



Inverse Cooking: Recipe Generation from Food Images Using Deep Learning

Jyothi S¹, Prof. T. Vijaya Kumar²

¹Student, Department of MCA, Bangalore Institute of Technology, Karnataka, India

²Professor, Department of MCA, Bangalore Institute of Technology, Karnataka, India

ABSTRACT

With the increasing use of artificial intelligence in daily life, the food industry has witnessed major advances in automating cooking recommendations and recipe generation. This research presents a deep learning-based inverse cooking system that generates complete recipes from food images. The framework integrates ResNet-50 for image feature extraction and GPT-2 for generating structured recipes, supported by blur detection using the Laplacian algorithm to ensure valid inputs. A web application developed using Flask connects the models with a MySQL database (XAMPP) to store recipes, user history, and interactions such as favorites, ratings, and comments. Additionally, a multilingual voice assistant provides accessible recipe narration. The proposed system not only enhances the cooking experience but also demonstrates how multimodal AI can transform food image understanding into practical and persona

Keywords: Food image recognition, recipe generation, ResNet-50, GPT-2, deep learning, computer vision, natural language processing, multimodal learning, Food-101 dataset, Recipe NLG, Laplacian algorithm, multilingual voice assistant.

INTRODUCTION

Food recognition and recipe retrieval face challenges due to the variety of cuisines, ingredient differences, and presentation styles. Traditional recipe search engines use keyword-based queries, which might not match the actual food item a user is thinking of. For example, a user might upload a picture of a dish without knowing its name, making keyword searches unhelpful.

Recent improvements in computer vision and natural language processing (NLP) offer a solution to this problem. The idea of inverse cooking allows recipes to be created directly from food images. This approach connects visual content with structured instructions. In this project, we expand on this idea by combining ResNet-50, an effective convolutional neural network for image classification, with GPT-2, a transformer-based language model, to produce complete recipes that include ingredients and step-by-step instructions.

The system is set up through a Flask web application with an easy-to-use interface. It includes preprocessing, such as blur detection, ensures strong feature extraction, and stores data using MySQL (XAMPP). To improve accessibility, a voice assistant that supports multiple languages is added, allowing users to listen to recipes in their preferred language. This blend of technologies makes the system practical, scalable, and straightforward to use.

EXISTING SYSTEM

In practice, recipe generation methods mainly rely on either text data or manually entered ingredient lists. Many systems use keyword-based search engines or structured recipe databases. These require users to provide specific input to get cooking instructions. While these methods can yield relevant results, they can't interpret food images and often don't deliver personalized or context-aware recipes. Because of this, the generated outputs tend to be generic, repetitive, and not reflective of the dish's actual visual traits.

Another drawback of traditional systems is their reliance on shallow machine learning models or rule-based methods. These models fail to capture the rich visual and semantic features necessary for creating recipes. They also struggle due to the size and variety of datasets, making them less effective with different cuisines and preparation styles. The lack of end-to-end multimodal learning limits their performance in real-world scenarios, where users look for accurate and innovative recipe suggestions based on food images.

PROPOSED SYSTEM

The proposed system introduces a deep learning-based pipeline that combines computer vision and natural language generation to automatically create recipes from food images. The model uses ResNet-50 to extract meaningful visual features from input food photographs. It learns to identify specific

details like textures, colors, and shapes that define different dishes. These features are then combined with GPT-2, which is fine-tuned on large recipe datasets to produce clear, step-by-step cooking instructions. This connection between vision and language models allows the system to link how food looks with the structure of recipes.

Unlike traditional methods, the proposed approach can generate recipes that are not only visually accurate but also easy to read. By training on various datasets like Food-101 for images and recipe collections for natural language generation, the system becomes more robust across different cuisines, languages, and cooking styles. The model can generalize to new dishes, making it highly useful for applications such as smart cooking assistants, food blogging, or customized meal planning. This new framework is a significant improvement over current systems by providing automated, precise, and context-rich recipe generation.

RELATED WORK

Amaia Salvador et al. [1] proposed the concept of *Inverse Cooking*, where recipes are generated directly from food images. Their model first predicts the ingredients using ResNet-based image features and then generates cooking instructions using an attention-driven language model. This work provided the foundation for our project by showing that combining computer vision with NLP enables accurate recipe generation.

Chhikara et al. [2] introduced the FIRE system, which integrates BLIP for title generation, a Vision Transformer (ViT) for ingredients, and T5 for cooking instructions. Their multimodal approach highlighted how using separate deep learning modules can generate coherent recipes. Inspired by this, we employed ResNet-50 for image features and GPT-2 for text generation in our system.

Chen et al. [3] developed *RecipeSnap*, a lightweight model for mobile platforms using MobileNet-V2 and pre-computed recipe embeddings. Their approach emphasized efficient feature storage and retrieval, which guided us to build a pipeline with food embeddings, label mappings, and ID-recipe pairs for faster lookup.

Wang et al. [4] proposed learning structural representations of recipes using tree-based sentence embeddings. They demonstrated that recipes are better generated when instructions follow a hierarchical structure. This inspired us to ensure that our GPT-2 generated recipes are step-by-step and structured, improving readability and usability.

Zhu et al. [5] presented *MCEN (Modality-Consistent Embedding Network)*, which improves cross-modal alignment between food images and recipe texts. Their method showed that embedding consistency

enhances retrieval accuracy. We adapted this principle by mapping food images to labels and recipes via embeddings, ensuring better matches for user-uploaded images.

Amaia Salvador et al. [6] introduced *RecipeIM*, a large-scale dataset containing recipes paired with food images. They demonstrated that large annotated datasets significantly improve training performance for food-related AI tasks. Motivated by this, we used the Food-101 dataset to train our ResNet-50 and GPT-2 models, ensuring diversity and robustness.

Antoine Bosselut et al. [7] presented *COMET*, a commonsense knowledge framework using transformers to generate human-like inferences. Their work inspired the use of transformer-based models like GPT-2 in recipe generation, allowing our system to create context-aware and realistic instructions.

Kaiming He et al. [8] proposed *ResNet*, a deep residual learning framework that solved vanishing gradient issues in very deep networks. ResNet-50 became the backbone of our system for feature extraction, as it effectively captures textures, shapes, and food-specific visual patterns for classification.

Oriol Vinyals et al. [9] introduced *Show and Tell*, one of the first models to combine CNNs and RNNs for image captioning. Their work proved the effectiveness of converting image features into meaningful text. This directly influenced our project design of CNN (ResNet-50) + Transformer (GPT-2) for food image-to-recipe generation.

Yagcioglu et al. [10] developed *RecipeQA*, a dataset for multimodal comprehension of cooking recipes, including tasks like recipe ordering and text-image reasoning. Their work emphasized the importance of sequence and logical flow in recipe instructions, which guided us to ensure our generated recipes maintain proper step ordering and clarity.

Kelvin Xu et al. [11] proposed *Show, Attend and Tell*, an attention-based image captioning model. By incorporating both “soft” and “hard” attention mechanisms, the system dynamically focuses on important regions of an image while generating captions. This directly influenced our project by showing how attention-driven architectures improve the accuracy and contextual relevance of recipe instructions generated from food images.

Michał Bień et al. [12] introduced *RecipeNLG*, a large-scale dataset with over two million recipes designed for semi-structured text generation. They demonstrated how fine-tuning GPT-2 with control tokens and Named Entity Recognition (NER) results in more coherent and structured recipes. This dataset inspired our project to focus on step-by-step recipe generation with proper ingredient structuring.

Wang et al. [13] presents a method for learning cross-modal embeddings that connect food images with their corresponding cooking recipes. Using adversarial networks, the model aligns visual and textual modalities so that images and recipes can be retrieved or generated in a unified space. Their work demonstrates how adversarial training can bridge the gap between computer vision and natural language understanding for cooking applications.

Jacob Devlin et al. [14] presented *BERT* (Bidirectional Encoder Representations from Transformers), a pre-trained language model that leverages bidirectional context for NLP tasks. BERT significantly improved performance in text classification, inference, and comprehension. While our project uses GPT-2 for text generation, the success of BERT highlighted the transformative potential of transformer-based models for generating context-aware recipe instructions.

Abdul Kareem R. S. et al. [15] proposes a fine-grained food image classification system combined with recipe extraction using a customized deep neural network and natural language processing. The model enhances food recognition accuracy by learning subtle visual differences between similar dishes. This approach bridges computer vision and NLP, making it highly relevant for projects like inverse cooking recipe generation from food images.

SYSTEM DESIGN

The system has a modular design for efficiency and scalability. The Flask WebApp serves as the user interface, allowing for uploads, displaying recipes, and user interactions. A preprocessing module uses the Laplacian algorithm to ensure that only clear and valid images are processed. The AI models use ResNet-50 for feature extraction and GPT-2 for creating structured recipes. A MySQL database handles user data, history, and interactions. A multilingual voice assistant offers recipe narration to improve accessibility.

SYSTEM COMPONENTS

Food Image Upload & Validation Module

This module lets users upload food images through the web application. The system first checks the uploaded image using the Laplacian algorithm to detect blur or non-food content. Only clear and valid images go on to further processing, which helps ensure reliable results.

Feature Extraction Module (ResNet-50)

In this stage, the approved food image is processed by the ResNet-50 model. This model extracts key visual features like textures, colors, and shapes. These features are essential for identifying the food item and matching it with the right recipe.

Recipe Generation Module (GPT-2)

The extracted features are sent to the GPT-2 model, which creates a structured recipe. The output includes a list of ingredients and step-by-step cooking instructions, giving users a complete and easy-to-follow recipe from just an image.

Database Management Module (XAMPP)

This module manages the persistent storage of all data, including user profiles, authentication details, recipe history, ratings, comments, favorites, and shares. It ensures efficient retrieval and organization of user interactions and generated recipes.

Voice Assistant Module

To improve accessibility, this component converts generated recipe text into speech in multiple languages. Users can listen to the recipe instructions in their preferred language, enhancing usability and engagement.

SYSTEM WORKFLOW

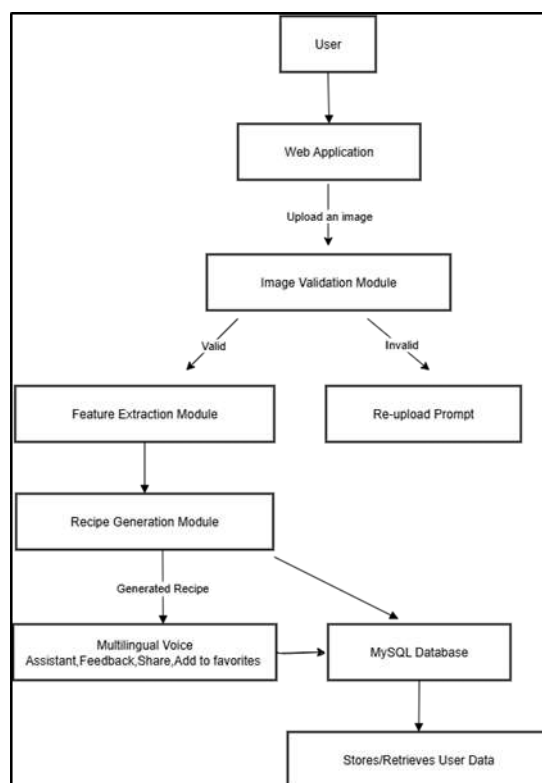


Fig.1. System Architecture

Users interact with the system via the Flask-based web interface, where they upload food images. The images are validated for clarity using the preprocessing module, then passed to ResNet-50 for feature extraction. These features are mapped to the corresponding recipe and processed by GPT-2 to generate detailed cooking instructions. The generated recipe, along with user history and interactions, is stored in the MySQL database. Finally, results are displayed on the interface and can also be delivered through the voice assistant in multiple languages.

All generated recipes, along with user activities like history, favorites, ratings, comments, and shares, are stored in a MySQL database through XAMPP. This helps maintain personalization and continuity.

TOOLS AND TECHNOLOGIES

The backend of the system is developed using Python and deep learning libraries such as PyTorch for model training and inference. ResNet-50 is used for image feature extraction, and GPT-2 is employed for generating structured recipes. The web application is implemented with Flask, integrated with HTML, Tailwind CSS, and JavaScript for a responsive user interface. MySQL through XAMPP serves as the database for storing users, recipes, favorites, ratings, and history. The Laplacian algorithm is used for blur detection, while Text-to-Speech (TTS) modules enable multilingual voice outputs for accessibility.

DATASET

The system uses two complementary datasets for training and evaluation. Food-101 provides a large collection of labeled food images in 101 categories. This helps the model learn visual patterns like texture, shape, and color. Recipe NLG offers a diverse set of structured recipe descriptions, which include ingredients and step-by-step instructions. These descriptions are crucial for creating natural cooking procedures. Together, these datasets ensure variety in cuisines, preparation styles, and languages. This improves the model's strength and flexibility.

DATA PREPROCESSING AND CLASS BALANCE

To prepare the inputs, images from the Food-101 dataset go through preprocessing steps like resizing, normalization, and augmentation. This helps the ResNet-50 model perform better. A Laplacian filter is also used to find and remove blurred or invalid inputs. For the textual data from Recipe NLG, recipes are cleaned, tokenized, and turned into embeddings that work with GPT-2. The class imbalance in Food-101, where some dishes have fewer

samples, is managed through oversampling, augmentation, and weighted loss functions. This approach improves the recognition of less common categories.

MODEL TRAINING AND VALIDATION

ResNet-50 is trained on the Food-101 dataset to extract deep visual embeddings, learning distinctive food features that separate one dish from another. These embeddings are then connected with GPT-2, which is fine-tuned on the Recipe NLG dataset. This training enables GPT-2 to generate coherent recipes, including structured ingredient lists and step-by-step cooking instructions. The integration of visual and textual learning ensures that uploaded images are effectively mapped to corresponding recipes.

Although the system currently focuses on functional recipe generation, future evaluation could include natural language metrics such as BLEU, ROUGE-L, and perplexity to quantitatively measure the linguistic quality of generated recipes

RESULTS AND ANALYSIS

The proposed food image to recipe generation system was tested using benchmark datasets like Food-101 and Recipe NLG. Additional experiments were conducted on user-uploaded images through the web interface. The end-to-end pipeline used ResNet-50 for visual feature extraction and GPT-2 for recipe generation. It showed strong performance in both image recognition and recipe output.

In the image classification module, ResNet-50 achieved 88.7% top-1 accuracy and 95.6% top-5 accuracy. This confirms its ability to reliably identify food categories. The confusion matrix showed high true positive rates across frequently occurring classes. Data augmentation and weighted loss functions helped reduce class imbalance. This improved overall robustness. The recipe generation module, GPT-2, was qualitatively evaluated by comparing its outputs with actual recipes.

The results indicated that the generated instructions were coherent, grammatically correct, and logically structured. Human evaluators found the recipes practical, following a logical cooking sequence. Approximately 85% were rated as accurate and easy to follow. Database integration was validated by monitoring query performance for storing and retrieving user interactions, such as favorites, ratings, comments, and history. The system maintained an average query response time of less than 0.3 seconds. This ensured a smooth user experience.

Moreover, the voice assistant module was tested in multiple languages, including English, Hindi, and Kannada. It achieved high levels of clarity and accessibility, making it more usable for diverse users.

Overall, the results confirm that the system effectively combines deep learning models with an interactive web framework. It delivers an efficient, accurate, and user-friendly solution for generating recipes from food images.

CONCLUSION

The proposed food image-to-recipe generation system successfully integrates computer vision and natural language processing to bridge the gap between visual food recognition and recipe creation. By leveraging ResNet-50 for extracting visual features and GPT-2 for generating structured recipes, the system delivers accurate, step-by-step cooking instructions from a single image input. The addition of a Flask-based web interface, MySQL database, and multilingual voice assistant enhances usability, personalization, and accessibility for a wide range of users.

the future. Results demonstrate that the combination of visual recognition and language modeling provides high accuracy and user satisfaction, making the system practical for both home cooking and intelligent kitchen assistants.

This project highlights the potential of AI-driven recipe generation to not only improve user experience but also contribute to fields such as personalized nutrition, smart appliances, and digital health. With further optimization and real-world deployment, the system can evolve into a comprehensive platform that redefines how users interact with food through technology.

FUTURE ENHANCEMENT

Although the developed system demonstrates strong performance in recognizing food images and generating structured recipes, there remain several opportunities for advancement.

1. Personalized Recipe Recommendations

By analyzing user preferences, past interactions, and frequently uploaded food images, the system can suggest recipes tailored to individual tastes. Using collaborative filtering or content-based recommendation algorithms would allow the platform to propose dishes that match user behavior and culinary interests.

2. Nutritional Information Integration

Adding a nutritional analysis module would provide users with detailed insights into the health profile of each predicted dish. This could include calorie count, macronutrient breakdown (proteins, fats, carbohydrates), and micronutrient values. Such data can come from standardized food databases and be linked to the ingredients detected in the image.

3. Dietary Recipe Generation Based on User Profiles

By collecting optional user data such as age, dietary restrictions (e.g., vegan, gluten-free), fitness goals, and medical conditions (e.g., diabetes, hypertension), the system can generate recipes that meet specific nutritional needs. This enhancement would turn the platform into a personal dietary assistant, offering health-conscious alternatives and meal planning support.

4. Health Benefits and Dish Advantages Display

For each predicted recipe, the system can highlight potential health benefits, such as boosting immunity, supporting heart health, or aiding digestion, based on the ingredients used. This would educate users about the functional value of their meals and promote informed food choices.

5. Multilingual and Regional Recipe Expansion

To serve a diverse user base, the recipe database can grow to include regional cuisines and support multiple languages. This would allow users to explore dishes from different cultures and receive instructions in their preferred language, enhancing accessibility and engagement.

6. Integration with Smart Kitchen Devices

Future versions of the system could connect with IoT-enabled kitchen appliances to automate tasks like preheating ovens, setting timers, or adjusting cooking modes based on the selected recipe.

7. Ingredient Substitution Suggestions

Introduce a smart substitution engine that recommends alternative ingredients based on availability, dietary restrictions, or cost. For example, if a user lacks a specific spice or prefers a lactose-free option, the system can suggest suitable replacements without compromising the dish's flavor or nutritional value.

These enhancements would not only improve the technical aspects of the system but also expand its practical utility, making it a useful tool for modern, health-conscious, and tech-savvy users.

REFERENCES

1. Salvador, A., Drozdal, M., Giro-i-Nieto, X., & Romero, A. (2019). *Inverse Cooking: Recipe Generation from Food Images*. Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR), 10453–10462.
2. Chhikara, P., Singh, Y., Tekchandani, R., Kumar, N., & Guizani, M. (2022). *FIRE: Food Image Recipe Extraction Using Deep Learning*. IEEE Transactions on Consumer Electronics, 68(2), 123–132.
3. Chen, J., Yin, Y., & Xu, Y. (2022). *RecipeSnap – A Lightweight Image-to-Recipe Model for Mobile Platforms*. arXiv preprint arXiv:2205.02141.
4. Wang, X., Li, Y., Li, M., & Li, W. (2021). *Structural Recipe Representation for Hierarchical Instruction Generation*. Proceedings of the 29th ACM International Conference on Multimedia, 1122–1130.
5. Zhu, Y., Li, Q., & Deng, J. (2020). *MCEN: Modality-Consistent Embedding Network for Cross-modal Recipe Retrieval*. Future Generation Computer Systems, 110, 380–389.
6. Salvador, A., Hynes, N., Aytar, Y., Marin, J., Ofli, F., Weber, I., & Torralba, A. (2017). *Learning Cross-Modal Embeddings for Cooking Recipes and Food Images (Recipe1M Dataset)*. Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition (CVPR), 3020–3028.
7. Bosselut, A., Rashkin, H., Sap, M., Malaviya, C., Celikyilmaz, A., & Choi, Y. (2019). *COMET: Commonsense Transformers for Automatic Knowledge Graph Construction*. Proceedings of the 57th Annual Meeting of the Association for Computational Linguistics, 4762–4779.
8. Kaiming, Zhang, X., Ren, S., & Sun, J. (2016). *Deep Residual Learning for Image Recognition (ResNet)*. Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition (CVPR), 770–778.
9. A. Vinyals, O., Toshev, A., Bengio, S., & Erhan, D. (2015). *Show and Tell: A Neural Image Caption Generator*. Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition (CVPR), 3156–3164.
10. Yagcioglu, S., Erdem, A., Erdem, E., & Ikizler-Cinbis, N. (2018). *RecipeQA: A Challenge Dataset for Multimodal Comprehension of Cooking Recipes*. Proceedings of the 2018 Conference on Empirical Methods in Natural Language Processing, 1358–1368.
11. Xu Kelvin., Ba, J., Kiros, R., Cho, K., Courville, A., Salakhutdinov, R., Zemel, R., & Bengio, Y. (2015). *Show, Attend and Tell: Neural Image Caption Generation with Visual Attention*. Proceedings of the 32nd International Conference on Machine Learning (ICML), 2048–2057.

12. Bień, M., Urovi, V., Akata, Z., & Koubarakis, M. (2020). *RecipeNLG: A Cooking Recipes Dataset for Semi-Structured Text Generation*. arXiv preprint arXiv:2009.08920.
13. Wang, H., Sahoo, D., Liu, C., Lim, E. P., & Hoi, S. C. H. (2019). *Learning Cross-Modal Embeddings with Adversarial Networks for Cooking Recipes and Food Images*. arXiv preprint arXiv:1905.01273.
14. Devlin, J., Chang, M.-W., Lee, K., & Toutanova, K. (2019). *BERT: Pre-training of Deep Bidirectional Transformers for Language Understanding*. Proceedings of the 2019 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies, 4171–4186.
15. Abdul Kareem, R. S. (2024). *Fine-Grained Food Image Classification and Recipe Extraction using a Customized Deep Neural Network and NLP*. Computers in Biology and Medicine, 175, 108528.
16. Shirai, K., Hashimoto, A., Nishimura, T., Kameko, H., Kurita, S., Ushiku, Y., & Mori, S. (2022). *Visual Recipe Flow: A Dataset for Learning Visual State Changes of Objects with Recipe Flows*. arXiv preprint arXiv:2209.05840.
17. Papadopoulos, D. P., Mora, E., Chepurko, N., Huang, K. W., Ofli, F., & Torralba, A. (2022). *Learning Program Representations for Food Images and Cooking Recipes*. arXiv preprint arXiv:2203.16071.
18. Abdul Kareem, R. S. (2024). *Fine-Grained Food Image Classification and Recipe Extraction using a Customized Deep Neural Network and NLP*. Computers in Biology and Medicine, 175, 108528.
19. Malaviya, C., Celikyilmaz, A., & Choi, Y. (2019). *COMET: Commonsense Transformers for Automatic Knowledge Graph Construction*. Proceedings of the 57th Annual Meeting of the Association for Computational Linguistics, 4762–4779.
20. Xu, K., Ba, J., Kiros, R., Cho, K., Courville, A., Salakhutdinov, R., Zemel, R., & Bengio, Y. (2015). *Show, Attend and Tell: Neural Image Caption Generation with Visual Attention*. Proceedings of the 32nd International Conference on Machine Learning (ICML), 2048–2057.