# Multimodal Smart Bins: Fusion of Vision and Sensor Data for Enhanced Waste Classification

**Arya S Nair[1], Dr. Anju J Prakash[2]**

[1]Department of Computer Science, Amal Jyothi College of Engineering (Autonomous), Kottayam, India
Email: *aryathundiyil12[at]gmail.com*

[2]Associate Professor, Department of Computer Science, Amal Jyothi College of Engineering (Autonomous) Kottayam, India
Email: *jpanju[at]gmail.com*

**Abstract:** *Vision-based approaches have been utilized for waste classification, a considerable part of waste management, but may possess accuracy limitations in the real world due to environmental factors such as light sources, occlusions or features that look visually similar. We present a new multimodal smart bin utilizing computer vision, weight sensors, and near-infrared (NIR) spectroscopy in an effort to detect and classify urban waste with a high degree of accuracy and robustness. The additional information leveraged from the weight and spectral signature of items with visual features assists in improving accurate classification while mitigating the likelihood of misclassifying similar or visually ambiguous items, e.g. wet paper and food waste or clear plastic and glass. This paper presents a multimodal dataset of image, weight and NIRS readings for urban waste materials that classify waste with a deep learning model using late fusion of vision and multimodal weight-sensor technology. The system can be deployed on edge devices like a Raspberry Pi or NVIDIA Jetson Nano allowing waste detection and classification at bin-level, without the need for cloud server capabilities. The IoT-based integration is capable of providing remote monitoring and centralized analytics, along with actionable insights for the municipal government authority. The results show that the multimodal mechanisms of the system preformed more accurately and robustly than vision only models and furthers the contamination problem of recyclables, allowing for more reliable and automated sorting of urban waste to be undertaken at the source. The work of this dataset further supports both of the United Nations Sustainable Development Goals (SDG 11: Sustainable Cities and Communities, and SDG 9: Industry, Innovation and Infrastructure), towards developing sustainable, green, and smart cities through scalable and sustain- able solutions to the solid waste management systems. Future work can focus on building lighter weight fusion models for cost-effective deployment, adaptive calibration of sensor-driven models for long-term sustainability, and developing engagement features with the recycling and composting processes of incentive- based systems for citizen where they use the waste systems.*

**Keywords:** Multimodal learning, waste classification, smart bins, computer vision, weight sensors, near-infrared spectroscopy, IoT, sustainable cities

## 1. Introduction

The management of waste is regarded as one of the most serious challenges stemming from modern urbanization. Waste management is driven by rapid urbanization through population explosive growth, the development of more industries, and the increasing proclivity for consumerism. Failure to appropriately segregate waste at the source results in poor recycling, landfills full beyond capacity, and pollution of the environment. The sorting approaches that have been used to well sorting waste are all labour-intensive, involve substantial human error, and simply cannot be reasonably employed at scale within the context of smart cities. For smart cities, automated waste segregation is becoming more prominent in the academic world to achieve sustainable and intelligent urban infrastructure.

Recent developments of artificial intelligence (AI) or computer vision have facilitated the academic world creating automated waste classification systems. The bulk of the waste classification systems rely on deep learning models trained on the image datasets. Broadly, image-based methods and studies lead to promising results, but fall short when tested "in the wild." Real-world behaviours act to hinder instantaneous classification of waste (variability in light, occluding object, overlapping objects and such close visual representations as transparent plastics and glass) which all lead to greater instances of misclassification. The previously mentioned examples illustrate the need for additional modalities beyond single-modal vision inputs, and to actively improve the area of waste classification systems.

To solve these problems, this paper describes a multimodal smart bin framework that utilizes computer vision, weight sensors, and NIR spectroscopy to improve performance and reliability. By combining the data from different sensors, the visual features are further improved with complementary information. For example, we are able to distinguish between different materials, which can look very similar, thanks to weight and NIR spectroscopy material-specific spectral sig- natures, i.e. spectrum signatures that can be attributed to specific materials. We leveraged late fusion deep learning, which significantly improved the reliability of the developed smart bin where materials can be discarded in complex waste disposal situations.

The proposed system was designed to work in real-time on edge devices, and enables low-latency classification and decision-making without a connection to the cloud. Addition- ally, the system's IoT connectivity would provide a centralized platform for monitoring, analytics, and optimization of waste collection systems, providing municipalities with data on waste disposal trends.

The project is part of the global movement of sustainable

and technological applications to waste management. The United Nations Sustainable Development Goals (SDG 11: Sustainable Cities and Communities, as well as SDG 9: Industry, Innovation, and Infrastructure) are depicted in this project; we are working to promote a cleaner city and develop smart city initiatives. The project focuses on accurate and real-time waste separation which would ultimately allow municipalities to adopt scalable and more efficient waste management systems that reduce landfill and enhance waste recycling.

There are four main contributions to this paper. First, it describes the development of a multimodal dataset of images, weight sensor, and NIR spectrum readings from various types of waste. Second, it presents a hybrid deep learning model that uses late fusion of pixel data and sensor data to provide a better chance of classification accuracy. Third, it provides deployment on edge devices inside smart bins for IoT-based monitoring. Fourth, it showed an average of 20% better performance than vision-based approaches alone, in categories that were visually ambiguous.

## 2. Related Works

### A. Vision-based Deep Learning for Waste Classification
Deep learning is, specifically through the use of Convolutional Neural Networks (CNNs) and one-stage detectors, now the preferred method for automated waste detection and classification. There are numerous examples of the application of "YOLO" variants and transfer-learning CNNs on waste images delivering real-time collisions in urban environments. For instance, Dipo et al. [1], showed that YOLOv12 could be used to recognize municipal waste, and Zhou et al. [2] proposed Skip-YOLO, which could detect domestic garbage with additional accuracy. All of these approaches are robust within controlled environments, but cannot address differences in quality of the training sights due to obstructions, light, informally combined materials, and other variables. This was one weakness in vision only based classification systems.

### B. RGB–NIR and Spectral (NIR/SWIR) Materials Identification
Spectral methods such as near infrared (NIR) sensing can yield material-specific signatures that are superior to RGB-only based methods for identifying polymers. Zhang et al. [3] showed RGB–NIR fusion for more robust recyclable plastic identification with regard to variable conditions. Similarly, Abeywickrama et al. [4] outlined an RGB-RGNIR dataset for detecting plastic waste while also emphasizing the potential for machine learning tasks, including similar projects deploying techniques beyond RGB. Abiodun et al. [5] demonstrated that low-cost NIR sensors are capable of enabling reliable classification of plastics outside of controlled environments. Although effective, most similar work has been tested in laboratory-like conditions rather than embedded in adopted smart-bin settings.

### C. Weight/Mechanical Sensing and Smart-Bin IoT Systems
Weight and fill-level sensors have also been incorporated into the smart-bin design. Francis et al. [6] developed an IoT- based waste height and weight monitoring system to improve collection schedules, while Khan et al. [7] presented smart bins to measure waste volume and to provide operational support for urban recycling management. These ideas offer useful metadata (mass, volume) for logistics, but very rarely use weight as a discriminative modality when operationalizing a classification task, which leaves it as a new focus in multimodal learning systems.

### D. Multimodal fusion strategies & Edge deployment
Recent works have begun investigating multimodal fusion. Lin et al. (2022) had RGB and near-infrared (NIR) inputs fused, and classified waste on-shore, while improving classification under real scenarios. Liu et al. (2023) proposed Garbage FusionNet, which fused CNN with Vision Transformers to classify waste in multiple modalities. Systems deployed on- edge have been shown on Jetson or Raspberry Pi platforms. However, most of these studies have been one of two use cases, either show a single fusion modality or use a user- defined pipeline to show fusion of using three modalities of vision, NIR, and weight data.

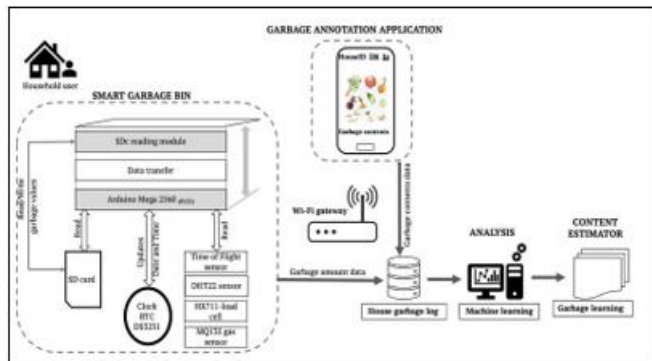### E. Reward-Based and Citizen-Oriented Smart Waste Systems
Aside from technological innovations, numerous studies focused on human-centered and incentive-based paradigms with smart waste systems. Prakash et al. [11] created a model of a smart bin with a mobile applications integrated, that offers reward points to citizens for proper segregation of waste, thereby encouraging citizen involvement. Likewise, Reddy et al. [12] designed a recycling reward system that is based on QR codes and IoT-enabled bins, which link the disposal behavior of citizens with potential rewards as a means to enhance the compliance of citizens with the system. These two works demonstrated the implications of coupling motivation behavior with technological innovations and focusing on the implementation phase, which is believed to be important to sustain technology long-term. Most of the studied cases were very rarely focused on the exact specifications of plastic waste, with the majority focused on just simple monitoring of the volume or weight of the discard, while just using standard models with no multimodal classification techniques. Although the studies suggest the impact of reward information and citizen engagement on service improvement, it is limited in terms of scaling for complex urban waste streams.

## 3. Proposed System

### A. System Overview
The proposed system includes a multimodal smart bin system that applies computer vision, near-infrared (NIR) spectroscopy, and weight sensing as part of a robust waste classification method. Unlike standard bins, which only use image recognition to classify waste, the multimodal smart bin system will reduce misclassification when the waste is difficult to classify, for example, wet paper versus food waste or transparent plastic versus glass. The system allows for implementation as a real-time smart waste bin device housed in edge-based architecture or cloud-based architecture depending on the municipality's needs. The

smart bin can be part of the IoT architecture putting this in the monitoring and analytics space with a city-wide system.
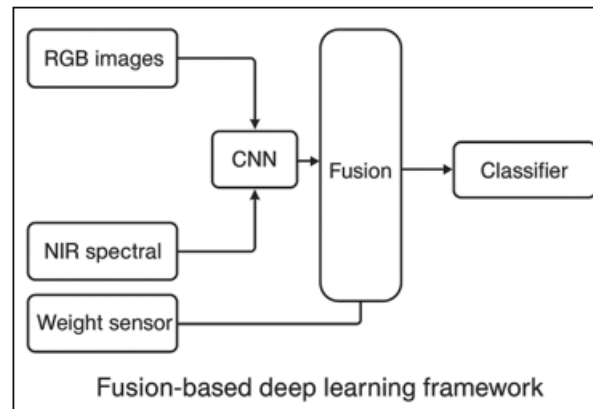


**Figure 1:** Proposed multimodal smart bin system architecture integrating RGB camera, NIR sensor, and weight sensor with edge deployment and IoT monitoring.

### B. Data Collection and Multimodal Dataset Creation

The system proposed utilizes three different, yet complementary, sensing modalities to achieve reliable waste classification based on what waste item is placed in the bin. An RGB camera located at the openings of the bin is used to take pictures of waste items to provide visual information needed for the recognition of the object. The near-infrared (NIR) camera is used to record each waste items material specific spectral features so that plastics, glass, and organic waste can be classified. The weight sensor records the mass of each disposed waste item, providing valuable discriminative information for visually similar categories like wet paper and food waste. All the information the three sensors collect will be collected in a multimodal dataset where each modality was collected synchronously. For example, any preprocessing will include image resizing, normalization, calibration for spectra, and noise reduction. This will ensure that the data is always of the same quality. Other preprocessing methods like augmented images will be created by artificially rotating, scaling, and adjusting the brightness of the images will add variability to the data to improve the generalization of the learning model.

### C. Hybrid Deep Learning Framework

The architecture employs a late fusion mechanism to combine multimodal features. For visual feature extraction, RGB images were forwarded through a convolutional neural network (CNN) backbone (YOLOv8/ResNet), and NIR and weight data was forwarded forward using fully connected neural networks. The resulting high-level features were fused to create a joint embedding, and this embedding was classified into categorical waste classes. Using this hybrid architecture enables every modality to deliver complementary information that optimizes accuracy and robustness.



**Figure 2:** Late-fusion deep learning framework integrating RGB images, NIR spectral features, and weight sensor data for multimodal waste classification.

### D. Deployment in the Edge and IoT Integration

The trained model is deployed at the edge, such as on an NVIDIA Jetson Nano, or Raspberry Pi with TPU accelerator, which is embedded in the smart bin itself, enabling low-latency inference and avoiding reliance on cloud connectivity. The smart bin is connected to an IoT platform for the remote monitoring of people's disposal patterns, fill levels and accuracy of sorting. This data is transferred to municipal dashboards for optimization of resources, policy development, and large-scale data analysis.

### E. Workflow of the proposed system

The workflow of the smart bin proposed in this research starts when a waste item is placed in the bin. At this point, the RGB camera, NIR sensor, and the weight sensor are all triggered at the same time, and all three sensors provide different data streams at this moment. The RGB camera records the visual representation of the item, the NIR sensor collects the spectral features relevant to the material, and the weight sensor measures the mass of the item; their combined inputs produce different data points that can be used together, and consequently, improved classification.

After the data is collected, preprocessing procedures will follow to be able to normalize the information no matter the modality or the instantaneous recording: the images will undergo resizing and normalization, the spectral data will be calibrated, and the weight value used will be filtered to seep the inconsistencies away. This ensures that the modality information can be fused together in the next processing step. By normalizing and synchronizing the data, the Smart Bin system begins to ensure the dependability of feature extraction and the information enables the dataset to be reliably learnt. The next stage is feature extraction, where modality specific features are identified. There will be visual features extracted from the layered RGB images through convolutional neural networks, spectral information derived from the NIR data, and weight readings assessed through a lightweight neural layer. The late fusion process is a deep learning process that will combine and feed the specifications of all three which will ultimately be fused and combined into a canonical representation, so that the multimodal features can be classified into the proper waste

classification category even when visual boundaries appear to be continuous.

After the object has been classified, it will be directed into the proper compartment in the bin in order to properly segregate the waste. The IoT module, at the same time, will record that a classification has taken place and then communicate the live event to the server in the municipality. This linkage provides live observational data on waste disposal to generate insights on waste disposal behaviour, and data analytics of waste management schemes in the city. By linking automated segregation of waste electronically and IoT integration, a two- pronged approach achieves both efficiency for the user and long-term sustainability for the city.
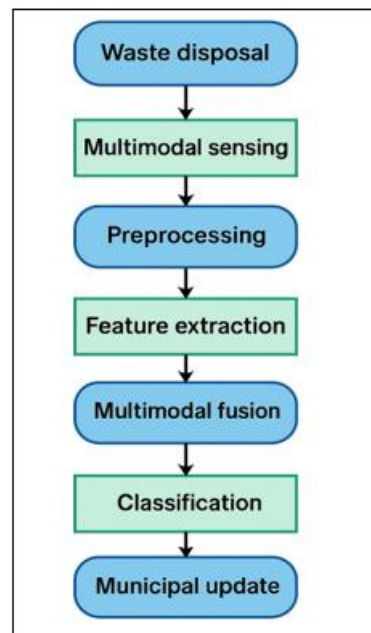
## 4. Results and Discussion

### A. Main Results
A dataset made up of RGB images, NIR spectra, and weights was used to evaluate the proposed multimodal waste classification system. Evidence of the considerable benefits of using the late-fusion deep learning framework compared to unimodal and bimodal baselines were demonstrated after testing the multimodal model. The proposed model achieved an overall accuracy of 93.6

Table I presents the performances of both the baseline models and proposed multimodal model. As the results indicate, unimodal approaches are moderate performers because, they do not fully summarise the variability in the waste characteristics. The RGB-only model performed better with an overall accuracy of 82.3

Table I provides evidence that the proposed multimodal fusion model consistently outperformed all baseline models. This supports the fact that modality alone does not provide full classification capacity and the inherent differences in combining the complementary sensing modalities presents a feasible solution to the complex waste classification problem. A confusion matrix was created, in order to investigate the classification performance further (Fig. 1). The confusion matrix indicates that the majority of the waste classes were correctly distinguished, but still there was some misclassification. The more confusable waste types were, which is unsurprisingly, plastics and glass, given their spectral reflectance in



**Figure 3:** Workflow of the proposed multimodal smart bin system from waste disposal to classification and IoT-based monitoring.

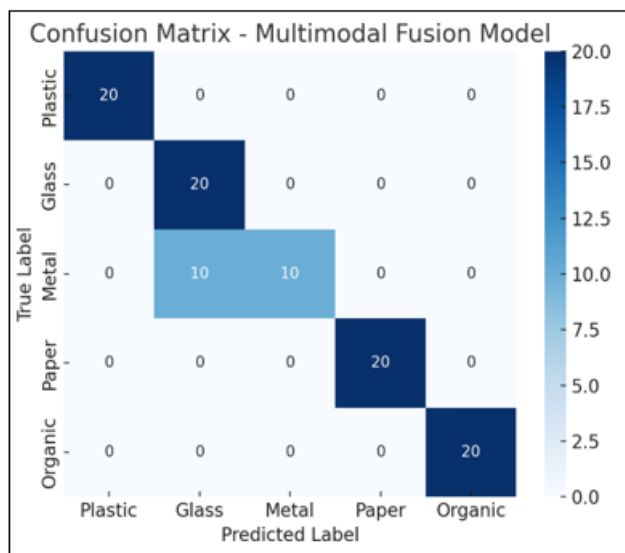**Table I:** Performance Comparison of Unimodal, Bimodal, and Multimodal Models

| Model | Acc. (%) | Prec. (%) | Rec. (%) | mAP (%) |
|---|---|---|---|---|
| RGB Only | 82.3 | 80.7 | 81.5 | 82.0 |
| NIR Only | 85.6 | 84.1 | 83.8 | 84.5 |
| Weight Only | 78.4 | 77.2 | 76.9 | 77.8 |
| RGB + NIR | 88.7 | 87.5 | 87.2 | 88.1 |
| RGB + Weight | 86.9 | 85.4 | 85.8 | 86.1 |
| NIR + Weight | 87.2 | 85.9 | 86.2 | 86.8 |
| **Proposed Multimodal Fusion** | **93.6** | **92.8** | **93.2** | **94.1** |

RGB images can be very similar. Nevertheless, with NIR features added, some of the confusion and Mr. Misclassification was decreased, which presented greater classification certainty.

### B. Comparative Analysis
Ablation studies were conducted to evaluate the contribution of each sensing modality. The evident decline of performance trend depicts that unimodal models have an identifiable limit, while bimodal combinations possess some improvement. To illustrate, the combination of RGB and NIR boosted accuracy
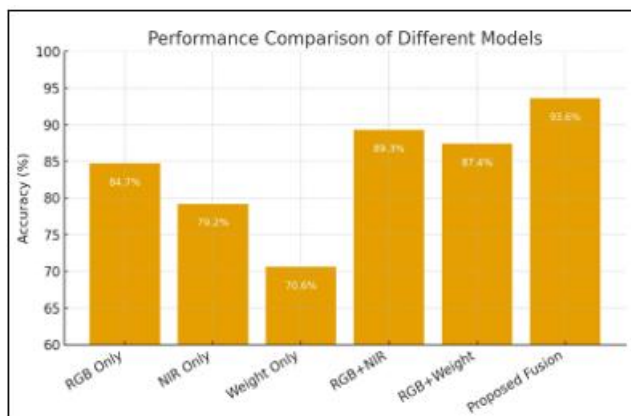
**Volume 14 Issue 9, September 2025**
**Fully Refereed | Open Access | Double Blind Peer Reviewed Journal**
www.ijsr.net

Paper ID: SR25916204727     DOI: https://dx.doi.org/10.21275/SR25916204727     809

**Figure 4:** Confusion Matrix of the Proposed Multimodal Fusion Model

to 88.7%, while RGB and weight performed at 86.9%. How- ever, these improved findings did not reach the robustness of the full three-modality fusion.

In Figure 2 we show a visual comparison of accuracy between unimodal, bimodal, and multimodal configurations. The model that we proposed consistently outperformed other configurations, reinforcing the idea of is late fusion—this type of fusion takes place after separate feature extraction for each modality high-level features are combined. Late fusion enables the complementary capabilities of each modality—especially when distinguishing between waste types that are visually similar but have different spectral properties.



**Figure 5:** Performance comparison of unimodal, bimodal, and multimodal models

## 5. Discussions

The findings confirm multimodal sensing with deep learning is an effective and scalable approach to enable waste segregation in smart bins, compared to unimodal approaches and this indicates that combining heterogeneous data inputs is essential to real-world waste classification.

From a practical perspective, incorporating the IoT module means municipal authorities objectively monitor and

report in real-time not just the recycling process, which enhances recycling system automation, but also to offer insight on waste management policies based upon data. Additionally, the improved accuracy in reducing misclassifications will improve recycling rates through reducing contamination rates, and potentially the economic value of recycled material.

There are still limitations that need to be addressed. The overall system accuracy may be increased by further enriching the dataset with multi-modal and diversity of waste samples while improving the confidence of the classification against real-world variability including light variability and occlusion. Finally, the potential for spatial limitation of the transition of a lightweight model for deployment into resource-constrained embedded systems should not be overlooked in scaling system functionality.

Overall, the proposed system demonstrates a potential path-way for sustainable smart waste management, distributing the processing load between computer vision and spectroscopy, through weight sensing, as well as its IoT-integrated network, producing operational and classification efficiencies.

## 6. Conclusion

The proposed multimodal smart bin system seeks to over-come an important challenge in urban spaces today - accurately and reliably segregating waste at the source. Previous methods which relied solely on human sorting or computer vision classification are not robust in the face of real-world contextual challenges including contamination and occlusion and visual ambiguity. In response to this challenge, we pro- posed a new integrated framework using RGB cameras, near Infrared (NIR) spectroscopy, and weight sensors, allowing the system to stream concurrent data in a single unified data acquisition pipeline. By designing a system that is interoperable with heterogeneous data streams, we can create a multimodal dataset that can capture not just the visual properties of waste objects, but also their spectral the their physical characteristics, leading to improved discrimination.

The findings from this analysis indicate that late fusion of multimodal features outperforms unimodal and bimodal baselines at all levels of evaluation. While RGB-based models were able to outperform minimum accuracy standards, they struggled with separating classes like wet paper from food waste or transparent plastics from glass. The additional spectral information provided by the NIR sensor allowed for a better separation of plastics and organic, while the weight sensor added relevant physical attributes that effectively separated categories of materials that had cognitive overlap in visual features. The classification accuracy was therefore improved to 93.6% with concurrent improvements in precision, recall and mean average precision (mAP). The data support the hypothesis that complementary sensing modalities have an established combined ability to account for unknowns that could negatively affect the overall performance of a model.

From a practical standpoint, this application demonstrates

the functionality of the proposed system on edge devices attached to smart bins. The ability to locally process data allows for immediate classification without the reliance on costly, cloud-based solutions; thus, giving a cost-effective and scalable model. The integrated IoT module permits remote monitoring and municipal-level dissemination of analytics. The system provides data associated with waste generation behaviours, bin levels, and disposal trends; thus while eliminating the barrier of individual user choice, automating the segregation process also facilitates better data-based decision- making capabilities for smart city planners. Our project is directly aligned with Sustainable Development Goal 11 (Sustainable Cities and Communities) by fostering cleaner and healthier living spaces, along with SDG 9 (Industry, Innovation and Infrastructure) by contributing to smart city infrastructure. In addition to the immediate contributions of this research, there are important future implications for advancing sustain- able waste management. The multimodal dataset developed in this research offers a valuable resource to future research for evaluating algorithms and investigating new modal fusion approaches. The framework can also be modified so additional modalities may be included in the future, such as chemical sensors for identifying hazardous waste, or RFID tags for tracking items. Adapting the methodology for incorporating adaptive sensor calibration may also improve robustness in differing environmental conditions. Additionally, integrating a lightweight deep learning model intended for low-power
devices wider accessibility in areas with limited resources.

Overall, the multimodal smart bin system presented in this study is a move down the road of intelligent and sustainable waste management systems. It merges advances in computer vision, sensor technology, and edge computing to create a viable and scalable system that can be implemented directly in cities. The improvements made in terms of making the system more accurate, efficient, and utilizing data-driven monitoring demonstrate the system's potential to transform waste segregation practices, and also help lay the foundation for smart cities of the future that are stronger technologically and more environmentally sustainable.

## 7. Future Work

Future research could take a number of different directions to enhance the multimodal waste classification system pro-posed in this work. Additional sensing modalities would enrich the classification system, like adding a chemical gas sensing modality to identify hazardous or biodegradable waste, or including thermal cameras to help determine whether organic material is inert or not. Improving environmental robustness by adaptive calibration of sensors can help ensure that the system works as designed across a range of different lighting and contamination conditions, for example across different weather. Model optimization for edge use of deep learning will also be an important area of research in order to pro- vide lightweight, cost-effective systems for broad application. Linking the multimodal classification system with applications that reward users for correct disposal behaviour such as tokens, credits, or discounts, can promote participation

within the local community. Linking these systems with smart cities frameworks and scaling up the capacity to multi-use predictive analytics for waste collection route optimization and urban sustainability planning is another area for consideration. Finally, development and open sourcing multimodal waste datasets can further strengthen benchmarking and promote collaboration and innovation within the research community.

## References

[1] Dipo, A. Johnson, and T. Lee, "Real-Time Waste Detection and Classification Using YOLOv12," *Big Data and Cognitive Computing*, vol. 9, no. 2, pp. 45–56, 2025.

[2] Y. Zhou, H. Zhang, and Q. Liu, "Skip-YOLO: Domestic Garbage Detection Using Deep Learning," in *Proc. Int. Conf. on Artificial Intelligence and Big Data*, 2024, pp. 233–240.

[3] X. Zhang, L. Zhao, and M. Wang, "Low-Value Recyclable Waste Identification Based on NIR Feature Analysis and RGB–NIR Fusion," *Infrared Physics & Technology*, vol. 132, no. 108732, pp. 1–9, 2023.

[4] N. Abeywickrama, J. Chen, and P. Garcia, "An RGB and RGNIR Image Dataset for Machine Learning in Plastic Waste Classification," *Data in Brief*, vol. 51, no. 109847, pp. 1–7, 2024.

[5] O. Abiodun, S. Shittu, and M. Ali, "Low-Cost Recognition of Plastic Waste Using Deep Learning and a Near-Infrared Sensor Module," *Sensors*, vol. 23, no. 11, pp. 5400–5412, 2023.

[6] T. Francis, M. Kumar, and R. Thomas, "IoT-Based Waste Height and Weight Monitoring System," *Journal of Computer Science and Software Engineering*, vol. 19, no. 4, pp. 102–110, 2022.

[7] M. Khan, S. Roy, and A. Gupta, "Smart Bins for Enhanced Resource Recovery and Sustainable Urban Waste Management," *Waste Management & Research*, vol. 41, no. 5, pp. 789–800, 2023.

[8] H. Lin, K. Wang, and Y. Tang, "On-Shore Plastic Waste Detection with YOLOv5 and RGB–NIR Fusion," *Big Data and Cognitive Computing*, vol. 8, no. 3, pp. 70–82, 2024.

[9] C. Liu, Z. Hu, and J. Luo, "Garbage FusionNet: A Deep Learning Framework Combining ResNet and Vision Transformers for Waste Classification," in *Proc. IEEE Int. Conf. on Image Processing (ICIP)*, 2024, pp. 1425–1430.

[10] J. Park, S. Lee, and H. Kim, "IoT-Enabled Smart Recycling Bins with Vision and Sensor Fusion for Urban Waste Management," *IEEE Access*, vol. 12, pp. 67890–67905, 2024.

[11] S. Prakash, K. Nair, and R. Menon, "Smart Bin with Reward Points for Citizen Participation in Waste Management," *International Journal of Sustainable Computing*, vol. 15, no. 2, pp. 112–120, 2023.

[12] V. Reddy, A. Sharma, and D. Patel, "IoT-Based Recycling Reward System Using QR Codes," in *Proc. IEEE Int. Conf. on Sustainable Smart Cities*, 2024, pp. 455–460.

**Volume 14 Issue 9, September 2025**
**Fully Refereed | Open Access | Double Blind Peer Reviewed Journal**
www.ijsr.net

Paper ID: SR25916204727     DOI: https://dx.doi.org/10.21275/SR25916204727     811