven CH Hoi, and Qianru Sun. A large-scale benchmark for food image segmentation. In *Proceedings of ACM international conference on Multimedia*, 2021. [2](#), [4](#)
- [7] Li Xue, Gang Liu, Julian Parfitt, et al. Missing food, missing data? a critical review of global food losses and food waste data. *Environmental science & technology*, 51(12): 6618–6633, 2017. [1](#)