

YOLOV5-BASED FRAMEWORK FOR REAL-TIME INDIAN FOOD CLASSIFICATION

Suhas H. Patel¹, Tejaskumar B. Sheth², Ujjval R. Dave³, Mitul B. Patel⁴

^{1,3}Assistant Professor, Electronics and Communication Engineering Department, Government Engineering College, Gandhinagar, Gujarat India.

²Associate Professor, Electronics and Communication Engineering Department, Government Engineering College, Gandhinagar, Gujarat India.

⁴Assistant Professor, Biomedical Engineering Department, Government Engineering College, Gandhinagar, Gujarat India.

ABSTRACT

This paper presents a YOLOv5-based framework for real-time Indian food classification, aimed at automating food recognition in various applications such as dietary tracking and restaurant management. The model was developed using a dataset comprising 4,298 images across 12 distinct Indian dishes, split into 2,990 training images, 833 validation images, and 475 testing images. The framework was trained over 10, 50 and 100 epochs, achieving an overall accuracy of 85%. The model's classification performance varies across different food items, with precision, recall, and F1-scores analyzed for each class. Notable results include an F1-score of 0.97 for Gulab Jamun and 0.94 for Rasgulla, indicating high reliability in recognizing these items. However, some classes, such as Jalebi and Momos, showed lower F1-scores of 0.65 and 0.71, respectively, suggesting areas for further model refinement. This study highlights the challenges of distinguishing visually similar food items and demonstrates the capability of YOLOv5 in handling real-time classification tasks. The findings suggest that with further optimization, this framework can be effectively deployed in real-world applications, enhancing the accuracy and efficiency of automated food classification systems.

Keywords: Indian Food Classification, YOLOv5, CNN, Deep Learning, VGG16.

1. INTRODUCTION

In recent years, the rapid advancement of computer vision and machine learning has significantly transformed various industries, including food technology[1]. Among the numerous applications, food recognition and classification have emerged as vital components in automating dietary tracking, enhancing restaurant management systems, and supporting health and wellness programs[2]. With the proliferation of mobile devices and the increasing availability of image data, there is a growing need for efficient and accurate food classification systems that can operate in real-time[3]. This paper introduces a YOLOv5-based framework tailored for real-time Indian food classification, addressing the unique challenges posed by the diversity and complexity of Indian cuisine[4].

Food classification plays a crucial role in several domains. In the healthcare sector, accurate food identification can assist in monitoring dietary intake, managing nutritional information, and supporting patients with specific dietary requirements[5]. In the hospitality industry, automated food recognition systems can streamline operations in restaurants, enabling faster service and reducing errors in food delivery. Moreover, food classification systems are increasingly being integrated into mobile applications that help users track their meals, calculate caloric intake, and maintain balanced diets[6]. The ability to classify food items accurately and in real-time has the potential to revolutionize these industries, providing convenience, efficiency, and enhanced user experiences.

Indian cuisine is renowned for its rich diversity, characterized by a vast array of dishes with unique appearances, textures, and ingredients. This diversity, while culturally enriching, presents significant challenges for automated food classification systems. Unlike standardized packaged foods, Indian dishes are often complex, with multiple ingredients and varied presentations that can make accurate classification difficult. Additionally, different regions of India have their own versions of similar dishes, further complicating the classification task. The visual similarities between certain dishes, combined with variations in preparation methods and serving styles, necessitate a robust and sophisticated approach to food classification. Figure 1 shows different Indian Food in (a), (b), (c), and (d). The advent of deep learning has opened new avenues for tackling complex image classification tasks. Convolutional Neural Networks (CNNs) have proven to be particularly effective in image recognition, leading to significant advancements in the field of food classification[2]. Early approaches relied on traditional machine learning techniques, which required extensive feature engineering and often struggled with the variability and complexity of food images. However, the introduction of deep learning models, such as AlexNet, VGG, and ResNet, marked a paradigm shift, enabling more accurate and automated feature extraction from images[7]. Among the various deep learning models, YOLO (You Only Look Once) has gained popularity for its speed and accuracy in object detection and classification. YOLOv5, the latest iteration in the YOLO series, builds on this foundation, offering improvements in performance and efficiency. YOLOv5's ability to perform real-time object detection makes it an ideal choice for applications requiring rapid classification, such as mobile-based food recognition systems[8].

The remainder of this paper is organized as follows: Section 2 reviews related work, focusing on prior research in food classification and the challenges specific to Indian cuisine. Section 3 explains the methodology, including the dataset used, details of the YOLOv5 model architecture, and the training and evaluation procedures. Section 4 presents the experimental results, offering a comparative analysis of YOLOv5 performance in Indian food classification. Finally, Section 5 concludes the study by summarizing the key findings and discussing their implications.

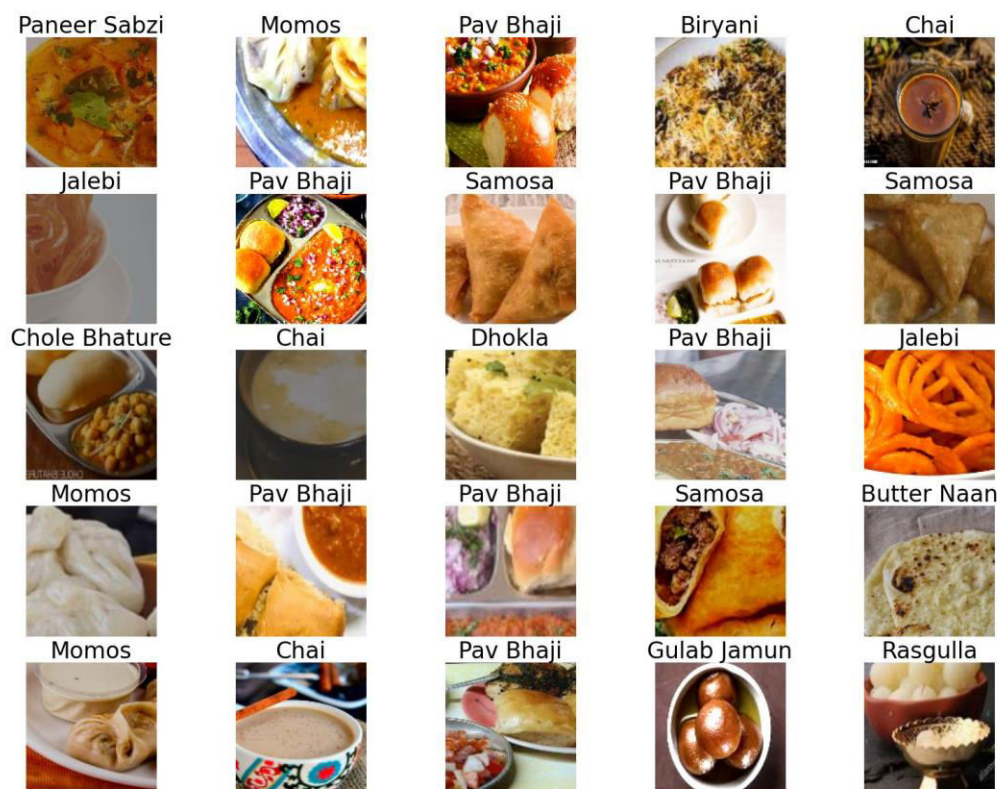


Figure 1. A Variety of Indian Food Dishes

2. RELATED WORK

Food classification is a complex task due to the vast variety of food items, their appearances, and the different presentation styles influenced by cultural practices. Traditional methods for food recognition relied heavily on handcrafted features, such as color histograms, texture descriptors, and shape-based features. These methods often involved the use of machine learning algorithms like Support Vector Machines (SVM), k-Nearest Neighbors, and decision trees [9]. However, these approaches had limitations in handling the variability and complexity of food images, especially when applied to diverse cuisines like Indian food. The dependence on manually engineered features often resulted in suboptimal performance, as these features could not capture the high-level abstractions necessary for accurate classification. The advent of deep learning has revolutionized image classification tasks, including food recognition. Convolutional Neural Networks (CNNs), with their ability to learn hierarchical feature representations directly from data, have become the foundation for modern food classification systems. Early applications of CNNs in food classification focused on popular architectures like AlexNet, VGG, and ResNet[10]. These models demonstrated significant improvements over traditional methods by achieving higher accuracy and robustness to variations in food presentation. AlexNet, introduced in 2012, marked a significant breakthrough in deep learning, winning the ImageNet Large Scale Visual Recognition Challenge (ILSVRC)[11]. Its architecture, consisting of multiple convolutional layers followed by fully connected layers, set the stage for deeper and more

complex networks. VGG, introduced later, further explored the impact of depth in CNNs by using smaller convolutional filters and increasing the number of layers[12]. This model became a popular choice for transfer learning in various image classification tasks, including food recognition. ResNet, introduced with the concept of residual learning, addressed the problem of vanishing gradients in very deep networks. By using residual connections, ResNet enabled the training of much deeper networks, achieving state-of-the-art performance in many image classification benchmarks[13]. These architectures laid the groundwork for the development of specialized models for food classification.

The need for real-time food classification has driven research towards models that balance accuracy with computational efficiency. YOLOv5 builds upon these advancements with a more streamlined and modular design, making it easier to train and deploy. YOLOv5's architecture includes a flexible backbone, a PANet-based neck for feature aggregation, and a head designed for precise object classification and localization[14]. The model's real-time performance and high accuracy make it an ideal choice for food classification tasks, especially in mobile and embedded applications where computational resources are limited.

2.1 PREVIOUS RESEARCH ON FOOD CLASSIFICATION

Despite the advancements in food classification, Indian food poses unique challenges due to its diverse and complex nature. Indian cuisine encompasses a wide range of dishes, each with its own distinct ingredients, preparation methods, and presentation styles. The visual similarities between certain dishes, combined with variations across different regions, make it difficult to achieve high classification accuracy. P. McAllister et al. [15] introduce a personalized food classification model that improves accuracy for specific dietary preferences using advanced feature extraction techniques like SURF, SFTA, and LBP. The study identifies a Neural Network as the best-performing classifier with 69.43% accuracy and suggests future research for further advancements. D. J. Attokaren et al. [16] present a novel food classification methodology using CNNs, achieving 86.97% accuracy on the Food-101 dataset. They emphasize the significance of real-time database creation for practical applications and outline future work, including developing a comprehensive dataset and exploring multi-level classification to reduce misclassifications. Few studies have specifically addressed the classification of Indian food. Research by J. R. Rajayogi et al. [17] focused on developing datasets and models specifically for Indian food, using a dataset of 20 classes with 500 images per class for training and validation. They employed models such as InceptionV3, VGG16, VGG19, and ResNet. Their experimentation revealed that Google InceptionV3 achieved the highest performance, with an accuracy of 87.9% and a loss rate of 0.5893. However, these studies often faced limitations in dataset size and diversity, impacting the generalization of the models. The lack of large, annotated datasets for Indian food has been a significant barrier to progress in this area. To address these challenges, this paper introduces a YOLOv5-based framework specifically designed for real-time Indian food classification. By leveraging the strengths of YOLOv5 in handling diverse object sizes and complex backgrounds, the proposed framework aims to

achieve high classification accuracy while maintaining the efficiency required for real-time applications.

3. METHODOLOGY

3.1 DATASET

To develop and evaluate the YOLOv5-based food classification framework, a diverse and publicly available dataset focused on 12 popular Indian dishes was utilized. As represented in Table 1, the dataset consists of 4,298 images, distributed across 2,990 training images, 833 validation images, and 475 testing images[17]. The images were sourced from various locations, capturing different presentation styles, lighting conditions, and backgrounds to simulate real-world scenarios. The selected dishes, such as Biryani, Butter Naan, Chai, Chole Bhature, Dhokla, Gulab Jamun, and more, represent a broad spectrum of Indian cuisine, offering a complex and varied dataset ideal for testing the model's robustness and accuracy.

Table 1. Indian Food Classification Dataset[17]

Total Images	4298
Classes	12
Training Set	2990 (70%)
Validation Set	833 (19%)
Testing Set	475 (11%)
Class Instances	
Chole Bhature	606
Pav Bhaji	535
Samosa	479
Dhokla	442
Chai	343
Paneer Sabzi	334
Momos	316
Butter Naan	298
Jalebi	287
Biryani	267
Gulab Jamun	233
Rasgulla	158

3.2 YOLOv5 ARCHITECTURE & MODEL TRAINING

As shown in figure 2 the flowchart visually outlines the process of Indian food classification using the YOLOv5 framework. It begins with data collection, where images of various Indian dishes are gathered. The next step is model selection, where YOLOv5 is chosen for object detection tasks. Following this, the model is trained on dataset. Once

training is complete, the model's performance is evaluated using various metrics. Finally, the model is applied to classify Indian food items, and the process is concluded. This step-by-step workflow ensures a systematic approach to developing and deploying a YOLOv5-based food classification model.

Figure 3 illustrates the architecture of YOLOv5, a single-stage object detection model designed to simplify the object detection task. Unlike traditional multi-stage approaches, YOLOv5 directly predicts bounding boxes and class probabilities from image pixels in a single pass through the network. This design treats object detection as a regression problem, optimizing for both spatial and categorical predictions simultaneously. The architecture balances speed and accuracy, making it effective for real-time applications where rapid detection and classification are crucial. The architecture of YOLOv5 consists of three main components: the backbone, the neck, and the head. The backbone, typically a CSPDarknet model, is responsible for feature extraction. The neck, which includes layers like the PANet (Path Aggregation Network), helps in enhancing feature fusion from different scales. The head, consisting of a series of convolutional layers, generates the final predictions for object classes and bounding boxes. This design allows YOLOv5 to effectively handle varying object sizes and complex backgrounds, which are common in food images.

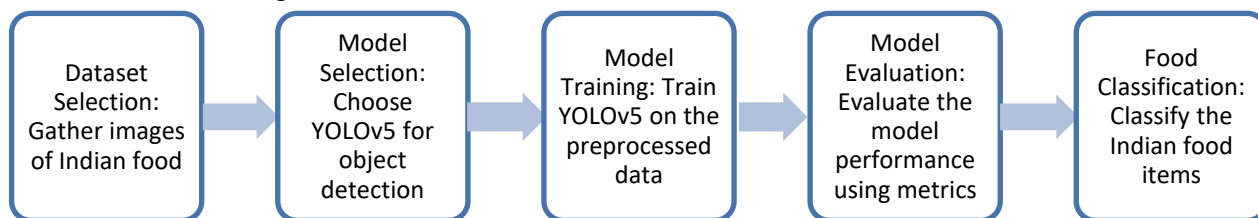


Figure 2. YOLOv5-based Indian Food Classification

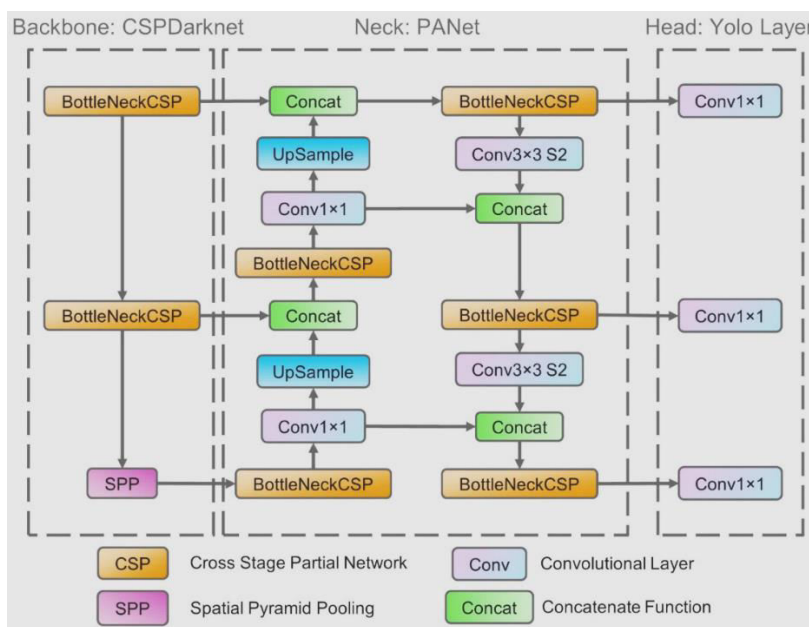


Figure 3. The network architecture of YOLOv5 [18]

The training for the food classification model was conducted on Google Colab using the YOLOv5s architecture (yolov5s-cls.pt) with image sizes set to 128x128 for both training and testing. The dataset was divided into training, validation, and test sets, with a batch size of 64. Training was carried out over 10, 50, and 100 epochs, with the AdamW optimizer managing learning rates starting at 0.001 and weight decay parameters set at 0.0 and 5e-05 for different weight groups. The model summary includes 149 layers, 4,187,852 parameters, and 4,187,852 gradients, with a total computational cost of 10.5 GFLOPs. The implementation involved several key steps, beginning with setting up the training environment and progressing through model training. It concluded with deploying the model for real-time food classification, ensuring the system could accurately identify and classify various Indian dishes in real-time scenarios.

4. Results and Discussion

The performance of the YOLOv5 model, as shown in Table 3, was evaluated across 10, 50, and 100 epochs, demonstrating a consistent improvement in classification accuracy over time. After 10 epochs, the model achieved a precision of 0.53, recall of 0.48, and an F1 score of 0.50, with an overall accuracy (mAP@0.5) of 47.00%, indicating early learning stages. By 50 epochs, the model's precision, recall, and F1 score significantly improved to 0.80, 0.79, and 0.79, respectively, with an overall accuracy of 79.00%. At 100 epochs, the model further refined its performance, achieving a precision of 0.87, recall of 0.85, F1 score of 0.86, and an overall accuracy of 85.00%. This progression illustrates the model's growing ability to accurately classify images as training continued, with the most notable gains occurring between 10 and 50 epochs. This suggests that 50 epochs may offer an optimal balance between computational efficiency and model performance for similar tasks.

Table 3. Accuracy of YOLOv5 Model for Indian Food Classification

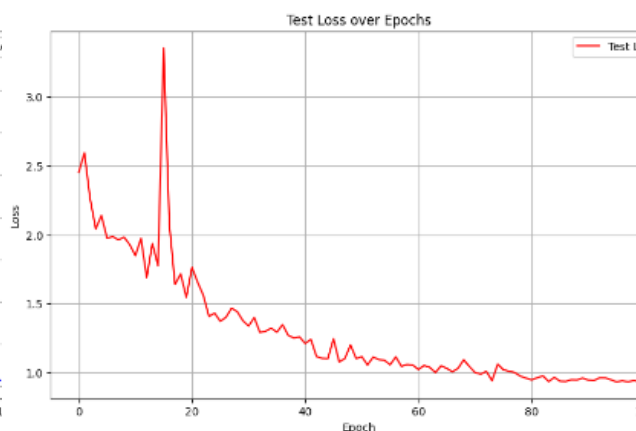
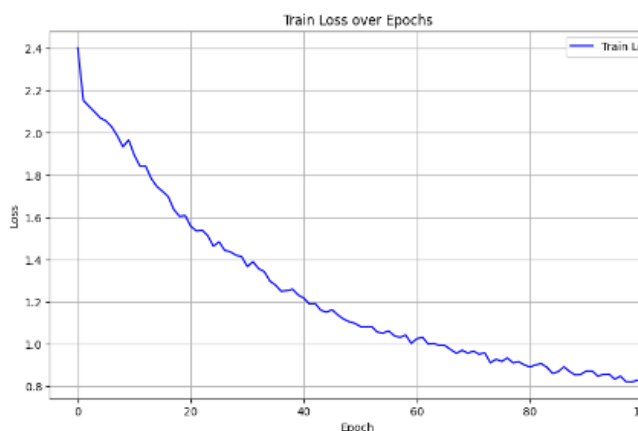
YOLOv5 Model	No. of epochs	CLASS	PRECISION	RECALL	F-1 SCORE	Overall Accuracy (mAP@0.5)
yolov5s-cls	10	ALL	0.53	0.48	0.50	47.00%
	50	ALL	0.8	0.79	0.79	79.00%
	100	ALL	0.87	0.85	0.86	85.00%

Table 4 illustrates the performance of the YOLOv5s-cls model for Indian food classification over 100 epochs, showing detailed metrics for each class. The precision values range from 0.63 to 0.89, indicating the model's ability to correctly identify positive instances across different food items. Recall values span from 0.45 to 0.91, reflecting how well the model retrieves all relevant instances. The F-1 score, which balances precision and recall, varies from 0.59 to 0.86, demonstrating overall robust performance but also highlighting areas for improvement, especially in classes like "Jalebi" and "Paneer Sabzi." The "Samosa" class shows the highest recall at 0.91, indicating that most "Samosa" instances were correctly identified, while "Jalebi" and "Paneer Sabzi" have relatively lower F-1 scores, suggesting difficulties in distinguishing

these items. The "Chole Bhature," "Dhokla," and "Gulab Jamun" classes exhibit balanced precision, recall, and F-1 scores, indicating stable model performance in these categories. Overall, the results emphasize the model's strong generalization capabilities but also reveal variability in performance across different food items, suggesting that certain classes are more challenging to classify accurately.

Table 4. Performance of YOLOv5s-cls Model for Indian Food Classification over 100 epochs

No. of Epoch	Class	PRECISION	RECALL	F-1 SCORE	Support
100	Biryani	0.77	0.85	0.81	40
	Butter Naan	0.81	0.84	0.82	25
	Chai	0.63	0.83	0.72	29
	Chole Bhature	0.86	0.86	0.86	64
	Dhokla	0.82	0.89	0.86	47
	Gulab Jamun	0.88	0.85	0.86	33
	Jalebi	0.87	0.45	0.59	29
	Momos	0.89	0.53	0.67	30
	Paneer Sabzi	0.74	0.5	0.6	28
	Pav Bhaji	0.78	0.85	0.82	60
	Rasgulla	0.83	0.62	0.71	24
	Samosa	0.71	0.91	0.79	66



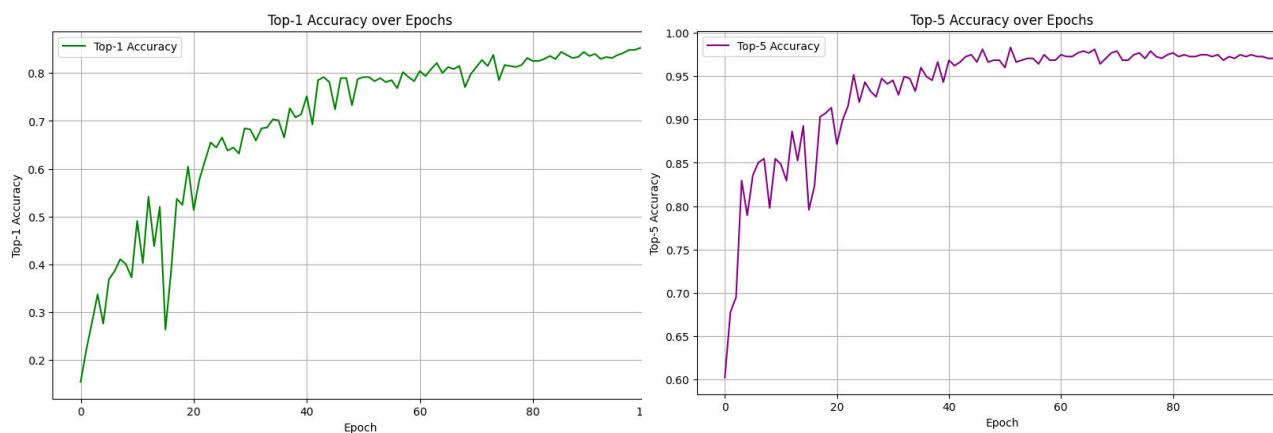


Figure 4. Simulation results obtained on YOLOv5s-cls Model for Indian Food Dataset

Figure 4 illustrates the model's performance over time, measured by epochs—each representing a full training cycle on the dataset. As training progresses, the train loss (indicating the model's fit to training data) generally decreases, showing improvement. The test loss (reflecting the model's performance on unseen data) also decreases, albeit with some fluctuations, suggesting how well the model generalizes. The accuracy metrics (top-1 and top-5) indicate the frequency of correct predictions, with top-5 accuracy being higher due to considering the model's top five guesses. Overall, the model shows improvement, with increasing accuracy and decreasing losses.

Figure 5 shows the confusion matrix of YOLOv5s-cls Model for Indian Food Dataset over 100 epochs. The matrix compares the true labels (on the vertical axis) with the predicted labels (on the horizontal axis) for various Indian dishes. The diagonal elements show the number of correct predictions for each class, while off-diagonal elements indicate misclassifications. For instance, "Chole Bhature" was correctly classified 61 times, while "Samosa" was classified correctly 59 times. The color intensity reflects the frequency of predictions, with darker shades representing higher values, indicating better performance for certain classes.

Figure 6 illustrates the results of the Indian Food Classification Model using YOLOv5s-cls after 100 epochs. The detection accuracy for various dishes is detailed as follows: Biryani at 0.61, Samosa at 0.38, Chole Bhature at 0.79, and Gulab Jamun at 0.39. These percentages reflect the model's proficiency in identifying these food items, with Chole Bhature achieving the highest detection rate among the four dishes, indicating strong model performance in recognizing this dish. On the other hand, Samosa has the lowest detection rate, suggesting that the model struggles with this class, possibly due to visual similarities with other items or inherent complexities in its features. Despite this variation, the overall results demonstrate the model's capability in classifying Indian foods with a reasonable degree of accuracy, highlighting both its strengths and areas for further improvement.

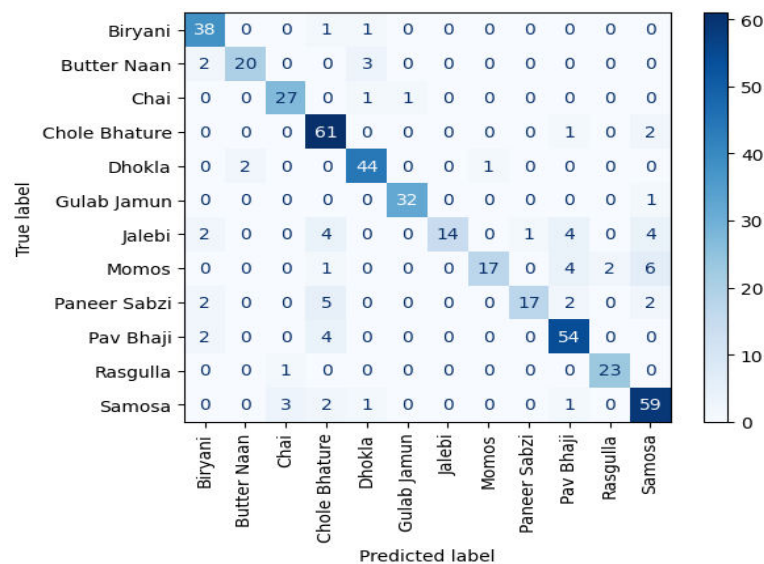


Figure 5. Confusion matrix of YOLOv5s-cls Model for 100 Numbers of epoch



(a)



(b)



(c)



(d)

Figure 6. Results of Indian food Classification Model: YOLOv5s-cls Epoch:100 (a) Biryani (b) Samosa (c) Chole Bhature (d) Gulab Jamun

5. CONCLUSION

This research paper highlights the effectiveness of the YOLOv5-based framework for real-time Indian food classification. Through extensive experimentation, the YOLOv5s-cls model achieved up to 85% mean Average Precision (mAP) after 100 epochs, showcasing its capability to accurately classify diverse Indian dishes. The study utilized a comprehensive dataset, ensuring the model's robustness under varying conditions. While the model excelled in classifying certain dishes like Biryani and Samosa, it faced challenges with others such as Jalebi and Paneer Sabzi due to visual similarities. Overall, this research provides a valuable framework for Indian food classification, contributing to the broader field of food recognition and classification using deep learning. It also paves the way for future work to refine the model, possibly by integrating multimodal data or exploring more advanced deep learning techniques to address the identified challenges.

REFERENCES

- [1] Z. Xiao, J. Wang, L. Han, S. Guo, and Q. Cui, "Application of Machine Vision System in Food Detection," *Frontiers in Nutrition*, vol. 9, no. May, pp. 1–7, May 2022, doi: 10.3389/fnut.2022.888245.
- [2] C. Kiourt, G. Pavlidis, and S. Markantonatou, "Deep Learning Approaches in Food Recognition," *Learning and Analytics in Intelligent Systems*, vol. 18, pp. 83–108, 2020, doi: 10.1007/978-3-030-49724-8_4.
- [3] Ş. Aktı, M. Qaraqe, and H. K. Ekenel, "A Mobile Food Recognition System for Dietary Assessment," in *Lecture Notes in Computer Science (including subseries Lecture Notes in Artificial Intelligence and Lecture Notes in Bioinformatics)*, vol. 13373 LNCS, 2022, pp. 71–81. doi: 10.1007/978-3-031-13321-3_7.
- [4] J. V. V. Kolla, P. C. Vemula, S. L. Chakravarthy, B. S. Naidu, and D. Patibandla, "Leveraging Transfer Learning to Identify Food Categories," *Advances in Science and Technology Research Journal*, vol. 15, no. 4, pp. 101–109, 2021, doi: 10.12913/22998624/142738.
- [5] D. Pawade, A. Dalvi, S. Irfan, M. Carvalho, P. Kotian, and G. Hima, "Cuisine Detection Using the Convolutional Neural Network," *International Journal of Education and Management Engineering*, vol. 10, no. 3, pp. 1–11, 2020, doi: 10.5815/ijeme.2020.03.01.
- [6] T. Theodoridis, V. Solachidis, K. Dimitropoulos, L. Gymnopoulos, and P. Daras, "A survey on AI nutrition recommender systems," in *Proceedings of the 12th ACM International Conference on Pervasive Technologies Related to Assistive Environments*, New York, NY, USA: ACM, Jun. 2019, pp. 540–546. doi: 10.1145/3316782.3322760.
- [7] L. Zhou, C. Zhang, F. Liu, Z. Qiu, and Y. He, "Application of deep learning in food: a review," *Comprehensive reviews in food science and food safety*, vol. 18, no. 6, pp. 1793–1811, 2019.

- [8] J. Du, "Understanding of Object Detection Based on CNN Family and YOLO," Journal of Physics: Conference Series, vol. 1004, no. 1, p. 012029, Apr. 2018, doi: 10.1088/1742-6596/1004/1/012029.
- [9] G. Ciocca, P. Napolitano, and R. Schettini, "Learning CNN-based Features for Retrieval of Food Images," in New Trends in Image Analysis and Processing -- ICIAP 2017, S. Battiato, G. M. Farinella, M. Leo, and G. Gallo, Eds., Cham: Springer International Publishing, 2017, pp. 426–434.
- [10] Y. Lu, "Food Image Recognition by Using Convolutional Neural Networks (CNNs)," pp. 1–6, 2016, [Online]. Available: <http://arxiv.org/abs/1612.00983>
- [11] A. Krizhevsky, I. Sutskever, and G. E. Hinton, "ImageNet Classification with Deep Convolutional Neural Networks," in Advances in Neural Information Processing Systems, F. Pereira, C. J. Burges, L. Bottou, and K. Q. Weinberger, Eds., Curran Associates, Inc., 2012. [Online]. Available: https://proceedings.neurips.cc/paper_files/paper/2012/file/c399862d3b9d6b76c8436e924a68c45b-Paper.pdf
- [12] K. Simonyan and A. Zisserman, "Very deep convolutional networks for large-scale image recognition," 3rd International Conference on Learning Representations, ICLR 2015 - Conference Track Proceedings, pp. 1–14, 2015.
- [13] K. He, X. Zhang, S. Ren, and J. Sun, "Deep Residual Learning for Image Recognition," in 2016 IEEE Conference on Computer Vision and Pattern Recognition (CVPR), IEEE, Jun. 2016, pp. 770–778. doi: 10.1109/CVPR.2016.90.
- [14] G. Yang et al., "Garbage Classification System with YOLOV5 Based on Image Recognition," in 2021 IEEE 6th International Conference on Signal and Image Processing (ICSIP), IEEE, Oct. 2021, pp. 11–18. doi: 10.1109/ICSIP52628.2021.9688725.
- [15] P. McAllister, H. Zheng, R. Bond, and A. Moorhead, "Towards Personalised Training of Machine Learning Algorithms for Food Image Classification Using a Smartphone Camera," in Lecture Notes in Computer Science (including subseries Lecture Notes in Artificial Intelligence and Lecture Notes in Bioinformatics), vol. 10069 LNCS, 2016, pp. 178–190. doi: 10.1007/978-3-319-48746-5_18.
- [16] D. J. Attokaren, I. G. Fernandes, A. Sriram, Y. V. S. Murthy, and S. G. Koolagudi, "Food classification from images using convolutional neural networks," IEEE Region 10 Annual International Conference, Proceedings/TENCON, vol. 2017-Decem, pp. 2801–2806, 2017, doi: 10.1109/TENCON.2017.8228338.
- [17] D. Rodge, "Indian Food Images - 12 Different Dishes." 2021.
- [18] R. Xu, H. Lin, K. Lu, L. Cao, and Y. Liu, "A forest fire detection system based on ensemble learning," Forests, vol. 12, no. 2, pp. 1–17, 2021, doi: 10.3390/f12020217.