

# Extending Maps with Semantic and Contextual Object Information for Robot Navigation: a Learning-Based Framework using Visual and Depth Cues

Renato Martins<sup>1,2</sup>, Dhiego Bersan<sup>1</sup>, Mario F. M. Campos<sup>1</sup>, and Erickson R. Nascimento<sup>1</sup>

<sup>1</sup>Universidade Federal de Minas Gerais (UFMG), Brazil    <sup>2</sup>INRIA, France

**Abstract.** This paper addresses the problem of building augmented metric representations of scenes with semantic information from RGB-D images. We propose a complete framework to create an enhanced map representation of the environment with object-level information to be used in several applications such as human-robot interaction, assistive robotics, visual navigation, or in manipulation tasks. Our formulation leverages a CNN-based object detector (Yolo) with a 3D model-based segmentation technique to perform instance semantic segmentation, and to localize, identify, and track different classes of objects in the scene. The tracking and positioning of semantic classes is done with a dictionary of Kalman filters in order to combine sensor measurements over time and then providing more accurate maps. The formulation is designed to identify and to disregard dynamic objects in order to obtain a medium-term invariant map representation. The proposed method was evaluated with collected and publicly available RGB-D data sequences acquired in different indoor scenes. Experimental results show the potential of the technique to produce augmented semantic maps containing several objects (notably doors). We also provide to the community a dataset composed of annotated object classes (doors, fire extinguishers, benches, water fountains) and their positioning, as well as the source code as ROS packages. <sup>1</sup>

## 1 INTRODUCTION

Scene understanding is a crucial factor for the deployment of intelligent agents in real-world scenes in order to perform and support humans in everyday tasks [1]. We have recently witnessed significant advances in the fields of scene understanding, human-robot interaction and mobile robotics, but they are still often challenged by typical adversities found in real environments. Surprisingly, these

---

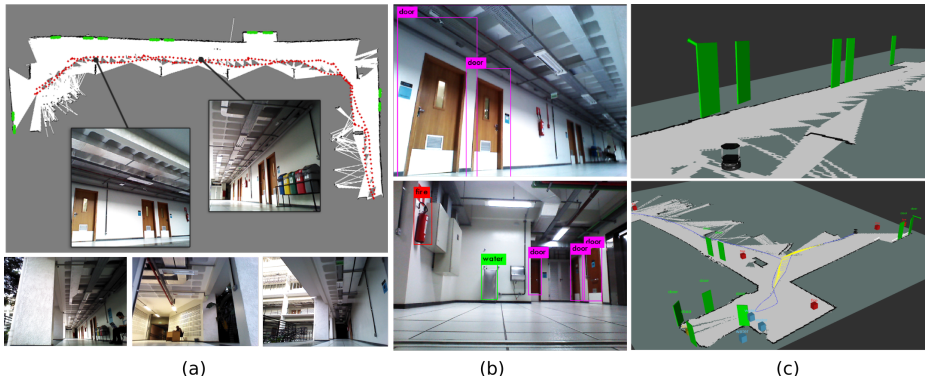
<sup>1</sup> Preprint paper version to appear at Journal of Intelligent & Robotic Systems, available online at: <https://doi.org/10.1007/s10846-019-01136-5>

adversities are always (and successfully) handled daily by humans using mostly vision as primary sense. From their tender age humans learn to recognize and to build more abstract representations of what they observe in their surroundings: we look at with the eyes, but we see with the brain.

In this context, embedding a higher level of scene understanding to identify particular objects of interest (including people), as well as to localize them, would greatly benefit intelligent agents to perform effective visual navigation, perception and manipulation tasks. Notably, this is a desired capability in human-robot interaction or in autonomous robot navigation tasks in daily-life scenes since it can provide “situation awareness” by distinguishing dynamic entities (e.g., humans, vehicles) from static ones (e.g., door, bench) [2, 3], or to recognize unsafe situations. This competence is also instrumental in the development of personal assistant robots, which need to deal with different objects of interest for guiding visually impaired people to cross a door, to find a bench, or a water fountain. Moreover, this higher level representation can provide awareness of dangerous situations (such as the presence ahead of steps, stairs) and of other people for safe navigation and interaction [4, 5]. Recent advances of data-driven machine learning techniques and the increased computing capability of daily-use electronic devices have allowed envisaging transferring, to artificial agents, the human skills required to build these higher-level representations. It is then desirable to integrate these advances notably to mobile robotic systems, allowing them to perform more complex tasks, in safer conditions and in less specialized environments.

In this paper, we propose and evaluate a learning-based framework using visual and depth cues for building semantic augmented metric maps. The resulting representation combines both environment structure, appearance (metric map) and semantics (objects classes). The presented approach detects and generates models of objects in the surrounding environment using an RGB-D camera (or any stereo camera rig such as ZED 3D stereo cameras) as primary sensorial input. In a first moment, these RGB-D images are processed by a convolutional neural network to extract object classes as higher-level information, which is then leveraged by a localization and tracking system of the object instances over time. Finally, the environment representation is extended with the semantic information extracted using the object categories. In order to allow easy deployment in different robotic platforms, the full system is integrated in ROS (Robot Operating System). Moreover, we also provide a dataset acquired in indoor environments with corridors and offices, containing annotated objects positions to help assessing and evaluating 3D object detection and mapping techniques. A characteristic result of our framework is depicted in Figure 1, containing some image frames from one data sequence of the proposed dataset, as well as the respective object detections and augmented maps.

A preliminary conference version paper is introduced in our previous [6]. In this manuscript, we have made a number of major modifications that we summarize as follows:



**Fig. 1.** Augmented semantic mapping overview. (a) Bird’s-eye view of the 2D map and door locations (in green) and some image frames of the dataset; (b) Object detection examples; and (c) Visualizations of the augmented semantic map output.

- The localization and object tracking of the classes are improved to handle multiple objects per image and to support online pose updates, during loop closing of the localization and metric mapping back-end.
- The object extraction and positioning are extended to support object classes with more complex shapes than planar patches.
- Additional experiments and results are presented with object classes beyond “door”, such as “person”, “bench”, “trash bin”, among others.
- We present the training strategy and protocols used with the neural network for taking into account custom object classes. We also describe and provide the code for performing data augmentation and labeling, given a small sample of images from the additional object classes.
- We also provide the image training samples, source code, dataset sequences and video demos of the project <sup>2</sup>.

The rest of this paper is structured as follows. In Section 2, we discuss some recent related work on semantic object information and augmented map representations. Section 3 presents the main stages of our semantic map augmentation. Then, we describe in Section 4 the experimental setup, implementation details and the obtained results using real image sequences. The proposed dataset, that includes three data sequences collected in indoor scenes, is introduced in Section 4.1. Finally, Section 5 presents concluding remarks and discusses some perspectives of the work.

## 2 RELATED WORK

There has been a great interest from the computer vision and robotics communities to exploit object-level information since from the perspective of many

<sup>2</sup> <https://www.verlab.dcc.ufmg.br/semantic-mapping-for-robotics/>

applications, it is beneficial to explore the awareness that object instances can provide for assistive computer vision [7, 8, 9], tracking/SLAM [10, 11], or place categorization/scene recognition and life-long mapping [12, 13].

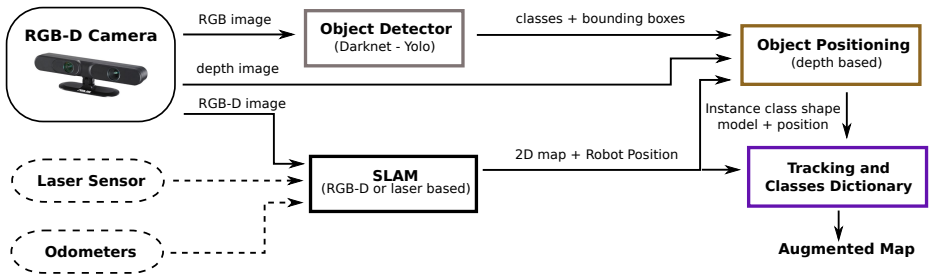
## 2.1 Object Detection and Segmentation

In order to build our extended map representation, we perform “instance semantic” segmentation of objects leveraging an object detector with geometric priors. Object-level representations are, in general, gathered from solving the challenging problems of object detection and semantic segmentation/labeling. An extensive amount of work have previously been reported to tackle these problems, employing a plethora of formulations ranging from graph-cuts, belief-propagation, or convex relaxation optimization/variational optimization, to name a few (the reader is referred to the survey [14]). However, the majority of recent state-of-the-art techniques are grounded on neural networks [15, 16, 17, 18].

Most recent object detection techniques are based on the generation of image region proposals, i.e., bounding boxes, and then predicting the most likely object class for each region. Commonly used benchmarks to evaluate object detection algorithms are the PASCAL Visual Object Classes (VOC) datasets [19], ImageNet [20] for object detection or the Multiple Object Tracking benchmarks (MOT) conceived specially to the evaluation of detection of humans in video. On the other hand, semantic segmentation is often done in the level of pixels and outputs different object classes in the image, but without object instance level information. Recent works as Mask-RCNN [21] and YOLACT [22] perform “instance semantic” segmentation by combining several nets to simultaneously detect object instances and their semantic segmentation. The bottleneck of aforementioned approaches adopting supervised semantic segmentation is the user effort required to annotate pixel-wise a large number of images containing the classes of interest. Furthermore, these approaches have a high computational requirement, which limits the application to mobile robotic systems for real-time operation.

Recent works in the area of intelligent vehicles [23, 24, 25] presented databases of pixel-wise semantic segmented images with object classes such as pedestrians, road, sidewalk, car, sky. Also, some realistic image proxy engines have been proposed to overcome the annotation effort to segment some object classes in indoor scenes, as with the ScanNet dataset [26] or the Stanford 2D-3D-S dataset [27]. However, concerning indoor visual navigation and assistive computer vision, the classes of interest such as doors, stairs or other path anomalies are not present or are not segmented in these datasets [9, 7, 28]. Unfortunately, the majority of available datasets for both semantic segmentation and object detection do not explicitly consider these objects [29]. In order to overcome this limitation, we acquired and labeled several images containing these custom objects of interest. In this paper, we adopted the object detection trend as a backbone for building our semantic-object augmented representation because the object information level met the expected augmented map requirements, but also because of the computational effort when making online inference with fully instance-segmentation





**Fig. 2.** Visualization of the formulation pipeline, showing the main modules and some of the information exchanged between them.

networks. Furthermore, the required user annotation effort for pixel-wise labeling is also higher when compared to box object annotation.

## 2.2 SLAM and Augmented Semantic Representations

The combination of semantic information to support mapping and localization has been also explored by several recent works. For instance, Nascimento et al. [30] applied a binary RGB-D descriptor to feed an Adaboost learning method to classify objects in a navigation task. McCormac et al. [11] proposed a method for semantic 3D mapping. Their work combined the formulation of Whelan et al. [31], an RGB-D based SLAM system for building a dense point cloud of the scene, with an encoder-decoder convolutional network for pixel-wise semantic segmentation. The segmented labels are then projected/registered into the 3D reconstructed point cloud. Similarly, Li and Belaroussi [10] provided a 3D semantic mapping system from monocular images. Their methodology is based on LSD-SLAM [32], which estimates a semi-dense 3D reconstruction of the scene and performs camera localization from monocular images. Similarly to McCormac et al. [11], the metric map and the semantic labeling are combined in order to obtain the semantic 3D map.

As discussed previously, due to runtime computational requirements and the amount of user effort to pixel-wise segmentation of the classes of interest for supervised semantic segmentation, we propose an instance semantic segmentation that leverages a lightweight data-driven object detection network with a model-based segmentation (object geometric shape priors). We also perform instance association and tracking through different time frames in order to build more complete and accurate extended maps.

## 3 METHODOLOGY

We divide our approach into four main components. An overview of the complete formulation is shown in Figure 2. The first component addresses the semantic categorization and location of objects in the image, which employs a neural

network to detect pre-trained object classes in real-time. This information is then used in a SLAM/localization step, which tracks the camera positioning in the scene and creates a projected 2D grid-based map of the environment using the available onboard robot sensors. Subsequently, we perform an efficient model-based object instance segmentation, from the object detection combined with 3D shape modeling priors. This component processes the information of the two previous components, together with the point cloud information, to localize the observed objects in the current frame and to segment pixels by fitting a primitive shape model (e.g., a planar patch for doors). Finally, the last component tracks previously localized objects in the map over time in order to combine multiple object measurements.

### 3.1 Visual Object Categorization and Detection

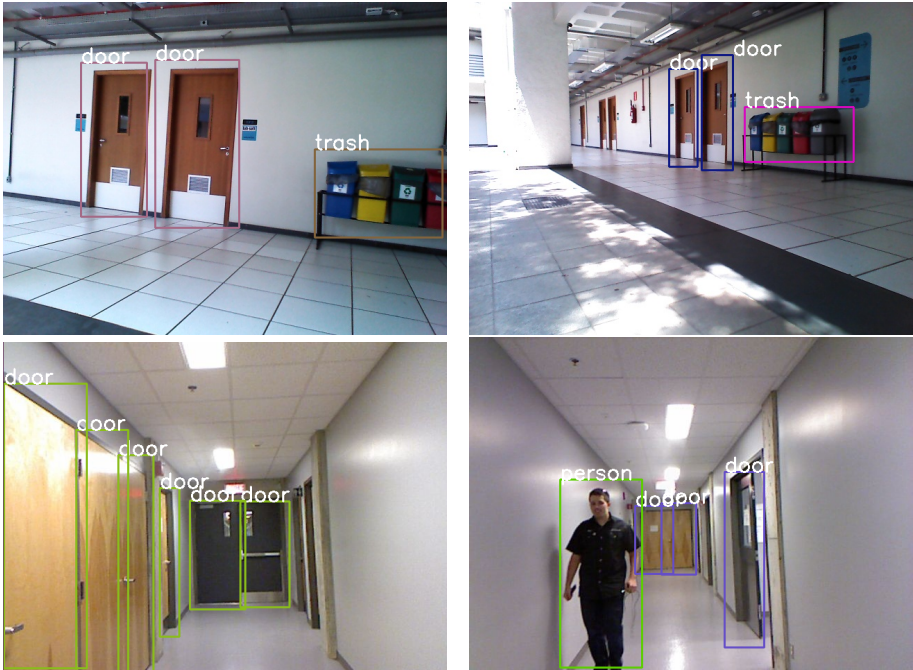
This section describes the first component of the augmented mapping framework. We start extracting a preliminary higher level representation of the scene with the detection of object classes of interest that are in the RGB-D camera field-of-view. For that, we profit of recent research progress on convolutional neural networks to reason from images to find objects and predict their semantic information, i.e., their location and category in the image. We selected the “You only look once” (Yolo) network [18] among the various available object detection techniques [15, 16, 33], because of its low computational effort and high precision-recall scores. The output of the network (as further described in the works [6, 18]) are bounding boxes modeled with four parameters: the center position coordinates  $(x, y)$ , their width  $w$  and height  $h$ .

In our context of understanding and reasoning mainly in indoor scenes, the training images of Yolo contained classes such as “door”, “bench”, “person”, “water fountain”, “trash bin” and “fire extinguisher”, as shown in Figure 3. Since the aforementioned available datasets [29] did not contain annotated images with these classes (“door” images are available on ImageNet but still with high appearance variability), we need to label and perform data augmentation in order to successfully detect these objects.

#### Training and Dataset Augmentation

We trained the network to detect a set of custom classes using a small amount of pictures taken from different objects in the environment, as well as using pre-labeled images from datasets available online (mainly for human detection). Our custom image labeling dataset consists of around 1,000 pictures of *doors*, *benches*, *trash bins*, *water fountains* and *fire extinguishers*, together with the labels of their locations and classes. The labels were manually added using a tool developed for this purpose, which we also provide with the code. A preview of the labeling process is shown in Figure 3. Our tool is structured to make easier the annotation and network training.

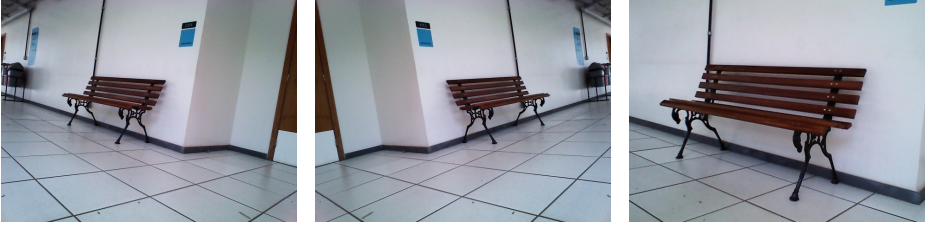
For the detection of people, we adopted the Pascal VOC 2017 and 2012 datasets [19]. This dataset comprises about 20,000 images of people, and their



**Fig. 3.** Visualization of the proposed labeling tool with different object classes.

corresponding bounding boxes. The configuration files generated for each subset (from our custom object images and the Pascal VOC dataset) were combined in order to train the network considering objects from both datasets. The network architecture was redefined to have the number of classes updated, as well as the number of filters.

One issue encountered after training and testing the network was that Pascal VOC datasets have an overwhelmingly more significant number of images than our custom built dataset of other objects. This led the network to become biased towards the *person* class, while rarely detecting other objects. To overcome this issue, we augmented the custom-built dataset using common dataset augmentation operations such as flipping, scaling & translating, and adding intensity noise. Each operation doubled the number of images. We then applied two random noise levels, two scale, and one flipping operations to the original images, which increased the number of training images by a factor of  $2^5 - 1 = 31$ . An example of the augmentation result can be seen in Figure 4. The noise operation, although increasing the network robustness, did not significantly affect the image appearance to the human eye. The final output of the system is the object boxes (encoded by five parameters) at an average mean frame-rate of 30Hz with a Nvidia GeForce GTX 1060. Some detection prediction examples can be seen in Figures 1 and 5.



**Fig. 4.** Dataset augmentation operations. The first image corresponds to the original frame, the resulting flipped image (in the center) and scaled & translated image (at right).



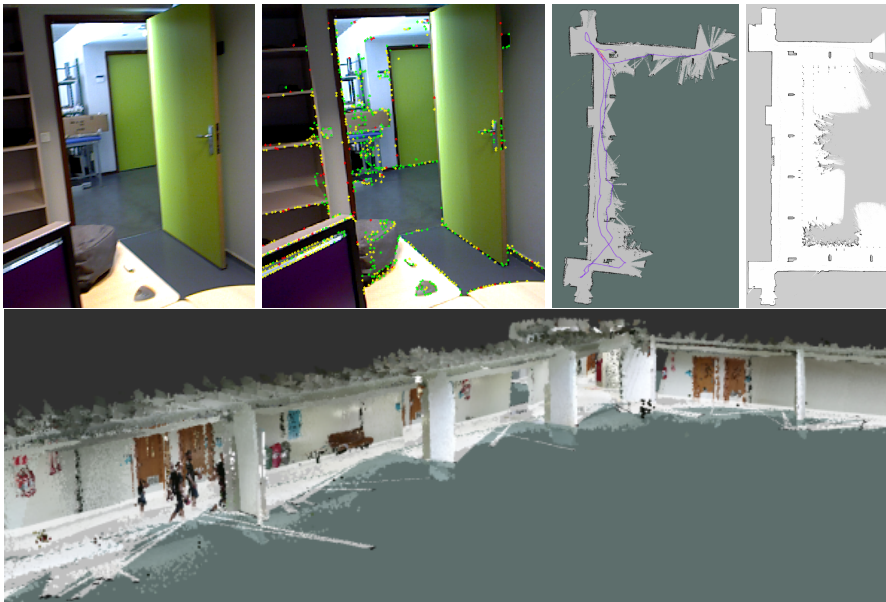
**Fig. 5.** Door detection: input image, object detection bounding box, RANSAC inliers for planar segmentation, object model represented in the map.

### 3.2 Localization and Mapping

Concurrently to the object detection described in Section 3.1, we generated an initial 2D projected map representation of the environment, along with the localization of the robot in this representation. A plethora of techniques can be used to localize the robot, depending mostly on the available sensors and computational requirements. We set as the minimal required sensor setup to our system as one RGB-D camera, which information is exploited in all stages of the formulation. However, it is worth noting that the proposed localization module is also designed to consider LIDAR and wheel odometers sensors when these are available in the robotic system. Thus, three main setups are supported:

- i) The laser scan is not available. In this case, the depth image is sampled from the RGB-D camera in order to create the scan stream.
- ii) Odometers are not available. The registration between the RGB-D frames is used in order to build the odometry information.
- iii) Both LIDAR and odometers are not available. We follow as indicated in the two previous i) and ii) settings.

In order to have an easy deployment system, we considered mostly SLAM/mapping techniques currently available on ROS. Also with the purposes of flexibility and portability, the adopted localization/mapping backbone is selected to produce



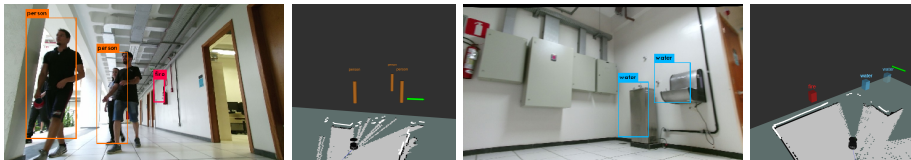
**Fig. 6.** 2D grid-based and textured point cloud of our first dataset sequence. Top row (from left to right): a sample input image, the extracted features from RTAB-Map, the resulting metric map and the map ground truth. Bottom row: the final textured 3D point cloud on top of the 2D grid-based representation.

an output consisting of a 2D grid-based representation of the environment  $\mathbf{M}$  along with the 2D projected location  $\mathbf{x}_r \in \mathbb{R}^3$  of the robot in the map:

$$\mathbf{x}_r = (x, y, \theta)^T, \quad (1)$$

where  $(x, y)$  is the position and  $\theta$  the orientation. This step can exploit any range-based or visual-based localization/SLAM algorithm, notably the provided framework supports and was tested with techniques already available on ROS as Gmapping SLAM [34], AMCL [35] and RTAB-Map library [36] which was initially developed for appearance-based loop closing and memory handling for large-scale scene mapping. These libraries provide localization and mapping techniques for several sensory modalities, including RGB-D, stereo or monocular camera settings for both 2D grid-based representation and the 3D textured point cloud of the scene.

We remark that other state-of-the-art image registration techniques such as the feature-based ORB-SLAM [37] and appearance-based RGBDSLAM [38] could also be used with minimal effort in the system, as long as the system provides camera localization and the 2D projected grip map of the scene. The required changes are then mainly in adjusting the API and ROS message exchange (subscribing and publishing topics) as done for AMCL, Gmapping and RTAB-Map algorithms. After performing several experiments, we adopted RTAB-Map



**Fig. 7.** “Person”, “water fountain” and “fire extinguisher” object detection and model fitting: input images and object model represented in the map.

for giving the most accurate and complete map results, as shown in the metric maps generated from the provided dataset sequences in Figure 6. Finally, it is worth noting that one could also leverage the redundancy given by the available sensor settings, especially concerning laser information with the depth provided by the RGB-D camera. Also, the odometry information can be either gathered from encoders, range or visual information, which of these having their complementary properties, advantages and cons to the localization and mapping.

### 3.3 Model Fitting and Positioning

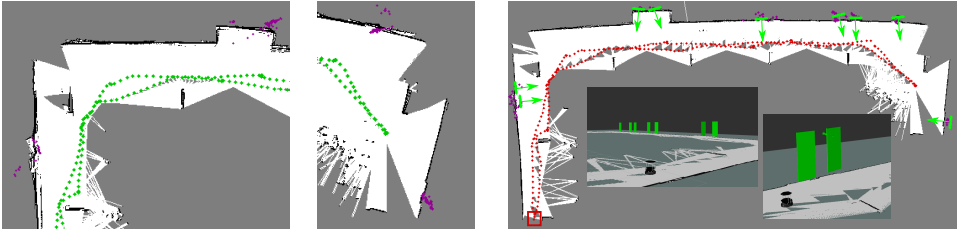
Given a set of detected objects, we perform efficient object instance segmentation of nearly thin or flat objects by adopting primitive 3D shape priors. For instance, a plane is a reasonable primitive for representing “doors”. From the RGB-D camera calibration parameters, we then reconstruct and find all 3D points inside the detected box, where the primitive model of the objects is fitted. The clustering technique adopted in the shape model fitting, to all classes except doors, was the Euclidean region growing segmentation technique [39], returning the centroid and respective convex hull dimensions. Whenever the detected objects are labeled as “door”, we fitted a planar patch using RANSAC [40] for estimating the position and orientation.

The projected pose of each object, denoted by  $\mathbf{y} \in \mathbb{R}^3$ , is then represented by the 2D projected centroid from the camera coordinate system to the global map coordinate system and its orientation. These steps can be seen in Figures 5 and 7 for different objects.

### 3.4 Object Tracking and Final Augmented Representation

After observing and projecting objects onto a location on the map, the final step is to perform tracking of the captured objects. That is, given multiple observations of multiple classes of objects (doors, benches, trash bins, etc..) across different instants of time, we wish to infer which objects have already been observed before and which have not. Ideally, we want to associate every previously seen instances with the right stored instance, and unseen objects as new instances. This would allow us to augment the map with the correct information about the semantics of the environment. Erroneous associations on this step yield undesired results as multiples instances of the same object (false





**Fig. 8.** Door object observations during the robot navigation. The extracted positions  $\mathbf{y}$  of the doors observed over time are shown in pink and the final tracked/filtered instances  $\mathbf{x}$  are indicated in green. The left image also includes two visualizations of the augmented semantic map.

positives) or associating two different observations of two different objects as belonging to the same object (false negatives).

Although tracking people is essential in re-active navigation and situation awareness, in this work, we do not perform the dynamic tracking of people because we are mostly interested to the static objects in the final augmented map representation. Furthermore, it is worth noting that an association/correspondence strategy based only on object locations is likely to fail to track humans. In this case, more elaborated models considering explicitly the appearance should be taken into account as, for instance, using bi-directional long short-term memories to handle appearance changes [41, 42].

For any given frame, all the  $m$  observed objects' positions of a given class are stored as a set of observations  $\mathbf{Y} = \{\mathbf{y}_0, \mathbf{y}_1, \dots, \mathbf{y}_m\}$ . We want to compare and check if any of these observations match one in the dictionary of  $n$  already observed instances of that same object class,  $\mathbf{X} = \{\mathbf{x}_0, \mathbf{x}_1, \dots, \mathbf{x}_n\}$ . The association cost matrix  $\mathbf{D}(\mathbf{x}_i, \mathbf{y}_j)$  between both sets is computed using the Mahalanobis distance for every possible match:

$$\mathbf{D}(\mathbf{x}_i, \mathbf{y}_j) = \sqrt{(\mathbf{y}_j - \mathbf{x}_i)^T \mathbf{S}_i^{-1} (\mathbf{y}_j - \mathbf{x}_i)}, \quad (2)$$

where  $\mathbf{x}_i$  is the  $i$ -th model of the  $n$  matched instances ( $i = \{1, 2, 3, \dots, n\}$ ) of  $\mathbf{X}$  and  $\mathbf{S}_i$  is its related covariance matrix. Once the cost matrix is computed, the association between the observations and the dictionary instances are gathered from the Hungarian-Algorithm [43]. All the resulted associations which distances are smaller than a threshold ( $\mathbf{D}(\mathbf{x}_i, \mathbf{y}_j) < \delta$ ) are assumed to correspond to previously seen objects; otherwise, new object instances representing the remaining observations are included in the dictionary.

In order to track and to increase the accuracy of detected instances, each stored semantic object is modeled with a constant state Kalman filter [6], since we are interested in storing mostly static classes in the final augmented map, to maintain its state up-to-date and combine different objects observations. Each filter combines the information of the different observations temporally as shown in Figure 8. The filter initialization and tuning details are described in Section 4.2.

The advantage of this simple tracking approach is that it pays the way for the integration of different object models that can be sufficiently described from a positioning/geometric point of view in the scene. Specifically, the positional properties of the object models of interest to this work were sufficiently discriminant to perform the tracking, as long as the accuracy of the localization/mapping system, described in Section 3.2, was below the Mahalanobis distance threshold.

## 4 EXPERIMENTS

The experiments were performed online with a mobile robot navigating indoor scenes and offline using previously acquired indoor dataset sequences. We also present qualitative results with a publicly available RGB-D dataset. We first detail the parameters setup considered in the experiments, and then we present some extended mapping results.

### 4.1 Dataset and Object Training Samples

Apart from online experiments and to evaluate the performances of the proposed method in controlled conditions, we collected a dataset containing three data sequences of different indoor places. These sequences were acquired while a robot was teleoperated in indoor environments, as depicted in Figure 9. Each sequence contains raw sensor streams recorded using the rosbag toolkit from two different RGB-D cameras, LiDAR and odometry. All data sequences contain different classes of objects: *person*, *door*, *bench*, *water fountain*, *trash bin*, *fire extinguisher*, as shown in the images of Figures 1, 3 and 9. Every class considered static (i.e., all, except for person and chair) have their location specified in a ground truth map we provide, as shown in Figure 10. An overview of these three sequences is depicted in Figures 10 and 11, which also contain the projected object positions. The RGB-D sequences used two different RGB-D sensing cameras: Microsoft Kinect (**sequence1-Kinect**) and Orbbec Astra (**sequence2-Astra** and **sequence3-Astra**). Further details of time duration, data statistics and information parsers is given in the project page<sup>3</sup>.

As previously mentioned, this dataset was built since the majority of available datasets for both semantic segmentation and 3D object detection did not consider doors and the other objects of interest to the navigation in our indoor scenes. Unfortunately, a motion capture system was not available to get the precise camera position along with the displacement in all the covered area of the scenes, nor a fine-detailed 3D mesh reconstruction of the environments due to their extension. To circumvent this limitation, we obtained the 3D robot position and of objects for each sequence performing a fine-level localization on the 2D CAD model of the scene. We then computed the 3D position of each object relative to the local image frames.

---

<sup>3</sup> <https://www.verlab.dcc.ufmg.br/semantic-mapping-for-robotics/>

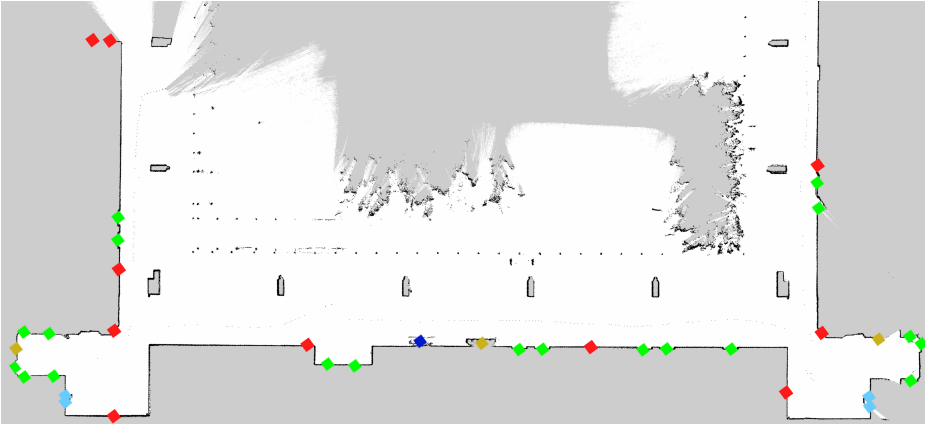




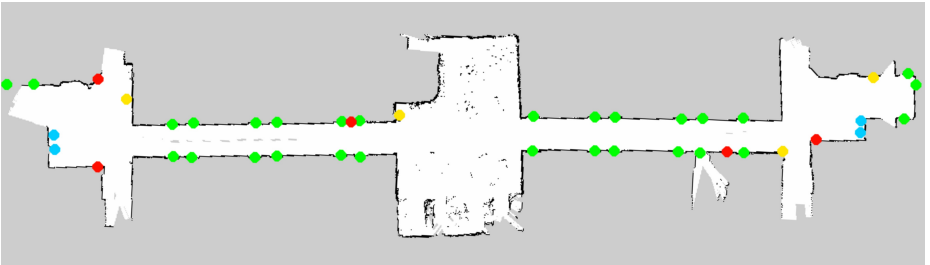
**Fig. 9.** Examples of scenes contained in the first sequence of the dataset and the Kobuki base robot. The first row displays an RGB frame and its corresponding point cloud visualization. The robot with on-board sensors (RGB-D camera and 2D LIDAR) is shown in the bottom left image.

## 4.2 System Setup and Implementation Aspects

We used a robotic platform containing a Kobuki base, where the different RGB-D cameras and LIDAR sensors were mounted, as described in the dataset Section 4.1. All the components of the formulation are integrated with ROS (Robot Operational System) and the output map generation runs at 15 Hz in a laptop with Ubuntu 16.04, Intel core i7 and Nvidia GeForce 1050 Ti. Since our main goal is to extend maps with relevant object information that do not usually change position over time, interesting candidates for navigation and user interaction available in your sequences were doors, bench, water fountain, and fire extinguisher. To this end, we trained the network following the protocol indicated in Section 3.1. In the robot localization and mapping, we adjusted few



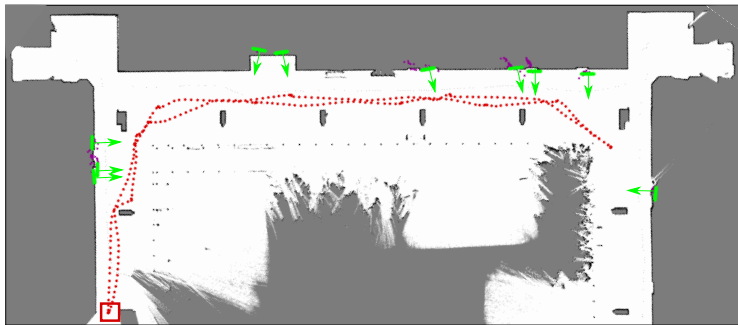
**Fig. 10.** 2D ground-truth map with the projected object positions used in the dataset sequences **sequence1-Kinect** and **sequence2-Astra**, and with a mapped area of  $42m \times 18.5m$ : doors (green squares), fire extinguisher (red squares), trash bin (in yellow, water fountain (in light blue) and bench (in dark blue).



**Fig. 11.** 2D ground-truth map (covered mapped area of  $54m \times 12m$ ) with the projected object positions used in the dataset sequence **sequence3-Astra**: doors (green squares), fire extinguisher (red squares), trash bin (in yellow, water fountain (in light blue) and bench (in dark blue).

parameters from the RTAB-Map default parameters (which are beyond 100), such as Reg/Strategy to use visual and depth information in the localization.

The geometric model fitting was performed with RANSAC [39]. We allow the point to plane fitting to optimize coefficients, and the distance threshold to 0.03. From our experiments, this value accounted for errors in the camera depth images, while allowing a correct segmentation of door points from the wall, in case these lie in different planes. In the object association and tracking, we adopted a constant uncorrelated noise affecting the process and observation measurements (i.e., the error covariance are diagonal matrices).



**Fig. 12.** Augmented 2D map with door instances with localization-only (figure axes dimensions  $55.7 \times 24.3m$ ). The red square indicates the starting and ending point of the robot trajectory (red-dotted). Purple dots are the unfiltered positions observations and green lines are the doors filtered results. The reconstructed map depicts with green arrows the position and orientation of objects.

**Table 1.** Object detection and tracking results of sequence **sequence1-Kinect**.

class	detection	FP	FN	avg. error [m]
door	19	1	3	0.78
bench	1	1	0	1.2
trash bin	3	1	0	1.04
fire exting.	9	1	3	0.53
water fount.	4	0	0	0.61

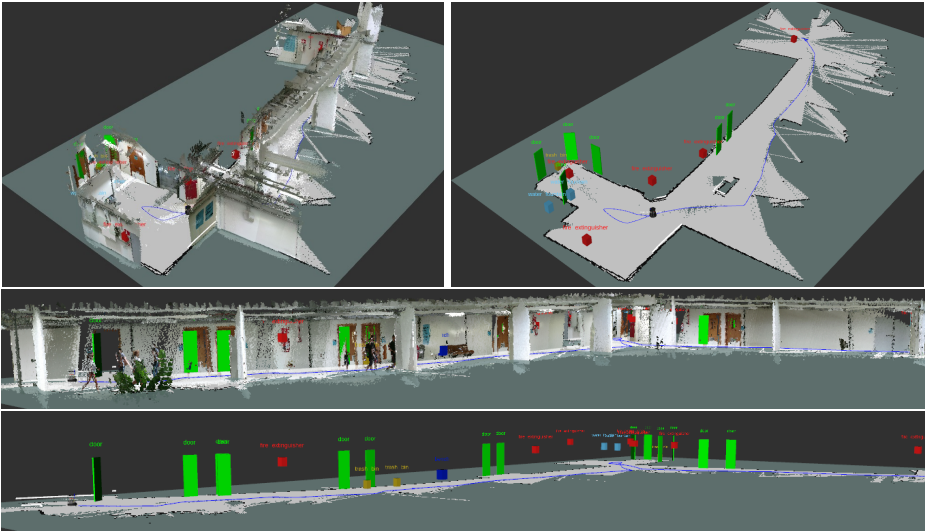
### 4.3 Augmented Mapping Results

To exemplify the flexibility of the approach, the evaluation is done with two different localization strategies: one performing 3D RGB-D SLAM (RTAB-Map) and one 2D probabilistic re-localization approach in a previously generated map (Adaptive Monte Carlo localization - AMCL[35]). We show some of the extended map results for both of these approaches in Figures 12, 13, 14 and 15. The detected objects are shown in green, red and blue representing “door”, “water fountain” and “fire extinguisher” respectively. The quantitative metrics adopted are the amount of false positives, which indicates percentage of objects that were wrongly instantiated, and the amount of false negatives that indicates the number of objects that were not integrated in the final map representation. The position errors of the objects is also considered. The adopted qualitative metric is the visual quality of the augmented semantic visualizations of the different scenes.

**Localization-based Mode** The first experiments were performed using the pure localization mode. This mode is useful for determining the accuracy of the final semantic representation since it allows the comparison of the estimated objects positions with the ground truth poses, by mitigating the undesired effects of

**Table 2.** Object detection and tracking results of sequence **sequence3-Astra**.

class	detection	FP	FN	avg. error [m]
door	18	1	12	0.67
bench	0	0	0	0
trash bin	2	0	2	0.47
fire exting.	4	0	1	7.62
water fount.	7	3	1	0.35

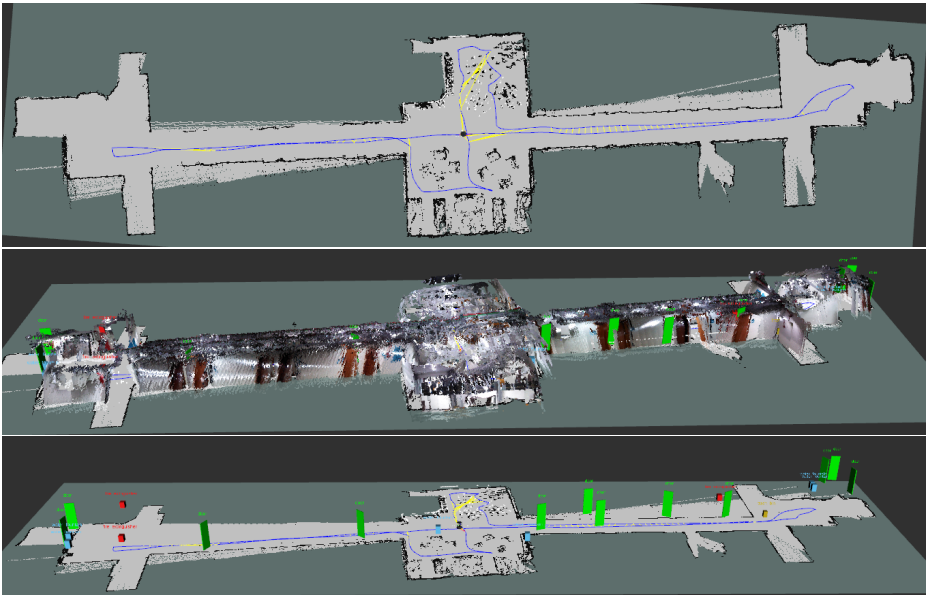
**Fig. 13.** Visualizations of the augmented map from sequence **sequence1-Kinect** with RTAB-Map. The geometric object primitives are shown in green, red and blue representing the “door”, “water fountain” and “fire extinguisher” respectively.

the errors from the mapping/SLAM system back-end. The required information in this mode is a previously acquired map and a starting robot position, such as the annotated maps shown in Figures 10 and 11. One obtained simplified map view, considering solely the door objects from sequence **sequence2-Astra**, is shown in Figure 12. We then computed the number of false positives, false negatives and projected position errors as shown in Tables 1 and 2. Note that due to hard illumination conditions, several false negatives occurred in the sequence **sequence3-Astra**, as illustrated in metrics on Table 2 and in the qualitative map visualization of Figure 14.

We then evaluate the sensibility of the main components to sensor noise, notably affecting the RGB-D camera. We also identified some key parameters that affects directly the final obtained representation. These are the association threshold and the image detection threshold. Ideally, we would desire that the framework performance to be stable from the effects of noise and with a

**Table 3.** Results varying the Mahalanobis distance threshold for the “door” class ( $\delta$ ) for sequence **sequence2-Astra**.

$\delta$ [m]	avg. error [m]	std [m]	FP	FN
0.9	0.46	0.25	27.2%	0%
1.0	0.70	0.49	18.2%	0%
1.2	0.54	0.45	11%	11%
1.5	0.87	0.63	0%	11%



**Fig. 14.** Visualizations of the augmented map from sequence **sequence3-Astra** with RTAB-Map. The geometric object primitives are shown in green, red and blue representing the “door”, “water fountain” and “fire extinguisher” respectively.

reasonable range of these parameters. The first performed parameter sensibility analysis is in the data association component, where we evaluated the influence of the Mahalanobis threshold to different distances as shown in Table 3, solely for the door objects on sequence **sequence2-Astra**. We observed that small distance association values tend to cause a smaller position error, but this also favors more false positives. This effect happens since some successive object measurements were corrupted with both positioning and model extraction errors. On the other hand, large distance association values affected close-by objects to be interpreted as the same instance.

We then realized experiments to evaluate the system robustness to different levels of noise in the RGB-D images, with errors following the properties:

- RGB:  $\tilde{\mathcal{I}}(\mathbf{p}) = \mathcal{I}(\mathbf{p}) + \mathbf{e}_I(\mathbf{p})$  and  $\mathbf{e}_I(\mathbf{p}) \sim \mathcal{N}(0, \sigma_I^2) \mathbf{I}_{1 \times 3}$ , for  $\sigma_I \in \{1, 5, 10, 20\}$ .

**Table 4.** Sensitivity experiments for different Gaussian noise levels for “door” and “fire extinguisher” objects using sequence **sequence1-Kinect**: **(Left)** Number of false positives (FP) and negatives (FN) of the final semantic representation by increasing noise in the RGB-D images. **(Right)** Sampled noise trial example for the highest variance level. Due to strong appearance changes, a “fire extinguisher” object, appearing in the left region of the image, was not detected over all frames and thus not included in the final representation.

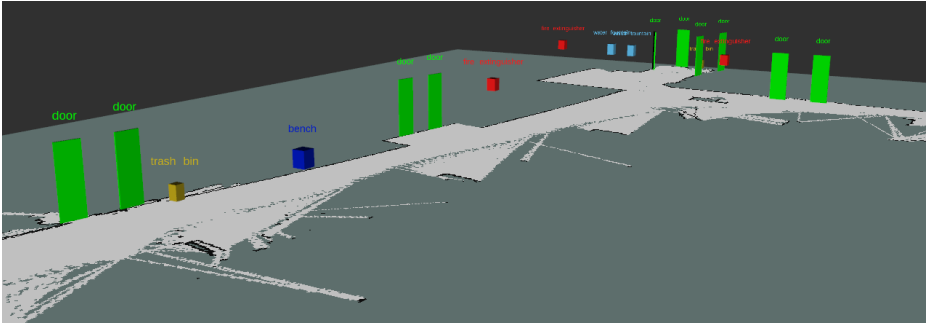
$\sigma_I$	door (FP, FN)	fire extinguisher (FP, FN)
1	(0,0)	(3,2)
5	(1,0)	(3,2)
10	(0,2)	(2,1)
20	(1,1)	(0,4)



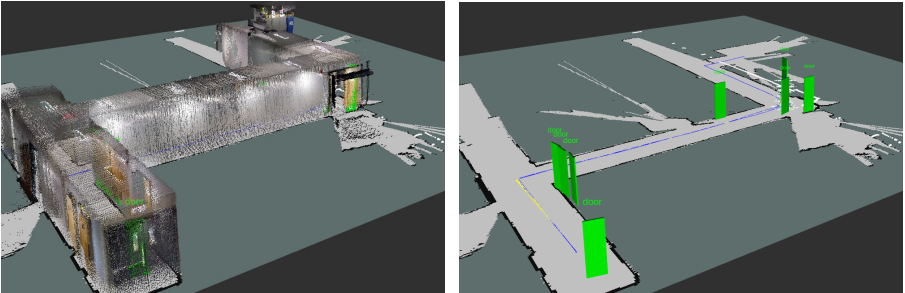
– Depth:  $\tilde{\mathcal{D}}(\mathbf{p}) = \mathcal{D}(\mathbf{p}) + e_D(\mathbf{p})$  and  $e_D(\mathbf{p}) \sim \mathcal{N}(0, \sigma_D^2)$ , for  $\sigma_D \in 0.1\sigma_I$ .

The effects of the corrupted data in the detection, tracking and positioning components were analyzed taking “door” and “fire extinguisher” classes in the **sequence1-Kinect** data sequence. This choice is due the larger number of these objects that could be observed in the scene (19 and 10 respectively). These results are presented in Table 4. We noted that the system was affected mainly for the larger errors with variance  $\sigma_I^2 = 400$ , with the increase of false positives, notably for “fire extinguisher” class. This indicated that the detection, object filtering and tracking were capable of handling these source of errors, but the performance was depreciated for the higher noise level, as shown in the image shown on right of Table 4.

**SLAM-based Mode** In the SLAM mode, the extended semantic and metric maps are built concurrently, while the robot explores the scene. As described in Section 3.2, our formulation is adapted to use the output of some commonly employed SLAM algorithms such as Gmapping and RTAB-Map. Some obtained semantic map representations are shown in Figures 13, 14 and 15. One noticed drawback of using the SLAM mode is that the formulation needs to explicitly handle the loop closing and bundle adjustment in the map generation of large spaces. While this is often done for reducing drift in trajectory errors, the tracking components should be aware of past pose adjustments to avoid misplaced objects. This undesired effect happens notably with Gmapping, which does not provide a public API of the pose graph nodes of the robot trajectory. However, this effect was greatly reduced when using RTAB-Map since we could recover the pose graph nodes directly, as illustrated in Figures 13, 14 and 15.



**Fig. 15.** Augmented 2D map of objects using **sequence2-Astra** with RTAB-Map. The geometric object primitives are shown in green, red and blue representing the “door”, “water fountain” and “fire extinguisher” respectively.



**Fig. 16.** Augmented map results from the sequence made available in RTAB-Map. Due to appearance gap only few door objects were successfully modeled into the semantic representation.

Finally, we also considered the publicly available data sequence from RTAB-Map<sup>4</sup> as shown in Figure 16. We note, however, that the available objects appearance were significantly altered from the trained ones and, thus, only a few door instances were observed and retained.

**Discussion and Limitations** We reduced the influence of the threshold association, presented in Table 3, by taking into account all observed instances simultaneously with the Hungarian algorithm assignment. Still, some scenarios were affected by this parameter, as when the objects were observed while being revisited after the robot had traveled long distances in and out of the object’s surrounding area in the SLAM mode. In these cases, the localization component was not capable of correcting the trajectory and mapping drifts.

Another parameter affecting the system performance was the latency in the object detection step, which is mainly linked with the image processing step and

<sup>4</sup> demo.mapping.bag sequence provided at [http://wiki.ros.org/rtabmap\\_ros](http://wiki.ros.org/rtabmap_ros).



ROS internal inter-process communication delay. We noticed in these cases that the objects frequently were projected into bad map locations, notably when the robot performed fast rotations. We reduced the effects of this practical limitation by storing the robot pose states at the moment when the network image processing started.

Describing each object as a two dimensional point on the map, the localization and tracking steps are greatly simplified. We might note, however, that this approach disregards object's dimensions, yielding a higher average error for larger objects, providing only a rough estimate of their position. Also, for dynamic classes (e.g., humans), tracking becomes harder, requiring more robust filtering approaches, faster processing speeds, and likely taking the object's appearance into consideration [44] [45].

We also found that object localization component was improved by setting a threshold on the maximum distance of projected objects from the robot. Conversely, objects that are visible from far away (greater than six to eight meters) were not taken into account, although this also incurred some loss in the object detection scores. Finally, let us conclude with an overview of the qualitative experimental results. Although the previously discussed limitations, the generated augmented semantic maps indicate desired characteristics for robot navigation and interaction tasks, as shown in the representations of the different dataset sequences in Figures 1, 13, 14, 15 and 16. In the case of the SLAM mode, the object localization error was affected by the robot position error itself. When bundle adjustments of the map were performed or the robot position was corrected, previously localized instances had to be corrected as well. This requires an additional object association step for every correction, which is sometimes hard to be successfully done as shown in the augmented maps shown in Figure 14.

## 5 CONCLUSIONS

This paper proposes a complete methodology and framework for building augmented maps with object-level information. This mapping framework is flexible and can be used with different sensor configurations, where the minimal required sensor setup consists of an RGB-D camera or a stereo camera rig. The first part of formulation leverages object detection with a shape segmentation strategy to perform instance semantic segmentation. This showed suitable for real-time operation in mobile robotic systems with limited computational resources, being an alternative to recent instance segmentation frameworks [21, 22]. The gathered information of the objects is improved overtime with a Kalman filtering tracking strategy, where the instances' associations are done using the Hungarian algorithm. The system was built on top of ROS, and it is highly modular, i.e., it can be easily modified without the need of changing other independent modules. The evaluation of the formulation was done in different indoor data sequences acquired in real conditions, containing people and objects as doors and other commonly found public space furnitures. This extended map representation can



be used then with motion planning algorithms and to provide situation awareness for navigation tasks. We also provide the code and a dataset composed of three data sequences, with annotated object classes (doors, fire extinguishers, benches, water fountains) and their positioning.

A possible extension to the presented work is to consider simultaneously both color and depth information in the object instance segmentation and localization. Ideally, both the object's detection, shape and pose would be performed simultaneously, in the sense of recent formulations discussed in the works about 3D shape and pose learning from images [46, 47]. Note however that our application scenarios require efficient algorithms, ideally displaying real-time performance in resource limited platforms. Another exciting direction would be to consider the semantic map in the localization while the robot navigates, as well as adopting a motion planning policy using the knowledge of the observed objects in reactive or proactive manners, seamless to how humans navigate and operates in daily-life conditions.

## Acknowledgments

The authors thank PNPd-CAPES and FAPEMIG for financial support during this research. We also gratefully acknowledge NVIDIA for the donation of the Jetson TX2 GPU used in the online experiments of this research.

## References

1. Bozhinoski, D., Di Ruscio, D., Malavolta, I., Pelliccione, P., Crnkovic, I.: Safety for mobile robotic systems: A systematic mapping study from a software engineering perspective. *Journal of Systems and Software* **151** (2019)
2. Rehder, E., Wirth, F., Lauer, M., Stiller, C.: Pedestrian prediction by planning using deep neural networks. *CoRR* (2018)
3. Carneiro, R., Nascimento, R., Guidolini, R., Cardoso, V., Oliveira-Santos, T., Badue, C., Souza, A.D.: Mapping road lanes using laser remission and deep neural networks. *CoRR* (2018)
4. Pronobis, A., Jensfelt, P.: Large-scale semantic mapping and reasoning with heterogeneous modalities. In: *IEEE ICRA*. (2012)
5. Papadakis, P., Rives, P.: Binding human spatial interactions with mapping for enhanced mobility in dynamic environments. *Autonomous Robots* **41**(5) (2017)
6. Bersan, D., Martins, R., Campos, M., Nascimento, E.R.: Semantic map augmentation for robot navigation: A learning approach based on visual and depth data. In: *IEEE LARS*. (2018)
7. Pérez-Yus, A., López-Nicolás, G., Guerrero, J.: Detection and modelling of staircases using a wearable depth sensor. In: *ECCV*. (2014)
8. Leo, M., Medioni, G., Trivedi, M., Kanade, T., Farinella, G.M.: Computer vision for assistive technologies. *Computer Vision and Image Understanding* **154** (2017)
9. Wang, H., Sun, Y., Liu, M.: Self-supervised drivable area and road anomaly segmentation using rgb-d data for robotic wheelchairs. *IEEE Robotics and Automation Letters* **4**(4) (2019)

10. Li, X., Belaroussi, R.: Semi-dense 3d semantic mapping from monocular slam. arXiv preprint arXiv:1611.04144 (2016)
11. McCormac, J., Handa, A., Davison, A., Leutenegger, S.: Semanticfusion: Dense 3d semantic mapping with convolutional neural networks. In: IEEE ICRA. (2017)
12. Häne, C., Zach, C., Cohen, A., Pollefeys, M.: Dense semantic 3D reconstruction. *IEEE Trans. on Pattern Analysis and Machine Intelligence* **39**(9) (2016)
13. Wang, C., Hou, S., Wen, C., Gong, Z., Li, Q., Sun, X., Li, J.: Semantic line framework-based indoor building modeling using backpacked laser scanning point cloud. *ISPRS journal of photogrammetry and remote sensing* **143** (2018)
14. Lateef, F., Ruichek, Y.: Survey on semantic segmentation using deep learning techniques. *Neurocomputing* (2019)
15. Krizhevsky, A., Sutskever, I., Hinton, G.: Imagenet classification with deep convolutional neural networks. In: NIPS. (2012)
16. Ren, S., He, K., Girshick, R., Sun, J.: Faster r-cnn: Towards real-time object detection with region proposal networks. In: NIPS. (2015)
17. Liu, W., Anguelov, D., Erhan, D., Szegedy, C., Reed, S., Fu, C., Berg, A.: Ssd: Single shot multibox detector. In: IEEE ECCV. (2016)
18. Redmon, J., Divvala, S., Girshick, R., Farhadi, A.: You only look once: Unified, real-time object detection. In: IEEE CVPR. (2016)
19. Everingham, M., Van Gool, L., Williams, C.K.I., Winn, J., Zisserman, A.: The pascal visual object classes (voc) challenge. *International Journal of Computer Vision* **88**(2) (2010)
20. Russakovsky, O., Deng, J., Su, H., Krause, J., Satheesh, S., Ma, S., Huang, Z., Karpathy, A., Khosla, A., Bernstein, M., Berg, A.C., Fei-Fei, L.: ImageNet Large Scale Visual Recognition Challenge. *International Journal of Computer Vision* **115**(3) (2015)
21. He, K., Gkioxari, G., Dollár, P., Girshick, R.: Mask r-cnn. In: IEEE ICCV. (2017)
22. Bolya, D., Zhou, C., Xiao, F., Lee, Y.J.: Yolact: Real-time instance segmentation. In: ICCV. (2019)
23. Gaidon, A., Wang, Q., Cabon, Y., Vig, E.: Virtual worlds as proxy for multi-object tracking analysis. In: IEEE CVPR. (2016)
24. Cordts, M., Omran, M., Ramos, S., Rehfeld, T., Enzweiler, M., Benenson, R., Franke, U., Roth, S., Schiele, B.: The cityscapes dataset for semantic urban scene understanding. In: IEEE CVPR. (2016)
25. Zhan, X., Liu, Z., Luo, P., Tang, X., Loy, C.C.: Mix-and-match tuning for self-supervised semantic segmentation. In: AAAI Conference on Artificial Intelligence. (2018)
26. Dai, A., Chang, A.X., Savva, M., Halber, M., Funkhouser, T., Nießner, M.: Scannet: Richly-annotated 3d reconstructions of indoor scenes. In: IEEE CVPR. (2017)
27. Armeni, I., Sax, A., Zamir, A.R., Savarese, S.: Joint 2D-3D-Semantic Data for Indoor Scene Understanding. ArXiv e-prints (2017)
28. Salaris, P., Vassallo, C., Souères, P., Laumond, J.P.: The geometry of confocal curves for passing through a door. *IEEE Trans. on Robotics* **31**(5) (2015)
29. Firman, M.: RGBD Datasets: Past, Present and Future. In: CVPR Workshop on Large Scale 3D Data: Acquisition, Modelling and Analysis. (2016)
30. Nascimento, E.R., Oliveira, G., Campos, M., Vieira, A.: Improving Object Detection and Recognition for Semantic Mapping with an Extended Intensity and Shape based Descriptor. In: IEEE IROS Workshop on Active Semantic Perception. (2011)
31. Whelan, T., Leutenegger, S., Salas-Moreno, R.F., Glocker, B., Davison, A.: Elasticfusion: Dense SLAM without A pose graph. In: Robotics: Science and Systems. (2015)

32. Engel, J., Schöps, T., Cremers, D.: LSD-SLAM: large-scale direct monocular SLAM. In: IEEE ECCV. (2014)
33. He, K., Zhang, X., Ren, S., Sun, J.: Deep residual learning for image recognition. In: IEEE CVPR. (2016)
34. Grisetti, G., Stachniss, C., Burgard, W.: Improved techniques for grid mapping with rao-blackwellized particle filters. *IEEE Trans. on Robotics* **23**(1) (2007)
35. Fox, D., Burgard, W., Dellaert, F., Thrun, S.: Monte carlo localization: Efficient position estimation for mobile robots. *AAAI/IAAI* **1999** (1999)
36. Labbé, M., Michaud, F.: Rtab-map as an open-source lidar and visual simultaneous localization and mapping library for large-scale and long-term online operation. *Journal of Field Robotics* **36**(2) (2019)
37. Mur-Artal, R., Montiel, J.M., Tardos, J.D.: Orb-slam: a versatile and accurate monocular slam system. *IEEE Trans. on Robotics* **31**(5) (2015)
38. Endres, F., Hess, J., Sturm, J., Cremers, D., Burgard, W.: 3-d mapping with an rgb-d camera. *IEEE Trans. on Robotics* **30**(1) (2014)
39. Rusu, R., Cousins, S.: 3D is here: Point Cloud Library (PCL). In: IEEE ICRA. (2011)
40. Raguram, R., Frahm, J.M., Pollefeys, M.: A comparative analysis of ransac techniques leading to adaptive real-time random sample consensus. In: European Conference on Computer Vision. (2008)
41. Sadeghian, A., Alahi, A., Savarese, S.: Tracking the untrackable: Learning to track multiple cues with long-term dependencies. In: IEEE ICCV. (2017)
42. Kim, C., Li, F., Rehg, J.M.: Multi-object tracking with neural gating using bilinear lstm. In: IEEE ECCV. (2018)
43. Kuhn, H.W.: The hungarian method for the assignment problem. *Naval research logistics quarterly* **2**(1-2) (1955)
44. Kim, S.J., Nam, J.Y., Ko, B.C.: Online tracker optimization for multi-pedestrian tracking using a moving vehicle camera. *IEEE Access* **6** (2018) 48675–48687
45. Dimitrievski, M., Veelaert, P., Philips, W.: Behavioral pedestrian tracking using a camera and lidar sensors on a moving vehicle. *Sensors* **19**(2) (2019) 391
46. J., K., Smith, E., Lafleche, J.F., Fuji Tsang, C., Rozantsev, A., Chen, W., Xiang, T., Lebededian, R., Fidler, S.: Kaolin: A pytorch library for accelerating 3d deep learning research. *arXiv:1911.05063* (2019)
47. Kanazawa, A., Black, M.J., Jacobs, D.W., Malik, J.: End-to-end recovery of human shape and pose. In: Computer Vision and Pattern Recognition (CVPR). (2018)