

Edge AI Agent Design for Policy-Aware Urban Waste Management

Binrong Zhu, Ruxue Jin, Yang Liu, Guiran Liu, Qun Wang, and Phuong Mai Nguyen

Computer Science Department, San Francisco State University, San Francisco, CA 94132
{bzhu2, rjin, yliu68, gliu, qunwang, pnguyen27}@sfsu.edu

Abstract

Urban waste management confronts a crisis of escalating volume, unsustainable costs, and environmental degradation. While AI has shown promise in waste classification, existing systems are predominantly static “classifiers” that operate in isolation, unaware of the dynamic, real-world context in which they are deployed. This paper introduces a novel framework for an Edge AI agent that transforms smart bins from passive classifiers into active, context-aware managers capable of dynamic, policy-driven decision-making at the point of disposal. The agent integrates a dual-YOLO perception module, a dynamic policy database, IoT-driven state awareness, and a cognitive core powered by a quantized, on-device Large Language Model (LLM) (ODLLM). By leveraging a Retrieval-Augmented Generation (RAG) pipeline to ground its reasoning in municipal guidelines, the agent provides context-specific, actionable guidance to reduce contamination, optimize logistics, and enhance public engagement. This approach offers a scalable solution to advance the goals of the circular economy and build more sustainable urban futures.

1 Introduction

The escalating crisis of urban waste is not merely a logistical challenge but a complex information problem demanding intelligent, decentralized solutions. Global municipal solid waste (MSW) is projected to grow from 2.1 billion tonnes in 2023 to 3.8 billion by 2050. This explosion carries a staggering economic burden, with the true annual cost, including environmental externalities, projected to reach \$640.3 billion by 2050 without significant intervention [1]. Environmentally, landfills are a primary source of methane emissions, and mismanaged MSW is the origin of an estimated 80% of marine plastic pollution [1].

Current waste management paradigms are ill-equipped for this challenge. Manual sorting is error-prone and unscalable, leading to high contamination rates in municipal recycling programs [4]. First-generation AI systems, while promising, are typically static “classifiers” that lack awareness of local disposal policies or the real-time status of waste bins. They often suffer from “catastrophic forgetting,” where fine-tuning on specialized waste datasets degrades their ability to recognize general objects, making them brittle in diverse urban environments.

This research proposes a paradigm shift from passive classification to active, context-aware management. We introduce an “Edge AI Agent” designed to be “situated”—meaning it reasons

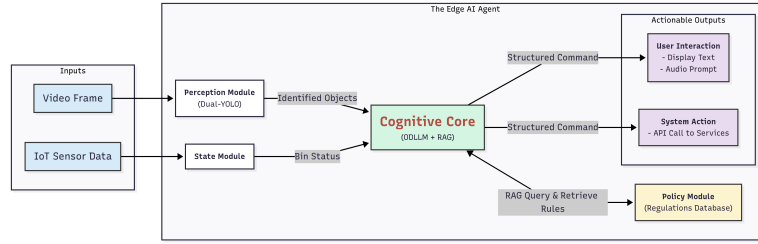


Figure 1: System Architecture of the Policy-Aware Edge AI Agent

over real-time perception, dynamic local policies, and the live physical state of its environment [3]. By integrating perception (what is this item?), policy (what are the rules?), and state (is the bin full?), the agent closes critical information gaps to provide actionable guidance at the point of disposal [2].

2 System Model

To achieve situated intelligence, the agent employs a modular, multi-layered architecture where each component is responsible for a distinct function. This design enhances robustness and maintainability, allowing for updates (e.g., to policy) without retraining the entire system. The four key modules are deployed on a resource-constrained edge device as shown in figure 1, transforming a standard bin into an intelligent node in the urban sanitation network.

Perception Module: A Dual-YOLO Framework The agent’s ”eyes” are a hybrid dual-model perception system designed to overcome the limitations of single-model classifiers. This architecture structurally separates specialized and generalized knowledge to prevent catastrophic forgetting. (1) **Specialized Model:** A YOLOv5 model is fine-tuned on the TACO (Trash Annotations in Context) dataset, which contains images of waste in diverse, real-world settings. This provides deep, task-specific expertise in identifying common waste items. (2) **Generalized Model:** A YOLOv8 model pre-trained on the broad COCO (Common Objects in Context) dataset provides robust detection for a wide array of everyday objects (e.g., ’laptop’, ’backpack’) that may not be in the TACO dataset but frequently appear in waste streams [6]. Outputs from both models are harmonized into a unified four-class system (Recycle, Compost, Landfill, No Bin/Hazardous) based on municipal guidelines. Redundant detections are fused using a Non-Maximum Suppression (NMS) algorithm, ensuring a clean, accurate input for the agent’s cognitive core.

Beyond perception, the agent integrates three modules to achieve true context-awareness and decision-making capabilities. **Policy Module:** This is a dynamic, queryable database containing local disposal guidelines, such as San Francisco’s Recology sorting requirements[4]. It decouples operational logic from the AI models, allowing municipal rules to be updated without retraining the system. The module can store nuanced rules like ”greasy pizza boxes must be composted, not recycled”.

State Module: This IoT-driven component provides real-time awareness of the physical waste infrastructure. It receives data from sensors (ultrasonic, infrared) in nearby bins, reporting crucial information like fill levels and operational status (e.g., ”full,” ”available”). This live data allows the agent to make decisions based on the practical, real-time capacity of the system.

Interaction Output Module: This is the agent’s cognitive core, powered by a lightweight, quantized Large Language Model (LLM) running on the edge device. Running the LLM locally

ensures low latency and data privacy [7] [8]. This "brain" fuses the data from the other three modules to reason over the complete context and generate an actionable output. The framework demonstrates practical on-device performance and adaptability to diverse urban waste management policies.

3 On-Device Reasoning with Retrieval-Augmented Generation

The agent's cognitive core uses an LLM to transform multimodal data into intelligent guidance, moving beyond rigid if-then logic [5]. This is achieved through a structured input process and a Retrieval-Augmented Generation (RAG) pipeline that grounds the LLM's reasoning in authoritative sources.

3.1 Multimodal Input and Domain-Informed Prompting

The LLM's reasoning begins with a structured JSON input that consolidates information from the other modules, providing a complete picture of the disposal scenario. This input includes the detected items, the status of nearby bins, and the relevant local policy query result [8]. The LLM is guided by a domain-informed prompt that establishes its goals and priorities, such as instructing it that preventing bin overflow takes precedence over a default sorting rule. This initial instruction sets the stage for the RAG pipeline to provide specific, timely information.

3.2 RAG for Policy Literacy

Municipal waste policies are often complex, text-based documents that are impractical to hard-code. The RAG pipeline makes the agent "policy-literate" by allowing it to dynamically interpret these regulations. When an object is identified, the agent queries a local vector database containing the full text of municipal waste manuals and guidelines. The RAG system retrieves the most relevant text snippets—for example, the specific rule for "greasy pizza boxes"—and appends them to the LLM's prompt. This mechanism ensures the LLM's response is grounded in the latest, verifiable sources, effectively allowing it to "read the manual" for each disposal event without requiring retraining.

3.3 Actionable Output Generation

After reasoning over the RAG-enhanced input, the LLM generates a structured JSON output that includes not just a classification, but a command that can trigger actions. For instance, if a user presents a recyclable bottle but the recycling bin is full, the agent's output can include: (1) A user-facing message: A clear instruction displayed on a screen or spoken via audio, such as, "The recycling bin is currently full. Please place this item in the landfill bin for now". (2) A system-level action: An automated command to the municipal waste service to request collection for the full bin, specifying the location and reason.

This closes the loop from perception to action, providing immediate guidance to the user while simultaneously initiating a back-end process to resolve the infrastructure issue.

4 Performance, Impact, and Future Directions

4.1 System Performance

The dual-model perception module was evaluated empirically. After training the YOLOv5 model for 50 epochs on the TACO dataset, key performance metrics were achieved as shown in figure 2. Performance testing revealed the combined approach reduced false negatives by 18% compared to either model used individually, directly addressing the recall limitations observed in single-model implementations.

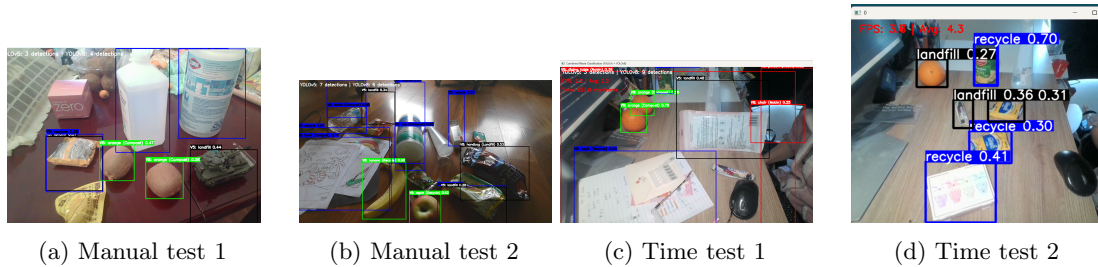


Figure 2: Overall Detection by dual-model

While the dual-model architecture has a higher computational overhead, reducing the frames per second (FPS) from 4.3 to 2.2, this trade-off is justified by the significant gain in detection diversity and robustness against catastrophic forgetting. Moreover, a validation framework using a simulated "digital twin" of an urban environment is proposed to test the decision-making accuracy of the complete agent under dynamic conditions. To validate the agent's decision-making, the system was tested on 40 waste detection cases across 10 simulated bin stations configured with San Francisco's Recology rules, demonstrating 100% policy compliance in all recommendations. Furthermore, the agent showcased robust context-awareness, adapting its output based on bin status in 15% of cases by correctly handling bin-full events and redirecting hazardous waste. The implementation is available at : <https://github.com/Binrongz/urban-edge-ai-agent>.

5 Conclusion

This research presents a paradigm shift from static waste classification to active, situated agency, establishing a comprehensive framework for a policy-aware Edge AI agent. The deployment of such a system offers notable social and economic benefits by reducing recycling contamination and optimizing collection logistics. It also promotes public engagement and sustainable habits through real-time, context-aware guidance. Future work will enhance adaptability via continual learning, model optimization, and multi-modal sensing to support large-scale deployment and alignment with evolving policies. These efforts aim to advance intelligent, resilient, and sustainable urban waste management infrastructures.

References

- [1] Priscilla Cristine Porto Leó Costa, Fábio Pinto Cardoso, and Maria Conceição Melo Silva Luft. Co-production and circularity: Integrating emerging international practices. *Revista Pensamento*

- Contemporâneo em Administração*, 19(2):19–42, 2025.
- [2] Lifu Gao, Joshua Sherwood, Nawwaf Aleisa, Andrews Damoah, Yingzhou Lu, and Xiaodong Qu. Human-centered ai agents for healthcare and education: A systematic literature review.
 - [3] Lincan Li, Jiaqi Li, Catherine Chen, Fred Gui, Hongjia Yang, Chenxiao Yu, Zhengguang Wang, Jianing Cai, Junlong Aaron Zhou, Bolin Shen, et al. Political-llm: Large language models in political science. *arXiv preprint arXiv:2412.06864*, 2024.
 - [4] SF Environment. Zero waste case study: San francisco. U.S. Environmental Protection Agency, 2022. Retrieved from <https://www.epa.gov/transforming-waste-tool/case-study-san-francisco>.
 - [5] Satvik Verma, Qun Wang, and E Wes Bethel. Intelligent iot attack detection design via odllm with feature ranking-based knowledge base. In *Proceedings of the AAAI Symposium Series*, volume 5, pages 188–195, 2025.
 - [6] M. Guru Vimal Kumar, Madde Kumar, K Narayana Rao, P Syamala Rao, Arepalli Tirumala, and Eswar Patnala. Advanced yolo-based trash classification and recycling assistant for enhanced waste management and sustainability. In *2024 Second International Conference on Intelligent Cyber Physical Systems and Internet of Things (ICoICI)*, pages 1238–1246, 2024.
 - [7] Yue Wang, Tianfan Fu, Yinlong Xu, Zihan Ma, Hongxia Xu, Bang Du, Yingzhou Lu, Honghao Gao, Jian Wu, and Jintai Chen. Twin-gpt: Digital twins for clinical trials via large language model. *ACM Trans. Multimedia Comput. Commun. Appl.*, July 2024. Just Accepted.
 - [8] Jiajun Xu, Zhiyuan Li, Wei Chen, Qun Wang, Xin Gao, Qi Cai, and Ziyuan Ling. On-device language models: A comprehensive review. *arXiv preprint arXiv:2409.00088*, 2024.