

Enhancing Graph Representation of the Environment through Local and Cloud Computation

Francesco Argenziano, Vincenzo Suriani, and Daniele Nardi

Dept. of Computer, Control, and Management Engineering, Sapienza University of Rome, Italy.

{lastname}@diag.uniroma1.it

Abstract—Enriching the robot representation of the operational environment is a challenging task that aims at bridging the gap between low-level sensor readings and high-level semantic understanding. Having a rich representation often requires computationally demanding architectures and pure point cloud based detection systems that struggle when dealing with everyday objects that have to be handled by the robot. To overcome these issues, we propose a graph-based representation that addresses this gap by providing a semantic representation of robot environments from multiple sources. In fact, to acquire information from the environment, the framework combines classical computer vision tools with modern computer vision cloud services, ensuring computational feasibility on onboard hardware. By incorporating an ontology hierarchy with over 800 object classes, the framework achieves cross-domain adaptability, eliminating the need for environment-specific tools. The proposed approach allows us to handle also small objects and integrate them into the semantic representation of the environment. The approach is implemented in the Robot Operating System (ROS) using the RViz visualizer for environment representation. This work is a first step towards the development of a general-purpose framework, to facilitate intuitive interaction and navigation across different domains.

I. INTRODUCTION

In recent years, the field of robotics has witnessed significant advancements in perception capabilities, thanks to the proliferation of sensors and computer vision techniques. However, bridging the gap between low-level sensor readings and high-level semantic understanding remains a challenge. To this end, we propose a framework, in its early stages, that tackles this gap using a graph representation to connect sensor data with a semantic representation of the environment.

To be able to have good precision in object detection and achieve computationally acceptable performance on resource-constrained robotic hardware, our approach combines classical computer vision tools with modern computer vision cloud services. By leveraging the power of cloud computing, we can offload intensive processing tasks and ensure real-time responsiveness even on limited onboard hardware. Since robots often need to deal with small objects in the environment, we adopted the cloud-vision system to have high precision on single objects in the environment.

One of the major advantages of our framework lies in its ability to be cross-domain and adaptable to different environments without requiring the development of specific customizations for each scenario. This is achieved through the incorporation of an ontology hierarchy, encompassing more than 800 object classes. We integrated the remote hierarchy

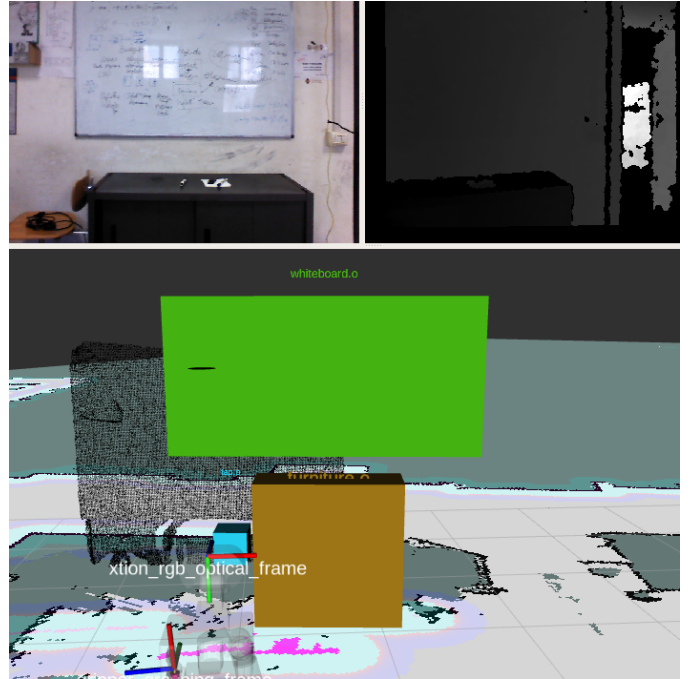


Fig. 1. The obtained 3D scene representation in RViz with the exploration of a small portion of the environment. The objects are delimited by the bounding boxes (obtained from the cloud computation) fused with the point cloud information. After this, the corresponding graph representation is generated.

with the local ontology to obtain a unified graph representation that can be used in the robot's tasks. By using such a comprehensive ontology, our framework can handle diverse environments and objects, facilitating seamless navigation and interaction across various contexts.

To semantically represent the robot's environment, we employ a set of entity classes to define objects and their attributes, while a graph representation is utilized to establish connections between the environment's entities. We implement the proposed architecture on the Robot Operating System (ROS) and use the RViz visualizer, allowing for intuitive visualization and interaction with the environment. An example of this visualization can be seen in Fig. 1. The platform used is the TIAGo robot, manufactured by PAL Robotics¹.

The rest of this paper is organized as follows: Section 2

¹<https://pal-robotics.com/>

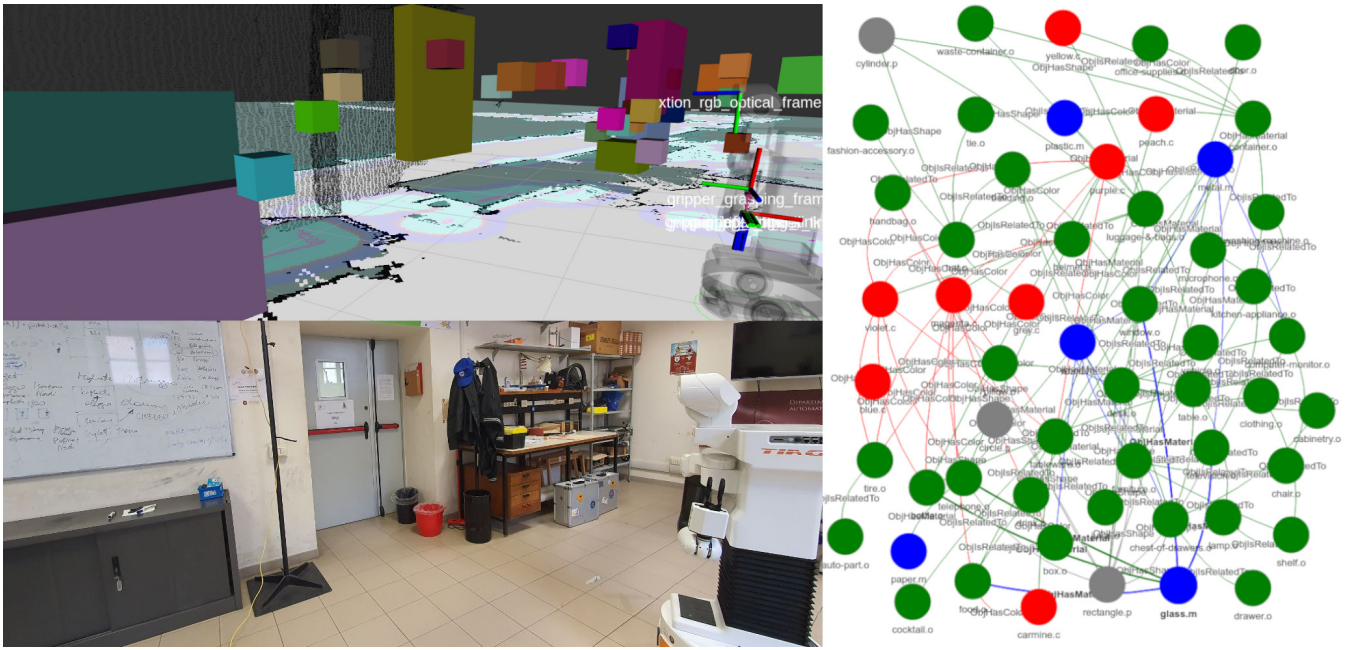


Fig. 2. On the left: population in the RViz visualizer of the objects detected thanks to Google Cloud’s API. On the right, the graph extracted from the underlying relations between objects, and between objects and their properties. Different colors mean different semantic groups of the nodes: materials, shapes, objects and colors.

provides a brief overview of related work in semantic representation in robotics. Section 3 elaborates on our proposed framework, highlighting the graph-based representation and ontology hierarchy, focusing on the integration with ROS and the chosen platform. Section 4 presents concluding remarks.

II. RELATED WORK

Semantic representation in robotics is a vital area of research, enabling robots to comprehend and interact with their environment effectively. Various approaches have been proposed to bridge the gap between sensor data and semantic understanding. Object recognition algorithms and scene understanding techniques [2] are commonly employed to extract high-level semantic information from sensor data, facilitating intelligent decision-making by robots.

In last years, graph-based approaches have gained popularity in robotics due to their ability to capture complex relationships and dependencies within the environment. By representing the environment as a graph, these frameworks provide a structured representation that facilitates semantic understanding and reasoning. Graph-based frameworks have been successfully applied to object detection [10] and scene parsing [6], enabling robots to perceive and interpret their surroundings effectively. 3D scene graphs have been presented and used in [1], [5] to represent 3D environments. In those, nodes represent spatial concepts at multiple levels of abstraction and edges represent relations between concepts. One of the limitations in building such a representation automatically has been represented by the computational costs. To overcome such issues, recently, in [4] the authors introduced Hydra, capable to build incrementally a 3D scene graph from sensor

data in real time thanks to the combination of novel online algorithms and a highly parallelized perception architecture. Another approach, capable of incrementally building the scene graph but also aggregating PointNet[8] features from primitive scene components using graph neural network has been proposed in [9]. In this, an attention mechanism has also been proposed to deal with missing graph data in incremental reconstruction scenarios. When dealing with local approaches, the set of objects that are detectable is quite limited due to the limited availability in hosting large neural network models. To this end, in order to guarantee cross-domain adaptability through a large set of detectable objects and computational sustainability on robot CPUs, cloud services have been adopted as an addition to the local perception pipeline. To address this challenge, cloud services for computer vision, such as Google Cloud Vision², have emerged as viable solutions in robotic applications [3]. By leveraging cloud services, robots can offload processing tasks, enabling real-time perception even on resource-constrained hardware [7]. By relying on this platform we fused the remote hierarchy with the local ontology to map the environment obtaining a unified graph representation that can be used for robot navigation and object localization tasks.

III. METHODOLOGY

A. Entity Representation and Hierarchy

Dealing with hierarchies of concepts with too many entries can be very difficult to handle, especially when cross-domain applications are involved. Hierarchical ontologies that include

²cloud.google.com/vision

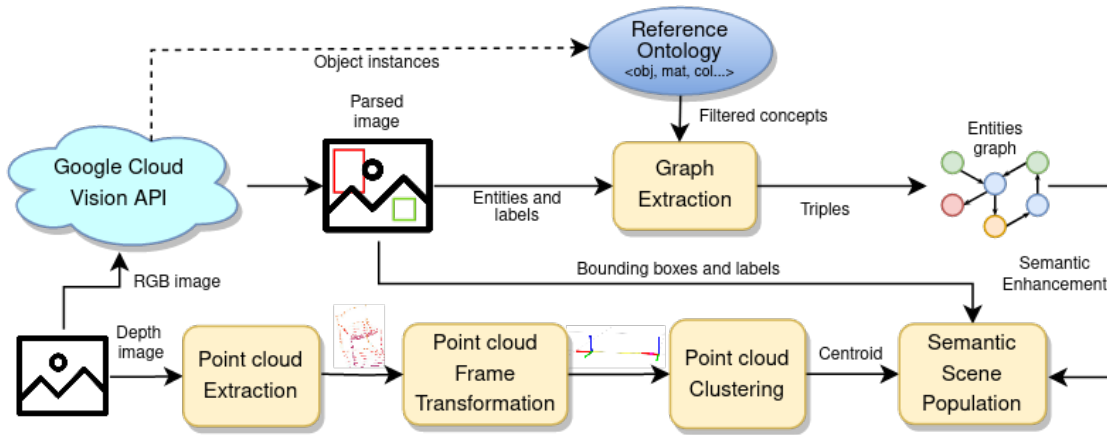


Fig. 4. The pipeline presented in this work. The local and remote perception branches lead to a unified graph representation of the environment.

approach that comprehends the full set of objects of the Google Cloud Vision detector, we integrate an ontology hierarchy, consisting of 800+ object classes. This has enabled cross-domain adaptability and eliminated the need for developing specific tools for different environments.

This is a first step in the implementation of a blended system, since, for example, there is plenty of room to improve the mapping between bounding boxes classified from the cloud service and the mapping on the point cloud depth. In fact, the majority of the issues come when dealing with partial occlusions in the detected objects. To this end, we aim to improve the integration between the local representation and the cloud-computed one, taking into account not only clusters of point clouds but preserving the shapes obtained from them. We also plan to explore the integration of machine learning techniques to improve the accuracy and robustness of semantic understanding in dynamic environments.

Overall, the proposed approach opens up avenues for enhanced perception and interaction capabilities of robots in various domains. By leveraging the power of semantic representation, we envision a future where robots can understand and navigate complex environments, enabling seamless human-robot interaction and collaboration in diverse settings.

ACKNOWLEDGMENTS

We acknowledge partial financial support from PNRR MUR project PE0000013-FAIR. This work has been carried out while Francesco Argenziano was enrolled in the Italian National Doctorate on Artificial Intelligence run by Sapienza University of Rome.

REFERENCES

- [1] Iro Armeni, Zhi-Yang He, JunYoung Gwak, Amir R. Zamir, Martin Fischer, Jitendra Malik, and Silvio Savarese. 3d scene graph: A structure for unified semantics, 3d space, and camera. *CoRR*, abs/1910.02527, 2019. URL <http://arxiv.org/abs/1910.02527>.
- [2] Saurabh Gupta, Ross Girshick, Pablo Arbeláez, and Jitendra Malik. Learning rich features from rgb-d images for object detection and segmentation, 2014.
- [3] Jhih-Yuan Huang and Wei-Po Lee. Enabling vision-based services with a cloud robotic system. In *2016 Asia-Pacific Conference on Intelligent Robot Systems (ACIRS)*, pages 84–88, 2016. doi: 10.1109/ACIRS.2016.7556193.
- [4] Nathan Hughes, Yun Chang, and Luca Carlone. Hydra: A real-time spatial perception engine for 3d scene graph construction and optimization. *CoRR*, abs/2201.13360, 2022. URL <https://arxiv.org/abs/2201.13360>.
- [5] Ue-Hwan Kim, Jin-Man Park, Taek-Jin Song, and Jong-Hwan Kim. 3-d scene graph: A sparse and semantic representation of physical environments for intelligent agents. *CoRR*, abs/1908.04929, 2019. URL <http://arxiv.org/abs/1908.04929>.
- [6] Xia Li, Zhisheng Zhong, Jianlong Wu, Yibo Yang, Zhouchen Lin, and Hong Liu. Expectation-maximization attention networks for semantic segmentation. *CoRR*, abs/1907.13426, 2019. URL <http://arxiv.org/abs/1907.13426>.
- [7] Ricardo C Mello, Moises RN Ribeiro, and Anselmo Frizzera-Neto. Introduction to cloud robotics. In *Implementing Cloud Robotics for Practical Applications: From Human-Robot Interaction to Autonomous Navigation*, pages 1–11. Springer, 2022.
- [8] Charles Ruizhongtai Qi, Hao Su, Kaichun Mo, and Leonidas J. Guibas. Pointnet: Deep learning on point sets for 3d classification and segmentation. *CoRR*, abs/1612.00593, 2016. URL <http://arxiv.org/abs/1612.00593>.
- [9] Shun-Cheng Wu, Johanna Wald, Keisuke Tateno, Nasir Navab, and Federico Tombari. Scenegrphfusion: Incremental 3d scene graph prediction from RGB-D sequences. *CoRR*, abs/2103.14898, 2021. URL <https://arxiv.org/abs/2103.14898>.
- [10] Shifeng Zhang, Longyin Wen, Xiao Bian, Zhen Lei, and Stan Z. Li. Single-shot refinement neural network for object detection. In *European Conference on Computer Vision (ECCV)*, pages 174–189, 2018.